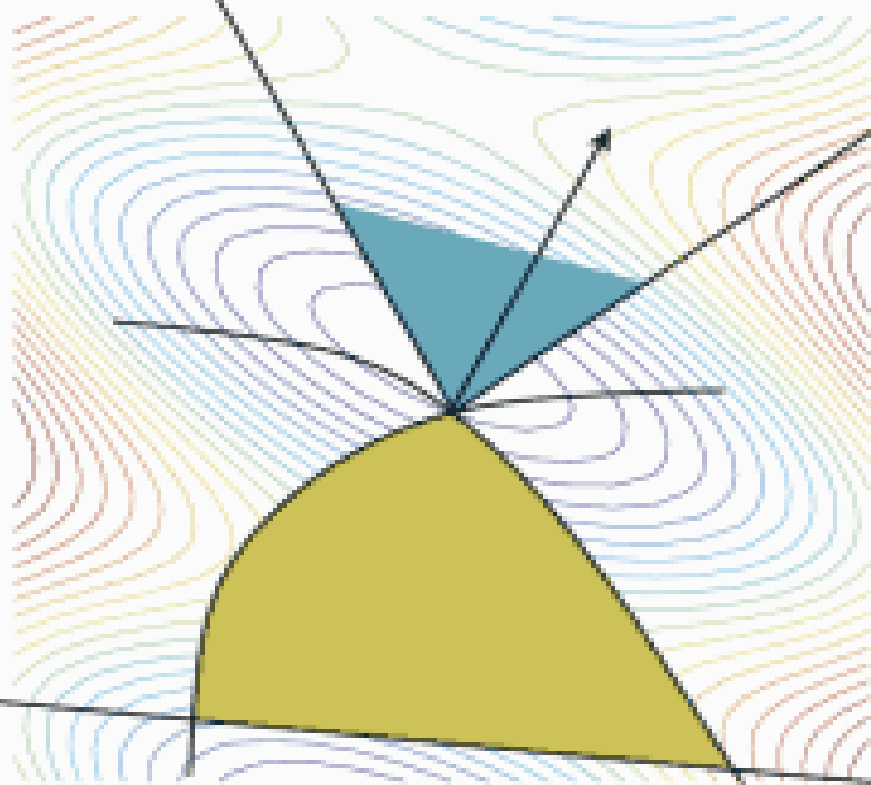


NICLAS ANDRÉASSON
ANTON EVGRAFOV
MICHAEL PATRIKSSON

An Introduction to Continuous Optimization



An Introduction to Continuous
Optimization: Foundations and
Fundamental Algorithms

Niclas Andréasson, Anton Evgrafov, and Michael Patriksson

Preface

The present book has been developed from course notes written by the third author, and continuously updated and used in optimization courses during the past several years at Chalmers University of Technology, Göteborg (Gothenburg), Sweden.

A note to the instructor: The book serves to provide lecture and exercise material in a first course on optimization for second to fourth year students at the university. The book's focus lies on providing a basis for the analysis of optimization models and of candidate optimal solutions, especially for continuous (even differentiable) optimization models. The main part of the mathematical material therefore concerns the analysis and algebra that underlie the workings of convexity and duality, and necessary/sufficient local/global optimality conditions for unconstrained and constrained optimization problems. Natural algorithms are then developed from these principles, and their most important convergence characteristics analyzed. The book answers many more questions of the form “Why/why not?” than “How?”

This choice of focus is in contrast to books mainly providing numerical guidelines as to how these optimization problems should be solved. The number of algorithms for linear and nonlinear optimization problems—the two main topics covered in this book—are kept quite low; those that are discussed are considered classical, and serve to illustrate the basic principles for solving such classes of optimization problems and their links to the fundamental theory of optimality. Any course based on this book therefore should add project work on concrete optimization problems, including their modelling, analysis, solution by practical algorithms, and interpretation.

A note to the student: The material assumes some familiarity with linear algebra, real analysis, and logic. In linear algebra, we assume an active knowledge of bases, norms, and matrix algebra and calculus. In real analysis, we assume an active knowledge of sequences, the basic

Preface

topology of sets, real- and vector-valued functions and their calculus of differentiation. We also assume a familiarity with basic predicate logic, since the understanding of proofs require it. A summary of the most important background topics is found in Chapter 2, which also serves as an introduction to the mathematical notation. The student is advised to refresh any unfamiliar or forgotten material of this chapter before reading the rest of the book.

We use only elementary mathematics in the main development of the book; sections of supplementary material that provide an outlook into more advanced topics and that require more advanced methods of presentation are kept short, typically lack proofs, and are also marked with an asterisk.

A detailed road map of the contents of the book's chapters are provided at the end of Chapter 1. Each chapter ends with a selected number of exercises which either illustrate the theory and algorithms with numerical examples or develop the theory slightly further. In Appendix A solutions are given to most of them, in a few cases in detail. (Those exercises marked "exam" together with a date are examples of exam questions that have been given in the course "Applied optimization" at Göteborg University and Chalmers University of Technology since 1997.)

In our work on this book we have benefited from discussions with Dr. Ann-Brith Strömberg, presently at the Fraunhofer-Chalmers Research Centre for Industrial Mathematics (FCC), Göteborg, and formerly at mathematics at Chalmers University of Technology, as well as Dr. Fredrik Altenstedt, also formerly at mathematics at Chalmers University of Technology, and currently at Carmen Systems AB. We thank the heads of undergraduate studies at mathematics, Göteborg University and Chalmers University of Technology, Jan-Erik Andersson and Sven Järner, respectively, for reducing our teaching duties while preparing this book. We also thank Yumi Karlsson for helping us by typesetting a main part of the first draft based on the handwritten notes; after the fact, we now realize that having been helped with this first draft made us confident that such a tremendous task as that of writing a text book would actually be possible. Finally, we thank all the students who gave us critical remarks on the first versions during 2004 and 2005.

Göteborg, May 2005

Niclas Andréasson

Anton Evgrafov

Michael Patriksson

Contents

I	Introduction	1
1	Modelling and classification	3
1.1	Modelling of optimization problems	3
1.2	A quick glance at optimization history	9
1.3	Classification of optimization models	11
1.4	Conventions	13
1.5	Applications and modelling examples	15
1.6	Defining the field	16
1.7	On optimality conditions	16
1.8	Soft and hard constraints	18
1.8.1	Definitions	18
1.8.2	A derivation of the exterior penalty function . . .	19
1.9	A road map through the material	20
1.10	On the background of this book and a didactics statement	25
1.11	Illustrating the theory	26
1.12	Notes and further reading	27
1.13	Exercises	28
II	Fundamentals	31
2	Analysis and algebra—A summary	33
2.1	Reductio ad absurdum	33
2.2	Linear algebra	34
2.3	Analysis	37
3	Convex analysis	41
3.1	Convexity of sets	41
3.2	Polyhedral theory	42
3.2.1	Convex hulls	42

3.2.2	Polytopes	45
3.2.3	Polyhedra	47
3.2.4	The Separation Theorem and Farkas' Lemma . . .	52
3.3	Convex functions	57
3.4	Application: the projection of a vector onto a convex set .	66
3.5	Notes and further reading	69
3.6	Exercises	69
 III Optimality Conditions		73
4	An introduction to optimality conditions	75
4.1	Local and global optimality	75
4.2	Existence of optimal solutions	78
4.2.1	A classic result	78
4.2.2	*Non-standard results	81
4.2.3	Special optimal solution sets	83
4.3	Optimality in unconstrained optimization	84
4.4	Optimality for optimization over convex sets	88
4.5	Near-optimality in convex optimization	95
4.6	Applications	96
4.6.1	Continuity of convex functions	96
4.6.2	The Separation Theorem	98
4.6.3	Euclidean projection	99
4.6.4	Fixed point theorems	100
4.7	Notes and further reading	106
4.8	Exercises	107
5	Optimality conditions	111
5.1	Relations between optimality conditions and CQs at a glance	111
5.2	A note of caution	112
5.3	Geometric optimality conditions	114
5.4	The Fritz John conditions	118
5.5	The Karush–Kuhn–Tucker conditions	124
5.6	Proper treatment of equality constraints	128
5.7	Constraint qualifications	130
5.7.1	Mangasarian–Fromovitz CQ (MFCQ)	131
5.7.2	Slater CQ	131
5.7.3	Linear independence CQ (LICQ)	132
5.7.4	Affine constraints	132
5.8	Sufficiency of the KKT conditions under convexity	133
5.9	Applications and examples	135
5.10	Notes and further reading	137

5.11 Exercises	138
6 Lagrangian duality	141
6.1 The relaxation theorem	141
6.2 Lagrangian duality	142
6.2.1 Lagrangian relaxation and the dual problem	142
6.2.2 Global optimality conditions	147
6.2.3 Strong duality for convex programs	149
6.2.4 Strong duality for linear and quadratic programs .	154
6.2.5 Two illustrative examples	156
6.3 Differentiability properties of the dual function	158
6.3.1 Subdifferentiability of convex functions	158
6.3.2 Differentiability of the Lagrangian dual function .	162
6.4 *Subgradient optimization methods	164
6.4.1 Convex problems	164
6.4.2 Application to the Lagrangian dual problem	170
6.4.3 The generation of ascent directions	173
6.5 *Obtaining a primal solution	174
6.5.1 Differentiability at the optimal solution	175
6.5.2 Everett's Theorem	176
6.6 *Sensitivity analysis	177
6.6.1 Analysis for convex problems	177
6.6.2 Analysis for differentiable problems	179
6.7 Applications	181
6.7.1 Electrical networks	181
6.7.2 A Lagrangian relaxation of the traveling salesman problem	185
6.8 Notes and further reading	189
6.9 Exercises	190
 IV Linear Programming	 195
7 Linear programming: An introduction	197
7.1 The manufacturing problem	197
7.2 A linear programming model	198
7.3 Graphical solution	199
7.4 Sensitivity analysis	199
7.4.1 An increase in the number of large pieces available	200
7.4.2 An increase in the number of small pieces available	201
7.4.3 A decrease in the price of the tables	202
7.5 The dual of the manufacturing problem	203
7.5.1 A competitor	203

7.5.2	A dual problem	203
7.5.3	Interpretations of the dual optimal solution	204
8	Linear programming models	205
8.1	Linear programming modelling	205
8.2	The geometry of linear programming	210
8.2.1	Standard form	211
8.2.2	Basic feasible solutions and the Representation Theorem	214
8.2.3	Adjacent extreme points	220
8.3	Notes and further reading	223
8.4	Exercises	223
9	The simplex method	225
9.1	The algorithm	225
9.1.1	A BFS is known	226
9.1.2	A BFS is not known: phase I & II	232
9.1.3	Alternative optimal solutions	236
9.2	Termination	237
9.3	Computational complexity	238
9.4	Notes and further reading	238
9.5	Exercises	239
10	LP duality and sensitivity analysis	241
10.1	Introduction	241
10.2	The linear programming dual	242
10.2.1	Canonical form	243
10.2.2	Constructing the dual	243
10.3	Linear programming duality theory	247
10.3.1	Weak and strong duality	247
10.3.2	Complementary slackness	250
10.4	The dual simplex method	254
10.5	Sensitivity analysis	257
10.5.1	Perturbations in the objective function	258
10.5.2	Perturbations in the right-hand side coefficients . .	259
10.6	Notes and further reading	260
10.7	Exercises	261
V	Algorithms	265
11	Unconstrained optimization	267
11.1	Introduction	267

11.2	Descent directions	269
11.2.1	Introduction	269
11.2.2	Newton's method and extensions	271
11.3	The line search problem	275
11.3.1	A characterization of the line search problem	275
11.3.2	Approximate line search strategies	276
11.4	Convergent algorithms	279
11.5	Finite termination criteria	281
11.6	A comment on non-differentiability	283
11.7	Trust region methods	284
11.8	Conjugate gradient methods	285
11.8.1	Conjugate directions	286
11.8.2	Conjugate direction methods	287
11.8.3	Generating conjugate directions	288
11.8.4	Conjugate gradient methods	289
11.8.5	Extension to non-quadratic problems	292
11.9	A quasi-Newton method: DFP	293
11.10	Convergence rates	296
11.11	Implicit functions	296
11.12	Notes and further reading	297
11.13	Exercises	298
12	Optimization over convex sets	303
12.1	Feasible direction methods	303
12.2	The Frank–Wolfe algorithm	305
12.3	The simplicial decomposition algorithm	308
12.4	The gradient projection algorithm	311
12.5	Application: traffic equilibrium	317
12.5.1	Model analysis	317
12.5.2	Algorithms and a numerical example	319
12.6	Notes and further reading	321
12.7	Exercises	322
13	Constrained optimization	325
13.1	Penalty methods	325
13.1.1	Exterior penalty methods	326
13.1.2	Interior penalty methods	330
13.1.3	Computational considerations	333
13.1.4	Applications and examples	334
13.2	Sequential quadratic programming	337
13.2.1	Introduction	337
13.2.2	A penalty-function based SQP algorithm	340
13.2.3	A numerical example on the MSQP algorithm	345

Contents

13.2.4 On recent developments in SQP algorithms	346
13.3 A summary and comparison	346
13.4 Notes and further reading	347
13.5 Exercises	348
VI Appendix	351
A Answers to the exercises	353
Chapter 1: Modelling and classification	353
Chapter 3: Convexity	356
Chapter 4: An introduction to optimality conditions	358
Chapter 5: Optimality conditions	360
Chapter 6: Lagrangian duality	361
Chapter 8: Linear programming models	363
Chapter 9: The simplex method	365
Chapter 10: LP duality and sensitivity analysis	366
Chapter 11: Unconstrained optimization	368
Chapter 12: Optimization over convex sets	370
Chapter 13: Constrained optimization	371
References	373
Index	385

Part I

Introduction

Modelling and classification



1.1 Modelling of optimization problems

The word “optimum” is Latin, and means “the ultimate ideal;” similarly, “optimus” means “the best.” Therefore, to *optimize* refers to trying to bring whatever we are dealing with towards its ultimate state, that is, towards its optimum. Let us take a closer look at what that means in terms of an example, and at the same time bring the definition of the term *optimization* forward, as the scientific field understands and uses it.

Example 1.1 (a staff planning problem) Consider a hospital ward which operates 24 hours a day. At different times of day, the staff requirement differs. Table 1.1 shows the demand for reserve wardens during six work shifts.

Table 1.1: Staff requirements at a hospital ward.

Shift	1	2	3	4	5	6
Hours	0–4	4–8	8–12	12–16	16–20	20–24
Demand	8	10	12	10	8	6

Each member of staff works in 8 hour shifts. The goal is to fulfill the demand with the least total number of reserve wardens. ■

Consider now the following interpretation of the term “to optimize:”

To optimize = to do something as well as is possible.

Modelling and classification

We utilize this description to identify the mathematical problem associated with Example 1.1; in other words, we create a *mathematical model* of the above problem.

To do something: We identify activities which we can control and influence. Each such activity is associated with a *variable* whose value (or, activity level) is to be decided upon (that is, optimized). The remaining quantities are constants in the problem.

As well as: How good a vector of activity levels is is measured by a real-valued function of the variable values. This quantity is to be given a highest or lowest value, that is, we minimize or maximize, depending on our goal; this defines the *objective function*.

Is possible: Normally, the activity levels cannot be arbitrarily large, since an activity often is associated with the utilization of resources (time, money, raw materials, labour, etcetera) that are limited; there may also be requirements of a least activity level, resulting from a demand. Some variables must also fulfill technical/logical restrictions, and/or relationships among themselves. The former can be associated with a variable necessarily being integer-valued or non-negative, by definition. The latter is the case when products are blended, a task is performed for several types of products, or a process requires the input from more than one source. These restrictions on activities form *constraints* on the possible choices of the variable values.

Looking again at the problem described in Example 1.1, this is then our declaration of a mathematical model thereof:

Variables We define

x_j := number of reserve wardens whose first shift is j , $j = 1, 2, \dots, 6$.

Objective function We wish to minimize the total number of reserve wardens, that is, the objective function, which we call f , is to

$$\text{minimize } f(\mathbf{x}) := x_1 + x_2 + \dots + x_6 = \sum_{j=1}^6 x_j.$$

Constraints There are two types of constraints:

Demand The demand for wardens during the different shifts can be written as the following inequality constraints:

$$\begin{aligned}x_6 + x_1 &\geq 8, \\x_1 + x_2 &\geq 10, \\x_2 + x_3 &\geq 12, \\x_3 + x_4 &\geq 10, \\x_4 + x_5 &\geq 8, \\x_5 + x_6 &\geq 6.\end{aligned}$$

Logical There are two physical/logical constraints:

Sign $x_j \geq 0$, $j = 1, \dots, 6$.

Integer x_j integer, $j = 1, \dots, 6$.

Summarizing, we have defined our first mathematical optimization model, namely, that to

$$\begin{aligned}\underset{\mathbf{x}}{\text{minimize}} \quad & f(\mathbf{x}) := \sum_{j=1}^6 x_j, \\ \text{subject to} \quad & x_1 + x_6 \geq 8, \quad (\text{last shift: 1}) \\ & x_1 + x_2 \geq 10, \quad (\text{last shift: 2}) \\ & x_2 + x_3 \geq 12, \quad (\text{last shift: 3}) \\ & x_3 + x_4 \geq 10, \quad (\text{last shift: 4}) \\ & x_4 + x_5 \geq 8, \quad (\text{last shift: 5}) \\ & x_5 + x_6 \geq 6, \quad (\text{last shift: 6}) \\ & x_j \geq 0, \quad j = 1, \dots, 6, \\ & x_j \text{ integer}, \quad j = 1, \dots, 6.\end{aligned}$$

This problem has an *optimal solution*, which we denote by \mathbf{x}^* , that is, a vector of decision variable values which gives the objective function its minimal value among the *feasible solutions* (that is, the vectors \mathbf{x} that satisfy all the constraints). In fact, the problem has at least two optimal solutions: $\mathbf{x}^* = (4, 6, 6, 4, 4, 4)^T$ and $\mathbf{x}^* = (8, 2, 10, 0, 8, 0)^T$; the *optimal value* is $f(\mathbf{x}^*) = 28$. (The reader is asked to verify that they are indeed optimal.)

The above model is of course a crude simplification of any real application. In practice, we would have to add requirements on the individual's competence as well as other more detailed restrictions, the planning horizon is usually longer, employment rules and other conditions apply, etcetera, which all contribute to a more complex model. We mention a few successful applications of staffing problems below.

Example 1.2 (applications of staffing optimization problems) (a) It was reported in 1990 that following a staffing problem application for the Montreal municipality bus company, employing 3,000 bus drivers and 1,000 metro drivers and ticket salespersons and guards, the municipality saved some 4 million Canadian dollars per year.

(b) Together with the San Francisco police department a group of operations research scientists developed in 1989 a planning tool based on a heuristic solution of the staff planning and police vehicle allocation problem. It has been reported that it gave a 20% faster planning and savings in the order of 11 million US dollars per year.

(c) In an 1986 application, scientists collaborating with United Airlines considered their crew scheduling problem. This is a complex problem, where the time horizon is long (typically, 30 minute intervals during 7 days), and the constraints that define a feasible pattern of allocating staff to airplanes are defined by, among others, complicated work regulations. The savings reported then was 6 million US dollars per year. Carmen Systems AB in Gothenburg develop and market such tools; customers include American Airlines, Lufthansa, SAS, and SJ; Carmen Systems has one of the largest concentrations of optimizers in Sweden. ■

Remark 1.3 (on the complexity of the variable definition) The variables x_j defined in Example 1.1 are *decision variables*; we say that, since the selection of the values of these variables are immediately connected to the decisions to be made in the decision problem, and they also contain, within their very definition, a substantial amount of information about the problem at hand (such as shifts being eight hours long).

In the application examples discussed in Example 1.2 the variable definitions are much more complex than in our simple example. A typical decision variable arising in a crew scheduling problem is associated with a specific staff member, his/her home base, information about the crew team he/she works with, a current position in time and space, a flight leg specified by flight number(s), additional information about the staff member's previous work schedule and work contract, and so on. The number of possible combinations of work schedules for a given staff member is nowadays so huge that not all variables in a crew scheduling problem can even be defined! (That is, the complete problem we wish to solve cannot be written down.) The philosophy in solving a crew scheduling problem is instead to algorithmically *generate* variables that one believes may receive a non-zero optimal value, and most of the computational effort lies in defining and solving good variable generation problems, whose result is (part of) a feasible work schedule for given staff members. The term *column generation* is the operations researcher's term for this process of generating variables. ■

Remark 1.4 (non-decision variables) Not all variables in a mathematical optimization model are decision variables:

In linear programming, we will utilize *slack variables* whose role is to take on the difference between the left-hand and the right-hand side of an inequality constraint; the slack variable thereby aids in the transformation of the inequality constraint to an equality constraint, which is more appropriate to work with in linear programming.

Other variables can be introduced into a mathematical model simply in order to make the model more easy to state or interpret, or to improve upon the properties of the model. As an example of the latter, consider the following simple problem: we wish to minimize over \mathbb{R} the special one-variable function $f(x) := \text{maximum}\{x^2, x + 2\}$. (Plot the function to see where the optimum is.) This is an example of a non-differentiable function: at $x = 2$, for example, both the functions $f_1(x) := x^2$ and $f_2(x) := x + 2$ define the value of the function f , but they have different derivatives there. One way to turn this problem into a differentiable one is by introducing an additional variable. We let z take on the value of the largest of $f_1(x)$ and $f_2(x)$ for a given value of x , and instead write the problem as that to minimize z , subject to $z \in \mathbb{R}$, $x \in \mathbb{R}$, and the additional constraints that $x^2 \leq z$ and $x + 2 \leq z$. Convince yourself that this transformation is equivalent to the original problem in terms of the set of optimal solutions in x , and that the transformed problem is differentiable. ■

Figure 1.1 illustrates several issues in the modelling process, which are forthwith discussed.

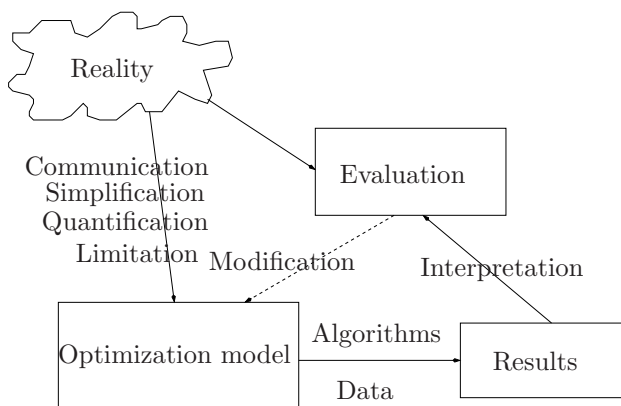


Figure 1.1: Flow chart of the modelling process

The decision problem faced in the “fluffy” reality is turned into an optimization model, through a process with several stages. By communicating with those who have raised the issue of solving the problem in the first place, one reaches an understanding about the problem to be solved. In order to identify and describe the components of a mathematical model which is also tractable, it is often necessary to simplify and also limit the problem somewhat, and to quantify any remaining qualitative statements.

The modelling process does not come without difficulties. The communication can often be difficult, simply because the two parties speak different languages in terms of describing the problem. The optimization problem quite often has uncertainties in the data, which moreover are not always easy to collect or to quantify. Perhaps the uncertainties are there for a purpose (such as in financial decision problems), but it may be that data is uncertain because not enough effort has been put into providing a good enough accuracy. Further, there is often a conflict between problem solvability and problem realism.

The problem actually solved through the use of an optimization methodology must be supplied with data, providing model constants and parameters in functions describing the objective function and perhaps also some of the constraints. For this optimization problem, an optimization algorithm then yields a result in the form of an optimal value and/or optimal solution, if an optimal solution exists. This result is then interpreted and evaluated, which may lead to alterations of the model, and certainly to questions regarding the applicability of the optimal solution. The optimization model can also be altered slightly in order to answer sensitivity analysis (“what if?”) type questions concerning the effect of small variations in data.

The final problems that we will mention come at this stage: it is crucial that the interpretation of the result makes sense to those who wants to use the solution, and, finally, it must be possible to transfer the solution back into the “fluffy” world where the problem came from.

The art of forming *good* optimization models is as much an art as a science, and an optimization course can only really cover the latter. On the other hand, this part of the modelling process should not be glossed over; it is often possible to construct more than one form of a mathematical model that represents the same problem equally accurately, and the computational complexity can differ substantially between them. Forming a good model is in fact as crucial to the success of the application as the modelling exercise itself.

Optimization problems can be grouped together in classes, according to their properties. According to this classification, the staffing problem

is a *linear integer optimization problem*. In Section 1.3 we present some major distinguishing factors between different problem classes.

1.2 A quick glance at optimization history

At Chalmers, courses in optimization are mainly given at the mathematics department. “Mainly” is the important word here, because courses that have a substantial content of optimization theory and/or methodology can be found also at other departments, such as computer science, the mechanical, industrial and chemical engineering departments, and at the Gothenburg School of Economics. The reason is that optimization is so broad in its applications.

From the mathematical standpoint, optimization, or *mathematical programming* as it is sometimes called, rests on several legs: analysis, topology, algebra, discrete mathematics, etcetera, build the foundation of the theory, and applied mathematics subjects such as numerical analysis and mathematical parts of computer science build the bridge to the algorithmic side of the subject. On the other side, then, with optimization we solve problems in a huge variety of areas, in the technical, natural, life and engineering sciences, and in economics.

Before moving on, we would just like to point out that the term “program” has nothing to do with “computer program;” a program is understood to be a “decision program,” that is, a strategy or decision rule. A “mathematical program” therefore is a mathematical problem designed to produce a decision program.

The history of optimization is very long. Many, very often geometrical or mechanical, problems (and quite often related to warfare!) that Archimedes, Euclid, Heron, and other masters from antiquity formulated and also solved, are optimization problems. For example, we mention the problem of maximizing the volume of a closed three-dimensional object (such as a sphere or a cylinder) built from a two-dimensional sheet of metal with a given area.

The masters of two millenia later, like Bernoulli, Lagrange, Euler, and Weierstrass developed variational calculus, studying problems in applied physics (and still often with a mind towards warfare!) such as how to find the best trajectory for a flying object.

The notion of *optimality* and especially how to *characterize* an optimal solution, began to be developed at the same time. Characterizations of various forms of optimal solutions are indeed a crucial part of any basic optimization course. (See Section 1.7.)

The scientific subject *operations research* refers to the study of decision problems regarding operations, in the sense of controlling complex

systems and phenomena. The term was coined in the 1940s at the height of World War 2 (WW2), when the US and British military commands hired scientists from several disciplines in order to try to solve complex problems regarding the best way to construct convoys in order to avoid, or protect the cargo ships from, enemy (read: German) submarines, how to best cover the British isles with radar equipment given the scarce availability of radar systems, and so on. The multi-disciplinarity of these questions, and the common topic of maximizing or minimizing some objective function subject to constraints, can be seen as being the defining moment of the scientific field. A better term than operations research is *decision science*, which better reflects the scope of the problems that can be, and are, attacked using optimization methods.

Among the scientists that took part in the WW2 effort in the US and Great Britain, some were the great pioneers in placing optimization on the map after WW2. Among them, we find several researchers in mathematics, physics, and economics, who contributed greatly to the foundations of the field as we now know it. We mention just a few here. George W. Dantzig invented the *simplex method* for solving linear optimization problems during his WW2 efforts at Pentagon, as well as the whole machinery of modelling such problems.¹ Dantzig was originally a statistician and famously, as a young Ph.D. student, provided solutions to some then unsolved problems in mathematical statistics that he found on the blackboard when he arrived late to a lecture, believing they were (indeed hard!) home work assignments in the course. Building on the knowledge of duality in the theory of two-person zero-sum games, which had been developed by the world-famous mathematician John von Neumann in the 1920s, Dantzig was very much involved in developing the theory of duality in linear programming, together with the various characterizations of an optimal solution that is brought out from that theory. A large part of the duality theory was developed in collaboration with the mathematician Albert W. Tucker.

Several researchers interested in national economics studied transportation models at the same time, modelling them as special linear optimization problems. Two of them, the mathematician Leonid W. Kantorovich and the statistician Tjalling C. Koopmans received The Bank of Sweden Prize in Economic Sciences in Memory of Alfred Nobel in 1975 “for their contributions to the theory of optimum allocation of resources.” They had, in fact, both worked out some of the basics

¹As Dantzig explains in [Dan57], linear programming formulations in fact can first be found in the work of the first theoretical economists in France, such as F. Quesnay in 1760; they explained the relationships between the landlord, the peasant and the artisan. The first practical linear programming problem solved with the simplex method was the famous Diet problem.

of linear programming, independently of Dantzig, at roughly the same time. (Dantzig stands out among the three especially for creating an efficient algorithm for solving such problems, but also as being the most important developer of the theory of linear programming.)²

1.3 Classification of optimization models

We here develop a subset of problem classes that can be set up by contrasting certain aspects of a general optimization problem. We let

$$\begin{aligned} \mathbf{x} &\in \mathbb{R}^n : \text{vector of decision variables } x_j, \quad j = 1, 2, \dots, n; \\ f : \mathbb{R}^n &\rightarrow \mathbb{R} \cup \{\pm\infty\} : \text{objective function}; \\ X &\subseteq \mathbb{R}^n : \text{ground set defined logically/physically}; \\ g_i : \mathbb{R}^n &\rightarrow \mathbb{R} : \text{constraint function defining restriction on } \mathbf{x} : \\ g_i(\mathbf{x}) &\geq b_i, \quad i \in \mathcal{I}; \quad (\text{inequality constraints}) \\ g_i(\mathbf{x}) &= d_i, \quad i \in \mathcal{E}. \quad (\text{equality constraints}) \end{aligned}$$

We let $b_i \in \mathbb{R}$, $i \in \mathcal{I}$, and $d_i \in \mathbb{R}$, $i \in \mathcal{E}$, denote the right-hand sides of these constraints; without loss of generality, we could actually let them all be equal to zero, as any constants can be incorporated into the definitions of the functions g_i , $i \in \mathcal{I} \cup \mathcal{E}$.

The optimization problem then is to

$$\begin{aligned} \text{minimize} \quad & f(\mathbf{x}), & (1.1a) \\ \text{subject to} \quad & g_i(\mathbf{x}) \geq b_i, \quad i \in \mathcal{I}, & (1.1b) \\ & g_i(\mathbf{x}) = d_i, \quad i \in \mathcal{E}, & (1.1c) \\ & \mathbf{x} \in X. & (1.1d) \end{aligned}$$

(If it is really a maximization problem, then we change the sign of f .)

The problem type depends on the nature of the functions f and g_i , and the set X . Let us look at some examples.

(LP) Linear programming Objective function linear: $f(\mathbf{x}) := \mathbf{c}^T \mathbf{x} = \sum_{j=1}^n c_j x_j$, $\mathbf{c} \in \mathbb{R}^n$; constraint functions affine: $g_i(\mathbf{x}) := \mathbf{a}_i^T \mathbf{x} - b_i$, $\mathbf{a}_i \in \mathbb{R}^n$, $b_i \in \mathbb{R}$, $i \in \mathcal{I} \cup \mathcal{E}$; $X := \{\mathbf{x} \in \mathbb{R}^n \mid x_j \geq 0, \quad j = 1, 2, \dots, n\}$.

(NLP) Nonlinear programming Some functions f, g_i , $i \in \mathcal{I} \cup \mathcal{E}$, are nonlinear.

²Incidentally, several other laureates in economics have worked with the tools of optimization: Paul A. Samuelson (1970, linear programming), Kenneth J. Arrow (1972, game theory), Wassily Leontief (1973, linear transportation models), Gerard Debreu (1983, game theory), Harry M. Markowitz (1990, quadratic programming in finance), John F. Nash, Jr. (1994, game theory), William Vickrey (1996, econometrics), and Daniel L. McFadden (2000, microeconomics).

Continuous optimization $f, g_i, i \in \mathcal{I} \cup \mathcal{E}$, are continuous on an open set containing X ; X is closed and convex.

(IP) Integer programming $X \subseteq \{0, 1\}^n$ (binary) or $X \subseteq \mathbb{Z}^n$ (integer).

Unconstrained optimization $\mathcal{I} \cup \mathcal{E} := \emptyset$; $X := \mathbb{R}^n$.

Constrained optimization $\mathcal{I} \cup \mathcal{E} \neq \emptyset$ and/or $X \subset \mathbb{R}^n$.

Differentiable optimization $f, g_i, i \in \mathcal{I} \cup \mathcal{E}$, are continuously differentiable on X ; further, X is closed and convex.

Non-differentiable optimization At least one of $f, g_i, i \in \mathcal{I} \cup \mathcal{E}$, is non-differentiable.

(CP) Convex programming f is convex; $g_i, i \in \mathcal{I}$, are concave; $g_i, i \in \mathcal{E}$, are affine; and X is closed and convex. (See Section 3.3 for definitions.)

Non-convex programming The complement of the above.

In Figure 1.3 we show how the problems NLP, IP, and LP are related.

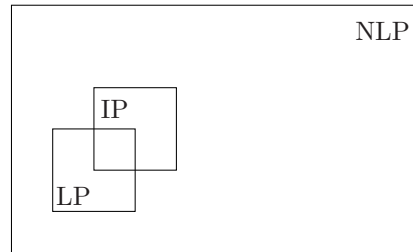


Figure 1.2: The relations among NLP, IP, and LP.

That LP is a special case of NLP is clear by the fact that a linear function is a special kind of nonlinear function; that IP is a special case of NLP can be illustrated by the fact that the constraint $x_j \in \{0, 1\}$ can be written as the nonlinear constraint $x_j(1 - x_j) = 0$.³

Last, there is a subclass of IP that is equivalent to LP, that is, a class of problems for which there exists at least one optimal solution which

³If a non-negative integer variable x_j is upper bounded by the integer M , it is also possible to write $\prod_{k=0}^M (x_j - k) = (x_j - 0)(x_j - 1) \cdots (x_j - M) = 0$, by which we restrict a *continuous* variable x_j to be integer-valued.

automatically is integer valued even without imposing any integrality constraints, provided of course that the problem has any optimal solutions at all. We say that such problems have the *integrality property*. An important example problem belonging to this category is the linear single-commodity network flow problem with integer data; this class of problems in turn includes as special cases such important problems as the linear versions of the assignment problem, the transportation problem, the maximum flow problem, and the shortest route problem.

Among the above list of problem classes, we distinguish, roughly only, between two of the most important ones, as follows:

- LP Linear programming \approx applied linear algebra. LP is “easy,” because there exist algorithms that can solve every LP problem instance efficiently in practice.
- NLP Nonlinear programming \approx applied analysis in several variables. NLP is “hard,” because there does *not* exist an algorithm that can solve every NLP problem instance efficiently in practice. NLP is such a large problem area that it contains very hard problems as well as very easy problems. The largest class of NLP problems that are solvable with some algorithm in reasonable time is CP (of which LP is a special case).

Our problem formulation (1.1) does not cover the following:

- infinite-dimensional problems (that is, problems formulated in function spaces rather than vector spaces);
- implicit functions f and/or g_i , $i \in \mathcal{I} \cup \mathcal{E}$: then, no explicit formula can be written down; this is typical in engineering applications, where the value of, say, $f(\mathbf{x})$ can be the result of a simulation (see Section 11.11 for more details);
- multiple-objective optimization:

$$\text{“minimize } \{f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_p(\mathbf{x})\}”;$$

- optimization under uncertainty, or, stochastic programming (that is, where some of f , g_i , $i \in \mathcal{I} \cup \mathcal{E}$, are only known probabilistically).

1.4 Conventions

Let us denote the set of vectors satisfying the constraints (1.1b)–(1.1d) by $S \subseteq \mathbb{R}^n$, that is, the set of *feasible solutions* to the problem (1.1). What exactly do we mean by solving the problem to

$$\underset{\mathbf{x} \in S}{\text{minimize}} \ f(\mathbf{x})? \tag{1.2}$$

Modelling and classification

Since there is no explicit operation involved here, the question is warranted. The following two operations are however well-defined:

$$f^* := \infimum_{x \in S} f(x)$$

denotes the infimum value of the function f over the set S ; if and only if the infimum value is attained at some point \mathbf{x}^* in S (and then both f^* and \mathbf{x}^* necessarily are finite) we can write that

$$f^* := \text{minimum}_{x \in S} f(x), \quad (1.3)$$

and then we of course have that $f(\mathbf{x}^*) = f^*$. (When considering maximization problems, we obtain the analogous definitions of the supremum and the maximum.)

The second operation defines the set of optimal solutions to the problem at hand:

$$S^* := \arg \text{minimum}_{x \in S} f(x);$$

the set $S^* \subseteq S$ is nonempty if and only if the infimum value f^* is attained. Finding at least one optimal solution,

$$\mathbf{x}^* \in \arg \text{minimum}_{x \in S} f(x), \quad (1.4)$$

is a special case which moreover defines an often much more simple task.

Consider the problem instance where $S = \{x \in \mathbb{R} \mid x \geq 0\}$ and

$$f(x) := \begin{cases} 1/x, & \text{if } x > 0, \\ +\infty, & \text{otherwise;} \end{cases}$$

here, $f^* = 0$ but $S^* = \emptyset$ —the value 0 is not attained for a finite value of x , so the problem has a finite infimum value but not an optimal solution.

These examples lead to our convention in reading the problem (1.2): the statement “solve the problem (1.2)” means “find f^* and an $\mathbf{x}^* \in S^*$, or conclude that $S^* = \emptyset$.”

Hence, it is *implicit* in the formulation that we are interested both in the infimum value and in (at least) one optimal solution if one exists. Whenever we are certain that only one of them is of interest we will state so explicitly. We are aware that the interpretation of (1.2) may be considered “vague” since no operation is visible; so, to summarize and clarify our convention, it in fact includes two operations, (1.3) and (1.4).

There is a second reason for stating the optimization problem (1.1) in the way it is, a reason which is computational. To solve the problem, we almost always need to solve a sequence of relaxations/simplifications of

the original problem in order to eventually reach a solution. (These manipulations include Lagrangian relaxation, penalization, and objective function linearization, to be developed later on.) When describing the particular relaxation/simplification utilized, having access to constraint identifiers [such as (1.1c)] certainly makes the presentation easier and clearer. That will become especially valuable when dealing with various forms of duality, when (subsets of) the constraints are relaxed.

A last comment on conventions: as it is stated prior to the problem formulation (1.1) the objective function f can in general take on both $\pm\infty$ as values. Since we are generally going to study minimization problems, we will only be interested in objective functions f having the properties that (a) $f(\mathbf{x}) > -\infty$ for every feasible vector \mathbf{x} , and (b) $f(\mathbf{x}) < +\infty$ for at least one feasible vector \mathbf{x} . Such functions are known as *proper* functions (which makes sense, as it is impossible to perform a proper optimization unless these two properties hold). We will some times refer to these properties, in particular by stating explicitly when f can take on the value $+\infty$, but we will assume throughout that f does *not* take on the value $-\infty$. So, in effect then, *we assume implicitly that the objective function f is proper*.

1.5 Applications and modelling examples

To give but a quick view of the scope of applications of optimization, here is a subset of the past few years of projects performed at Chalmers University of Technology or Gothenburg University:

- Planning schedules for snow removal machines, disabled persons transportation, and school transports
- Optimization of personnel planning for airlines
- Allocation of fault tolerances in complex assembly
- Scheduling production and distribution of electricity
- Scheduling paper cutting in paper mills
- Optimization of engine performance for aircraft, boats, and cars
- Engineering design by derivative-free optimization
- Maintenance optimization for aircraft jet engines
- Portfolio optimization under uncertainty for pension funds
- Policy optimization in financial planning
- Analysis of investment in future energy systems
- Optimal wave-length and routing in optical networks
- Intensity-modulated radiation therapy (IMRT)
- Optimal congestion pricing in urban traffic networks

1.6 Defining the field

To define what the subject area of optimization encompasses is difficult, given that it is connected to so many scientific areas in the natural and technical sciences.

An obvious distinguishing factor is that an optimization model always has an objective function and a group of constraints. On the other hand by letting $f \equiv 0$ and $\mathcal{I} = \emptyset$ then the generic problem (1.1) is that of a *feasibility problem* for equality constraints, covering the important topic of solving systems of linear equations, and by instead letting $\mathcal{I} \cup \mathcal{E} = \emptyset$ we obtain an *unconstrained optimization* problem. Both these special cases are classic problems in *numerical analysis*, which most often deal with the solution of a linear or non-linear system of equations.

We can here identify a distinguishing element between optimization and numerical analysis—that an optimization problem often involve *inequality constraints* while a problem in numerical analysis does not. Why does that make a difference? The reason is that while in the latter case the analysis is performed on a manifold—possibly even a linear subspace—the analysis of an optimization problem must deal with feasible regions residing in different dimensions because of the nature of inequality constraints being either active or inactive. As a result, there will always be some kind of *non-differentiability* present in some associated functionals, while numerical analysis typically is “smooth.”

As an illustration (albeit beyond the scope of this book), we ask the reader what the extension of the famous *Implicit Function Theorem* is when we replace the system $\mathbf{F}(\mathbf{u}, \mathbf{x}) = \mathbf{0}^k$ with, say, $\mathbf{F}(\mathbf{u}, \mathbf{x}) \leq \mathbf{0}^k$?

1.7 On optimality conditions

The most important topic of the book is the analysis of the local or global optimality of a given feasible vector \mathbf{x}^* in the problem (1.2), and its links to the construction of algorithms for finding such vectors. While locally or globally optimal vectors are the ones preferred, the types of vectors that one can expect to reach for a general problem are referred to as *stationary points*; we define what we mean by $\mathbf{x}^* \in S$ being a stationary point in the problem (1.2) in non-mathematical terms as follows:

$\mathbf{x}^* \in S$ is a *stationary point* in the problem (1.2) if, with the use only of first-order⁴ information about the problem at \mathbf{x}^* , we cannot find a feasible descent direction at \mathbf{x}^* .

⁴This means that we only utilize the values of f and ∇f at \mathbf{x}^* , and the same for any constraint functions defining the set S .

In mathematical terms, this condition can be written as follows:

$\mathbf{x}^* \in S$ is a *stationary point* in the problem (1.2) if $-\nabla f(\mathbf{x}^*) \in N_S(\mathbf{x}^*)$ holds, where $N_S(\mathbf{x}^*)$ is the *normal cone* to S at \mathbf{x}^* .

See Definition 4.25 for the definition of the normal cone.

In applications to all model problems considered in the book this condition collapses to something that is rather easy to check.⁵ In the most general case, however, its use in formulating necessary optimality conditions requires further that the point \mathbf{x}^* satisfies a regularity condition referred to as a *constraint qualification* (CQ).

The connection between local or global optimality, stationarity and regularity is given by the following two implications, which constitute the perhaps two most important ones in the entire book:

$$\left. \begin{array}{l} \mathbf{x}^* \text{ local min in (1.2)} \\ \mathbf{x}^* \text{ regular} \end{array} \right\} \implies \mathbf{x}^* \text{ stationary point in (1.2);} \quad (1.5)$$

$$\left. \begin{array}{l} \mathbf{x}^* \text{ stationary point in (1.2)} \\ \text{the problem (1.2) is convex} \end{array} \right\} \implies \mathbf{x}^* \text{ global min in (1.2).} \quad (1.6)$$

The logical implication $A \implies B$ is equivalent to $\neg B \implies \neg A$ and $\neg(A \wedge B)$ is equivalent to $(\neg A) \vee (\neg B)$.⁶ Hence, the implication (1.5) means that if \mathbf{x}^* is not stationary then it is not a local minimum in (1.2) or it is not a regular point. Since the latter case is rare, the typical case then is that a non-stationary point is not locally optimal, and in the process of investigating whether \mathbf{x}^* is stationary we quite often are able to generate a feasible descent direction if there is one. Investigating stationarity is therefore important for two reasons: if we are at a stationary point \mathbf{x}^* , then \mathbf{x}^* is an interesting candidate for an optimal solution (when the problem is convex then \mathbf{x}^* is even guaranteed to be a global minimum); if \mathbf{x}^* is not a stationary point, then we can generate a feasible descent direction from \mathbf{x}^* in order to move on to a better feasible solution, thus constructing an iterative sequence of improved solutions.

For the construction and analysis of optimality conditions we refer the reader to Section 4.3 for the unconstrained case ($S := \mathbb{R}^n$), to Section 4.4 for the case where the feasible set S is assumed to be convex, and to Chapter 5 for the most general case when S need not be convex and regularity issues become important.

⁵For example, if $S := \mathbb{R}^n$ then $N_S(\mathbf{x})$ is everywhere equal to $\mathbf{0}^n$; then the stationarity condition simply states that $\nabla f(\mathbf{x}^*) = \mathbf{0}^n$.

⁶Here, “ \neg ” means “not,” “ \wedge ” means “and,” and “ \vee ” means “or.”

1.8 Soft and hard constraints

1.8.1 Definitions

We have not yet discussed the role of different types of constraints. In the *set covering problem*, the constraints are of the form $\sum_{j=1}^n a_{ij}x_j \geq 1$, $i = 1, 2, \dots, m$, where $a_{ij} \in \{0, 1\}$. These, as well as constraints of the form $x_j \geq 0$ and $x_j \in \{0, 1\}$ are *hard constraints*, meaning that if they are violated then the solution does not make much sense. Typically, such constraints are technological ones; for example, if x_j is associated with the level of production, then a negative value has no meaning, and therefore a negative value is never acceptable. A binary variable, $x_j \in \{0, 1\}$, is often logical, associated with the choice between something being “on” or “off,” such as a production facility, a city being visited by a traveling salesman, and so on; again, a fractional value like 0.7 makes no sense, and binary restrictions almost always are “hard.”

Consider now a collection of constraints that are associated with the capacity of production, and which have the form $\sum_{j=1}^n u_{ij}x_{ij} \leq c_i$, $i = 1, 2, \dots, m$, where x_{ij} denotes the level of production of an item/product j using the production process i , u_{ij} is a positive number associated with the use of a resource (man hours, hours before inspection of the machine, etcetera) per unit of production of the item, and c_i is the available capacity of this resource in the production process. In some circumstances, it is natural to allow for the left-hand side to become larger than the capacity, because that production plan might still be feasible, provided however that additional resources are made available. We consider two types of ways to allow for this violation, and which give rise to two different types of solution.

The first, which we are not quite ready to discuss here from a technical standpoint, is connected to the *Lagrangian relaxation* of the capacity constraints. If, when solving the corresponding Lagrangian dual optimization problem, we terminate the solution process prematurely, we will typically have a terminal primal vector that violates some of the capacity constraints slightly. Since the capacity constraints are soft, this solution may be acceptable.⁷ See Chapter 6 for further details on Lagrangian duality.

Since it is natural that additional resources come at an additional cost, an increase in the violation of this *soft constraint* should have the effect of an increased cost in the objective function. In other words,

⁷One interesting application arises when making capacity expansion decisions in production and work force planning problems (e.g., Johnson and Montgomery [JoM74, Example 4-14]) and in forest management scheduling (Hauer and Hoganson [HaH96]).

violating a constraint should come with a *penalty*. Given a measure of the cost of violating the constraints, that is, the unit cost of additional resource, we may transform the resulting problem to an unconstrained problem with a *penalty function* representing the original constraint.

Below, we relate soft constraints to exterior penalties.

1.8.2 A derivation of the exterior penalty function

Consider the standard nonlinear programming problem to

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}), \quad (1.7a)$$

$$\text{subject to} \quad g_i(\mathbf{x}) \geq 0, \quad i = 1, \dots, m, \quad (1.7b)$$

where f and g_i , $i = 1, \dots, m$, are real-valued functions.

Consider the following relaxation of (1.7), where $\rho > 0$:

$$\underset{(\mathbf{x}, \mathbf{s})}{\text{minimize}} \quad f(\mathbf{x}) + \rho \sum_{i=1}^m s_i, \quad (1.8a)$$

$$\begin{aligned} \text{subject to} \quad & g_i(\mathbf{x}) \geq -s_i, \quad i = 1, \dots, m, \\ & s_i \geq 0, \quad i = 1, \dots, m. \end{aligned} \quad (1.8b)$$

We interpret this problem as follows: by allowing the variable s_i to become positive, we allow for extra slack in the constraint, at a positive cost, ρs_i , proportional to the violation.

How do we solve this problem for a given value of $\rho > 0$? We specialize the following result (see, for example, [RoW97, Proposition 1.35]): for a function $\phi : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$ one has in terms of $p(\mathbf{s}) := \inf_{\mathbf{x}} \phi(\mathbf{x}, \mathbf{s})$ and $q(\mathbf{x}) := \inf_{\mathbf{s}} \phi(\mathbf{x}, \mathbf{s})$ that

$$\inf_{(\mathbf{x}, \mathbf{s})} \phi(\mathbf{x}, \mathbf{s}) = \inf_{\mathbf{x}} q(\mathbf{x}) = \inf_{\mathbf{s}} p(\mathbf{s}).$$

In other words, we can solve an optimization problem in two types of variables \mathbf{x} and \mathbf{s} by “eliminating” one of them (in our case, \mathbf{s}) through optimization, and then determine the best value of the remaining one.

Suppose then that we for a moment keep \mathbf{x} fixed to an arbitrary value. The above problem (1.8) then reduces to that to

$$\begin{aligned} \underset{\mathbf{s}}{\text{minimize}} \quad & \rho \sum_{i=1}^m s_i, \\ \text{subject to} \quad & s_i \geq -g_i(\mathbf{x}), \quad i = 1, \dots, m, \\ & s_i \geq 0, \quad i = 1, \dots, m, \end{aligned}$$

which clearly separates into the m independent problems to

$$\begin{aligned} & \underset{s_i}{\text{minimize}} \quad \rho s_i, \\ & \text{subject to} \quad s_i \geq -g_i(\mathbf{x}), \\ & \quad \quad \quad s_i \geq 0. \end{aligned}$$

This problem is trivially solvable: $s_i := \text{maximum}\{0, -g_i(\mathbf{x})\}$, that is, s_i takes on the role of a slack variable for the constraint. Using this expression in the problem (1.8) we finally obtain the problem to

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{x}) + \rho \sum_{i=1}^m \text{maximum}\{0, -g_i(\mathbf{x})\}. \quad (1.9)$$

If the constraints instead are of the form $g_i(\mathbf{x}) \leq 0$, then the resulting penalty function is of the form $\rho \sum_{i=1}^m \text{maximum}\{0, g_i(\mathbf{x})\}$.

We note that the use of the linear penalty term in (1.8a) resulted in the penalty problem (1.9); other penalty terms than (1.8a) lead to other penalty problems. See Section 13.1 for a thorough discussion on and analysis of penalty functions and methods.

1.9 A road map through the material

Chapter 2 gives a short overview of some basic material from calculus and linear algebra that is used throughout the book. Familiarity with these topics is therefore very important.

Chapter 3 is devoted to the study of convexity, a subject known as *convex analysis*. We characterize the convexity of sets and real-valued functions and show their relations. We give an overview of the special convex sets called polyhedra, which can be described by linear constraints. Parts of the theory covered, such as the Representation Theorem, Farkas' Lemma, and the Separation Theorem, build the foundation of the study of optimality conditions in Chapter 5, the theory of strong duality in Chapter 6 and of linear programming in Chapters 7–10.

Chapter 4 gives an overview of topics associated with optimality, including the result that locally optimal solutions are globally optimal in a convex problem. We establish results regarding the existence of optimal solutions, including Weierstrass' Theorem, and establish basic logical relationships between locally optimal solutions and characterizations in terms of conditions of “stationarity.” Along the way, we define important concepts such as the normal cone, the variational inequality, and the Euclidean projection of a vector onto a convex set, and outline fixed point theorems and their applications.

Chapter 5 collects results leading up to the central Karush–Kuhn–Tucker (KKT) Theorem on the necessary conditions for the local optimality of a feasible point in a constrained optimization problem. Essentially, these conditions state that a given feasible vector \mathbf{x} can only be a local minimum if there is no descent direction at \mathbf{x} which simultaneously is a feasible direction. In order to state the KKT conditions in algebraic terms such that it can be checked in practice and such that as few interesting vectors \mathbf{x} as possible satisfy them, we must restrict our study to problems and vectors satisfying some regularity properties. These properties are called constraint qualifications (CQs); among them, the classic one is that “the active constraints are linearly independent” which is familiar from the Lagrange Multiplier Theorem in differential calculus. Our treatment however is more general and covers weaker (that is, better) CQs as well. The chapter begins with a schematic road map for these results to further help in the study of this material.

Chapter 6 presents a broad picture of the theory of Lagrangian duality. Associated with the KKT conditions in the previous chapter is a vector, known as the Lagrange multiplier vector. The Lagrange multipliers are associated with an optimization problem which is referred to as the Lagrangian dual problem.⁸ The role of the dual problem is to define a largest lower bound on the optimal value f^* of the original (primal) problem. We establish the basic properties of this dual problem. In particular, it is always a convex problem, and therefore appealing to solve in order to extract the optimal solution to the primal problem. This chapter is in fact much devoted to the topic of analyzing when it is possible to generate, from an optimal dual solution, in a rather simple manner an optimal primal solution. The most important term in this context then is “strong duality” which refers to the occasion when the optimal values in the two problems are equal—only then can the “translation” be relatively easy. Some of the results established are immediately transferable to the important case of linear programming, whose duality theory is analyzed in Chapter 10. The main difference is that in the present chapter we must work with more general tools, while for linear programming we have access to a more specialized analysis; therefore, proof techniques, for example in establishing strong duality, will be quite different. Additional topics include an analysis of optimization algorithms for the solution of the Lagrangian dual problem, and applications.

⁸The dual problem was first discovered in the study of (linear) matrix games by John von Neumann in the 1920s, but had for a long time implicitly been used also for nonlinear optimization problems before it was properly stated and studied by Arrow, Hurwicz, Uzawa, Everett, Falk, Rockafellar, etcetera, starting in earnest in the 1950s. By the way, the original problem is then referred to as the primal problem, a name given by George Dantzig’s father.

Chapters 7–10 are devoted to the study of linear programming (LP) models and methods. Its importance is unquestionable: it has been stated that in the 1980s LP problems was the scientific problem that ate the most computing power in the world. While the efficiency of LP solvers have multiplied since then, so has the speed of computers, and LP models still define the most important problem area in optimization in practice. (Partly, this is also due to the fact that integer optimization solvers use LP techniques.) It is not only for this reason, however, that we devote special chapters to this topic. Their optimal solutions can be found using quite special techniques that are not common to nonlinear programming. As was shown in Chapter 4 linear programs have optimal solutions at the extreme points of the polyhedral feasible set. This fact, together with the linearity of the objective function and the constraints, means that a feasible-direction (descent) method can be cleverly devised. Since we know that only extreme points are of interest, we start at one extreme point, and then only consider as candidate search directions those that point towards another (in fact, adjacent) extreme point. We can generate such directions efficiently by using a basis representation of the extreme points; the move from one extreme point to the other is associated with a simple basis change. This procedure is known as the simplex method, which was invented by George Dantzig in the 1940s.

In Chapter 7 a manufacturing problem illustrates the basics of linear programming. The problem is solved geometrically and shown to have an optimal extreme point. We investigate how the optimal solution changes if the data of the problem is changed, and the linear programming dual to the manufacturing problem is derived by using economical arguments.

Chapter 8 begins with a presentation of the axioms underlying the use of LP models, and a general modelling technique is discussed. The rest of the chapter deals with the geometry of LP models. It is shown that every linear program can be transformed into the *standard form* which is the form that the simplex method requires. We introduce the concept of *basic feasible solution* (BFS) and discuss its connection to extreme points. A version of the Representation Theorem adapted to the standard form is presented, and we show that if there exists an optimal solution to a linear program in standard form, then there exists an optimal solution among the basic feasible solutions. Finally, we define adjacency between extreme points and give an algebraic characterization of adjacency which proves that the simplex method at each iteration step moves from one extreme point to an adjacent one.

Chapter 9 presents the simplex method. We first assume that a BFS is known at the start of the algorithm, and then describe what to do when a BFS is not known. Termination characteristics of the

algorithm are discussed: it is shown that if all the BFSs of the problem are non-degenerate, then the algorithm terminates; if, however, there exist degenerate BFSs there is a possibility that the algorithm cycles between degenerate BFSs and hence never terminates. We introduce *Bland's rule* for choosing the adjacent BFS, which eliminates cycling. We close the chapter by discussing the computational complexity of the simplex algorithm.

In Chapter 10 linear programming duality is studied. We discuss how to construct the linear programming dual to a general linear program and present duality theory. The dual simplex method is developed, and we discuss how the optimal solution of a linear program changes if the right-hand sides or the objective function coefficients are modified.

Chapter 11 presents basic algorithms for differentiable, unconstrained optimization problems. The typical optimization algorithm is iterative, which means that a solution is approached through a sequence of trial vectors, typically such that each consecutive objective value is strictly lower than the previous one in a minimization problem. This improvement is possible because we can generate improving search directions—descent (ascent) directions in a minimization (maximization) problem—by means of solving an approximation of the original problem or the optimality conditions. This approximate problem (for example, the system of Newton equations) is then combined with a line search, which approximately solves the original problem over the half-line defined by the current iterate and the search direction. This idea of combining approximation (or, relaxation) with a line search (or, coordination) is the basic methodology also for constrained optimization problems. Also, while our opinion is that the subject of differentiable unconstrained optimization largely is a subject within numerical analysis rather than within the optimization field, its understanding is important because the approximations/relaxations that we utilize in constrained optimization often result in (essentially) unconstrained optimization subproblems. We develop a class of quasi-Newton methods in detail.

Chapter 12 presents classic algorithms for differentiable nonlinear optimization over polyhedral sets, which utilize LP techniques when searching for an improving direction. The basic algorithm is known as the *Frank–Wolfe algorithm*, or the *conditional gradient method*; it utilizes $\nabla f(\mathbf{x}_k)$ as the linear cost vector at iteration k , and the direction towards any optimal extreme point \mathbf{y}_k has already in Chapter 4 been shown to be a feasible direction of descent whenever \mathbf{x}_k is not stationary. We also present an improvement in which we utilize (possibly) all the previously generated extreme points to replace the line search with a multi-dimensional one over the convex hull of these vectors. The *gradi-*

ent projection method extends the steepest descent method for unconstrained optimization problem in a natural manner. The subproblems here are Euclidean projection problems which in this case are strictly convex quadratic programming problems that can be solved efficiently for some types of polyhedral sets. The convergence results reached show that convexity of the problem is crucial in reaching good convergence results—not only regarding the global optimality of limit points but regarding the nature of the set of limit points as well: Under convexity, the gradient projection algorithm converges to an optimal solution provided that one exists, even when the set of optimal solutions is unbounded; the result immediately specializes to the steepest descent method.

Chapter 13 begins by describing natural approaches to nonlinearly constrained optimization problems, wherein all (or, a subset of) the constraints are replaced by penalties. The resulting penalized problem is then possible to solve by using techniques for unconstrained problems or problems with convex feasible sets, like those we present in Chapters 11 and 12. In order to force the penalized problems to more and more resemble the original one, the penalties are more and more strictly enforced. There are essentially two types of penalty functions: exterior and interior penalties. Exterior penalty methods were devised mainly in the 1960s, and are perhaps the most natural ones; they are valid for almost every type of explicit constraints, and are therefore amenable to solving also non-convex problems. The penalty terms are gradually enforced by letting larger and larger weights be associated with the constraints in comparison with the objective function. Under some circumstances, one can show that a finite value of these penalty parameters are needed, but in general they must tend to infinity. Interior penalty methods are also amenable to the solution of non-convex problems, but are perhaps most naturally associated with convex problems, where they are quite effective. In particular, the best methods for linear programming in terms of their worst-case complexity are interior point methods which are based on interior penalty functions. In this type of method, the interior penalties are asymptotes with respect to the constraint boundaries; a decreasing value of the penalty parameters then allow for the boundaries to be approached at the same time as the original objective function come more and more into play. For both types of methods, we reach convergence results on the convergence to KKT points in the general case—including estimates of the Lagrange multipliers—and global convergence results in the convex case.

Chapter 13 also describes another popular class of algorithms for nonlinear programming problems with (twice) differentiable objective and constraint functions. It is called Sequential Quadratic Program-

ming (SQP) and is, essentially, Newton's method applied to the KKT conditions of the problem; there are, however, some modifications necessary. For example, because of the linearization of the constraints, it is in general difficult to maintain feasibility in the process, and therefore convergence cannot merely be based on line searches in the objective function; instead one must devise a measure of "goodness" that takes constraint violation into account. The classic approach is to utilize a penalty function so that a constraint violation comes with a price; as such the SQP method ties in with the penalty methods discussed above.

1.10 On the background of this book and a didactics statement

This book's foundation is the collection of lecture notes written by the third author and used in the basic optimization course "Applied Optimization" for nearly ten years at Chalmers University of Technology and Gothenburg University. The lecture notes have developed more and more from being based on algorithms to mainly covering the fundamentals of optimization. With the addition of the first two authors has come a further development of these fundamentals into the present book, in which also our didactic wishes has begun to come true; the present book significantly expands and improves upon the initial lecture notes.

The main inspiration in shaping the lecture notes and the book came from the excellent text book by Bazaraa, Sherali, and Shetty [BSS93]. In the book the authors separate the basic theory (convexity, polyhedral theory, separation, optimality, etcetera) from the algorithms devised for solving nonlinear optimization problems, and they develop the theory based on first principles, in a natural order. (The book is however too advanced to be used in a first optimization course, it does not cover linear programming, and some of the algorithmic parts are getting old.)

The main focus, as the title suggests, is the foundation of optimization models and methods. Hence, we have developed the chapters on convexity and optimality conditions in detail. On the other hand, with the exception of the classic topic of linear programming we have strived to keep the algorithmic chapters concise, yet rigorous; among the plentiful of possible choices of algorithms we have made those choices that appear the most natural given the appearance of the optimality conditions. The choices have therefore also become those of classic algorithms rather than the most advanced and modern ones; being an undergraduate text we find this to be appropriate, and our text therefore also serves as fundamental material that paves the way for more advanced

text books on optimization methods. Among those we mention especially that of Nocedal and Wright [NoW99], whose excellent graduate level book on numerical optimization also is developed through a careful selection of algorithms.

In writing the book we have also made a few additional didactic developments. In almost every text book on optimization the topic of linear programming is developed before that of nonlinear and convex optimization, and linear programming duality is developed before Lagrangian duality. Teaching in this order may however feel unnatural both for instructors and students: since Lagrangian duality is more general, but similar, to linear programming duality, the feeling is that more or less the same material is repeated, or (which is even worse) that linear programming is a rather strange special topic that we develop because we must, but not because it is interesting. We have developed the material in this book such that linear programming emerges as a natural special case of general convex programming, and having a duality theory which is even richer.

In keeping with this idea of developing nonlinear programming before linear programming, we should also have covered the simplex method last in the book. This is a possibly conflicting situation, because we believe that the simplex method should not be described merely as a feasible-direction method; its combinatorial nature is important, and the subject of degeneracy, for example, is more naturally treated and understood by developing the simplex method immediately following the development of the connections between the geometry and linear algebra of linear programming. This has been our choice, and we have consequently also decided that iterative algorithms for general nonlinear optimization over convex sets, especially polyhedra, should be developed before those for more general constraints, the reason being that linear programming is an important basis for these algorithms.

When teaching from this book, we have decided to stick to the chapter ordering with one exception: we introduce Chapter 11 as well as hands-on computer exercises on algorithms for unconstrained optimization immediately after teaching from Chapter 4 on optimality conditions for problems over convex sets. The motivation for doing so is our wish to integrate, in our teaching, algorithms with fundamental theory; the book itself separates the two topics.

1.11 Illustrating the theory

The subject of optimization, including both its basic theory and the natural, basic, algorithmic development that is associated with solving

different classes of optimization models, is special compared to many other mathematics subjects in that the ties between analysis/algebra and geometry are so strong. This means, particularly, that optimization can be learned, illustrated and revised (at least partially) by using geometric tools. We give a few such examples.

The various techniques available for checking the convexity of a set or a function can be illustrated by examples in one or two dimensions. All the necessary and sufficient conditions for local optimality in constrained and unconstrained optimization given in Chapters 4 and 5 can thus be illustrated. A simple method in \mathbb{R}^2 is as follows: choose a (suitably many times) differentiable function f such that a minimum over \mathbb{R}^2 is known. If the test problem should be unconstrained, one is immediately ready to work with the corresponding instruments; if the objective function should be minimized subject to constraints, then choose the feasible set such that the “constrained” optimum is different from the “unconstrained” one and use the corresponding optimality conditions to check that the optimal solution indeed satisfies them, or that an arbitrarily chosen non-optimal vector does not. The constraint qualifications (CQs), which play an important role for general sets, can also be investigated through such examples.

In linear programming much of the above is specialized, since duality and the KKT conditions have their correspondence in linear programming duality and optimality. A two-dimensional polyhedron, together with a suitable objective function, can illustrate primal–dual relationships such as the complementarity conditions, based on a problem with a known solution; it can also test one’s mastering of the simplex method.

The algorithmic chapters in Part V are similar with respect to these tests; for each problem class and algorithm, it is possible, and instrumental, to construct a two-dimensional example and check that the algorithm in question will reach a stationary point, if the convergence conditions are met, or disprove convergence when the conditions are not. This also provokes a revision of the optimality conditions of Chapters 4 and 5.

The variety of examples that can be thus constructed is immense. This is in fact one of the reasons why we have decided to limit the number of exercises; one can in fact create one’s own set of exercises, and will benefit greatly from doing so.

1.12 Notes and further reading

Extensive collections of optimization applications and models can be found in several basic text books in operations research, such as [Wag75, BHM77, Mur95, Rar98, Tah03]. The optimization modelling book by

Williams [Wil99] is a classic, now in its fourth edition. Modelling books also exist for certain categories of applications; for example, the book [EHL01] concerns the mathematical modelling and solution of optimization problems arising in chemical engineering applications. Further industrial applications are found in [AvG96, Casetal02].

Several accounts have been written during the past few years on the origins of operations research and mathematical programming, the reasons being that we recently celebrated the 50th anniversaries of the simplex method (1997), the creation of ORSA (the Operations Research Society of America) (2002), and the Operational Research Society (2003), as well as the 90th birthday of the inventor of the simplex method, George Dantzig (2004). The special issue of the journal *Operations Research*, vol. 50, no. 1 (2002), is filled with historical anecdotes, as are the books [LRS91, GaA05] on the history of mathematical programming and operations research.

1.13 Exercises

Exercise 1.1 (modelling, exam 980819) A new producer of perfume wish to get a break into a lucrative market. An exclusive fragrance, Chinelle, is to be produced and marketed. With the equipment available it is possible to produce the perfume using two alternative processes, and the company also consider utilizing the services of a famous model when launching it. In order to simplify the problem, let us assume that the perfume is manufactured by the use of two main ingredients—the first a secret substance called MO and the second a more well-known mixture of ingredients. The first of the two processes available provides three grams of perfume for every unit of MO and two units of the standard substance, while the other process gives five grams of perfume for every two (respectively, three) units of the two main ingredients. The company has at its disposal manufacturing processes that can produce at most 20,000 units of MO during the planning period and 35,000 units of the standard mixture. Every unit of MO costs three EUR (it is manufactured in France) to produce, and the other mixture only two EUR per unit. One gram of the new perfume will cost fifty EUR. Even without any advertising the company thinks they can sell 1000 grams of the perfume, simply because of the news value. A famous model can be contracted for commercials, costing 5,000 EUR per photo session (which takes half an hour), and the company thinks that a campaign using his image can raise the demand by about 200 grams per half hour of his time, but not exceeding three hours (he has many other offers).

Formulate an LP model of the best production strategy problem.

Exercise 1.2 (modelling) A computer company has estimated the service hours needed during the next five months; see Table 1.2.

Table 1.2: Number of service hours per month; Exercise 1.2.

Month	# Service hours
January	6000
February	7000
March	8000
April	9500
May	11,500

The service is performed by hired technicians; their number is 50 at the beginning of January. Each technician can work up to 160 hours per month. In order to cover the future demand of technicians new ones must be hired. Before a technician is hired he/she undergoes a period of training, which takes a month and requires 50 hours of supervision by a trained technician. A trained technician has a salary of 15,000 SEK per month (regardless of the number of working hours) and a trainee has a monthly salary of 7500 SEK. At the end of each month on average 5% of the technicians quit to work for another company.

Formulate a linear integer program whose optimal solution will minimize the total salary costs during the given time period, given that the number of available service hours are enough to cover the demand.

Exercise 1.3 (modelling, exam 010821) The advertising agency ZAP (Zetterström, Anderson, and Pettersson) is designing their new office with an open office space. The office is rectangular, with length l meters and width b meters. Somewhat simplified, we may assume that each working space requires a circle of diameter d and that the working spaces must not overlap. In addition, each working space must be connected to the telecom and computer network at one of the two possible connection points in the office. As the three telephones have limited cable lengths (the agency is concerned with the possible radiation danger associated with hands-free phones and therefore do not use cord-less phones)— a_i meters, respectively, $i = 1, \dots, 3$ —the work spaces must be placed quite near the connection points.⁹ See Figure 1.3 for a picture of the office.

For simplicity we assume that the phone is placed at the center of the work place. One of the office's walls is a large panorama window and the

⁹All the money went to other interior designs of the office space, so there is no money left to buy more cable.

Modelling and classification

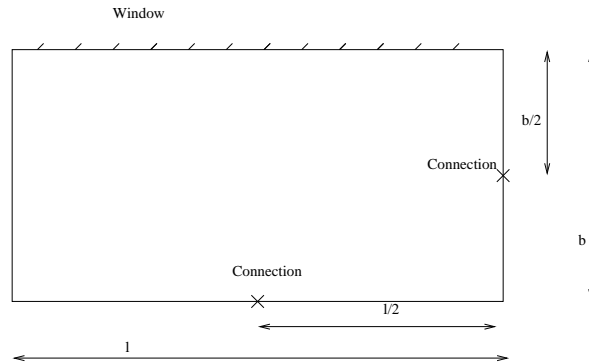


Figure 1.3: Image of the office; Exercise 1.3

three partners all want to sit as close as possible to it. Therefore, they decide to try to minimize the distance to the window for the workplace that is the furthest away from it.

Formulate the problem of placing the three work places so that the maximum distance to the panorama window is minimized, subject to all the necessary constraints.

Exercise 1.4 (modelling, exam 010523) A large chain of department stores wants to build distribution centers (warehouses) which will supply 30 department stores with goods. They have 10 possible locations to choose between. To build a warehouse at location i , $i = 1, \dots, 10$, costs c_i MEUR and the capacity of a warehouse at that location would be k_i volume units per week. Department store j has a demand of e_j volume units per week. The distance between warehouse i and department store j is d_{ij} km, $i = 1, \dots, 10$, $j = 1, \dots, 30$, and a certain warehouse can only serve a department store if the distance is at most D km.

One wishes to minimize the cost of investing in the necessary distribution centers.

(a) Formulate a *linear integer optimization model* describing the optimization problem.

(b) Suppose each department store must be served from *one* of the warehouses. What must be changed in the model?

Part II

Fundamentals

Analysis and algebra—A summary

II

The analysis of optimization problems and related optimization algorithms requires the basic understanding of formal logic, linear algebra, and multidimensional analysis. This chapter is not intended as a substitute for the basic courses on these subjects but rather to give a brief review of the notation, definitions, and basic facts which will be used in the subsequent chapters without any further notice. If you feel inconvenient with the limited summaries presented in this chapter, contact any of the abundant number of basic text books on the subject.

2.1 Reductio ad absurdum

Together with the absolute majority of contemporary mathematicians we accept proofs by contradiction. The proofs in this group essentially appeal to Aristotle's law of the excluded middle, which states that any proposition is either true or false. Thus, if some statement can be shown to lead to a contradiction, we conclude that the original statement is false.

Formally, proofs by contradiction amount to the following:

$$(A \implies B) \iff (\neg A \vee B) \iff (\neg\neg B \vee \neg A) \iff (\neg B \implies \neg A).$$

In the same spirit, when proving $A \iff B$, that is, $(A \implies B) \wedge (B \implies A)$, we often equivalently argue according to $(A \implies B) \wedge (\neg A \implies \neg B)$ (see, for example, the proof of Farkas' Lemma 3.30).

2.2 Linear algebra

We will always work with finite dimensional Euclidean vector spaces \mathbb{R}^n , the natural number n denoting the dimension of the space. Elements $\mathbf{v} \in \mathbb{R}^n$ will be referred to as *vectors*, and we will always think of them as of n real numbers stacked on top of each other, i.e., $\mathbf{v} = (v_1, \dots, v_n)^T$, v_i being real numbers, and T denoting the “transpose” sign. The basic operations defined for two vectors $\mathbf{a} = (a_1, \dots, a_n)^T \in \mathbb{R}^n$ and $\mathbf{b} = (b_1, \dots, b_n)^T \in \mathbb{R}^n$, and an arbitrary scalar $\alpha \in \mathbb{R}$ are as follows:

- addition: $\mathbf{a} + \mathbf{b} := (a_1 + b_1, \dots, a_n + b_n)^T \in \mathbb{R}^n$;
- multiplication by a scalar: $\alpha \mathbf{a} := (\alpha a_1, \dots, \alpha a_n)^T \in \mathbb{R}^n$;
- *scalar product* between two vectors: $(\mathbf{a}, \mathbf{b}) := \sum_{i=1}^n a_i b_i \in \mathbb{R}$. The scalar product will most often be denoted by $\mathbf{a}^T \mathbf{b}$ in the subsequent chapters.

A *linear subspace* $L \subseteq \mathbb{R}^n$ is a set enjoying the following two properties:

- for every $\mathbf{a}, \mathbf{b} \in L$ it holds that $\mathbf{a} + \mathbf{b} \in L$, and
- for every $\alpha \in \mathbb{R}, \mathbf{a} \in L$ it holds that $\alpha \mathbf{a} \in L$.

An *affine subspace* $A \subseteq \mathbb{R}^n$ is any set that can be represented as $\mathbf{v} + L := \{\mathbf{v} + \mathbf{x} \mid \mathbf{x} \in L\}$ for some vector $\mathbf{v} \in \mathbb{R}^n$ and some linear subspace $L \subseteq \mathbb{R}^n$.

We associate the *norm*, or length, of a vector $\mathbf{v} \in \mathbb{R}^n$ with the following scalar product:

$$\|\mathbf{v}\| := \sqrt{(\mathbf{v}, \mathbf{v})}.$$

We will sometimes write $|\mathbf{v}|$ in place of $\|\mathbf{v}\|$. The Cauchy–Bunyakowski–Schwarz inequality says that $(\mathbf{a}, \mathbf{b}) \leq \|\mathbf{a}\| \|\mathbf{b}\|$ for $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$; thus, we may define an angle θ between two vectors via $\cos \theta := (\mathbf{a}, \mathbf{b}) / (\|\mathbf{a}\| \|\mathbf{b}\|)$. We say that $\mathbf{a} \in \mathbb{R}^n$ is *orthogonal* to $\mathbf{b} \in \mathbb{R}^n$ if and only if $(\mathbf{a}, \mathbf{b}) = 0$ (i.e., when $\cos \theta = 0$). The only vector orthogonal to itself is the zero vector $\mathbf{0}^n := (0, \dots, 0)^T \in \mathbb{R}^n$; moreover, this is the only vector having a zero norm.

The scalar product is *symmetric* and *bilinear*, i.e., for every $\mathbf{a}, \mathbf{b}, \mathbf{c}, \mathbf{d} \in \mathbb{R}^n$, $\alpha, \beta, \gamma, \delta \in \mathbb{R}$ it holds that $(\mathbf{a}, \mathbf{b}) = (\mathbf{b}, \mathbf{a})$, and $(\alpha \mathbf{a} + \beta \mathbf{b}, \gamma \mathbf{c} + \delta \mathbf{d}) = \alpha \gamma (\mathbf{a}, \mathbf{c}) + \beta \gamma (\mathbf{b}, \mathbf{c}) + \alpha \delta (\mathbf{a}, \mathbf{d}) + \beta \delta (\mathbf{b}, \mathbf{d})$.

A collection of vectors $(\mathbf{v}_1, \dots, \mathbf{v}_k)$ is said to be *linearly independent* if and only if the equality $\sum_{i=1}^k \alpha_i \mathbf{v}_i = \mathbf{0}^n$, where $\alpha_1, \dots, \alpha_k$ are arbitrary real numbers, implies that $\alpha_1 = \dots = \alpha_k = 0$. Similarly, a collection of vectors $(\mathbf{v}_1, \dots, \mathbf{v}_k)$ is said to be *affinely independent* if and only if the collection $(\mathbf{v}_2 - \mathbf{v}_1, \dots, \mathbf{v}_k - \mathbf{v}_1)$ is linearly independent.

The largest number of linearly independent vectors in \mathbb{R}^n is n ; any collection of n linearly independent vectors in \mathbb{R}^n is referred to as a *basis*. The basis $(\mathbf{v}_1, \dots, \mathbf{v}_n)$ is said to be *orthogonal* if $(\mathbf{v}_i, \mathbf{v}_j) = 0$ for all $i, j = 1, \dots, n, i \neq j$. If, in addition, it holds that $\|\mathbf{v}_i\| = 1$ for all $i = 1, \dots, n$, the basis is called *orthonormal*.

Given the basis $(\mathbf{v}_1, \dots, \mathbf{v}_n)$ in \mathbb{R}^n , every vector $\mathbf{v} \in \mathbb{R}^n$ can be written in a unique way as $\mathbf{v} = \sum_{i=1}^n \alpha_i \mathbf{v}_i$, and the n -tuple $(\alpha_1, \dots, \alpha_n)^T$ will be referred to as *coordinates* of \mathbf{v} in this basis. If the basis $(\mathbf{v}_1, \dots, \mathbf{v}_n)$ is orthonormal, then the coordinates α_i are computed as $\alpha_i = (\mathbf{v}, \mathbf{v}_i)$, $i = 1, \dots, n$.

The space \mathbb{R}^n will typically be equipped with the *standard basis* $(\mathbf{e}_1, \dots, \mathbf{e}_n)$, where

$$\mathbf{e}_i := (\underbrace{0, \dots, 0}_{i-1 \text{ zeros}}, 1, \underbrace{0, \dots, 0}_{n-i \text{ zeros}})^T \in \mathbb{R}^n.$$

This basis is orthogonal, and for every vector $\mathbf{v} = (v_1, \dots, v_n)^T \in \mathbb{R}^n$ we have $(\mathbf{v}, \mathbf{e}_i) = v_i$, $i = 1, \dots, n$, which allow us to identify vectors and their coordinates.

Now, consider two spaces \mathbb{R}^n and \mathbb{R}^k . All linear functions from \mathbb{R}^n to \mathbb{R}^k may be described using a linear space of *real matrices* $\mathbb{R}^{k \times n}$ (i.e., with k rows and n columns). Given a matrix $\mathbf{A} \in \mathbb{R}^{k \times n}$ it will often be convenient to view it as a row of its columns, which are thus vectors in \mathbb{R}^k . Namely, let $\mathbf{A} \in \mathbb{R}^{k \times n}$ have elements a_{ij} , $i = 1, \dots, k, j = 1, \dots, n$, then we write $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_n)$, where $\mathbf{a}_i := (a_{1i}, \dots, a_{ki})^T \in \mathbb{R}^k$, $i = 1, \dots, n$. The addition of two matrices and scalar-matrix multiplication are defined in a straightforward way. For $\mathbf{v} = (v_1, \dots, v_n) \in \mathbb{R}^n$ we define $\mathbf{A}\mathbf{v} = \sum_{i=1}^n v_i \mathbf{a}_i \in \mathbb{R}^k$, where $\mathbf{a}_i \in \mathbb{R}^k$ are the columns of \mathbf{A} . We also define the *norm* of the matrix \mathbf{A} by

$$\|\mathbf{A}\| := \max_{\mathbf{v} \in \mathbb{R}^n: \|\mathbf{v}\|=1} \|\mathbf{A}\mathbf{v}\|.$$

Well, this is an example of an optimization problem already!

For a given matrix $\mathbf{A} \in \mathbb{R}^{k \times n}$ with elements a_{ij} we define $\mathbf{A}^T \in \mathbb{R}^{n \times k}$ as the matrix with elements $\tilde{a}_{ij} := a_{ji}$ $i = 1, \dots, n, j = 1, \dots, k$. We can give a more elegant, but less straightforward definition: \mathbf{A}^T is the unique matrix, satisfying the equality $(\mathbf{A}\mathbf{v}, \mathbf{u}) = (\mathbf{v}, \mathbf{A}^T \mathbf{u})$ for all $\mathbf{v} \in \mathbb{R}^n$, $\mathbf{u} \in \mathbb{R}^k$. From this definition it should be clear that $\|\mathbf{A}\| = \|\mathbf{A}^T\|$, and that $(\mathbf{A}^T)^T = \mathbf{A}$.

Given two matrices $\mathbf{A} \in \mathbb{R}^{k \times n}$ and $\mathbf{B} \in \mathbb{R}^{n \times m}$, we define the *product* $\mathbf{C} = \mathbf{A}\mathbf{B} \in \mathbb{R}^{k \times m}$ element-wise by $c_{ij} = \sum_{\ell=1}^n a_{i\ell} b_{\ell j}$, $i = 1, \dots, k, j = 1, \dots, m$. In other words, $\mathbf{C} = \mathbf{A}\mathbf{B}$ if and only if for all $\mathbf{v} \in \mathbb{R}^n$, $\mathbf{C}\mathbf{v} = \mathbf{A}(\mathbf{B}\mathbf{v})$. By definition, the matrix product is associative

(that is, $\mathbf{A}(\mathbf{BC}) = (\mathbf{AB})\mathbf{C}$) for matrices of compatible sizes, but *not* commutative (that is, $\mathbf{AB} \neq \mathbf{BA}$) in general. It is easy (and instructive) to check that $\|\mathbf{AB}\| \leq \|\mathbf{A}\|\|\mathbf{B}\|$, and that $(\mathbf{AB})^T = \mathbf{B}^T\mathbf{A}^T$. Vectors $\mathbf{v} \in \mathbb{R}^n$ can be (and sometimes will be) viewed as matrices $\mathbf{v} \in \mathbb{R}^{n \times 1}$. Check that this embedding is norm-preserving, that is, the norm of \mathbf{v} viewed as a vector equals the norm of \mathbf{v} viewed as a matrix with one column.

Of course, no discussion about norms could escape mentioning the *triangle inequality*: for all $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ it holds that $\|\mathbf{a} + \mathbf{b}\| \leq \|\mathbf{a}\| + \|\mathbf{b}\|$, as well as its consequence (check this!) that for all $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{k \times n}$, $\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|$ holds. It will often be used in a little bit different form: for all $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$, $\|\mathbf{b}\| - \|\mathbf{a}\| \leq \|\mathbf{b} - \mathbf{a}\|$ holds.

For square matrices $\mathbf{A} \in \mathbb{R}^{n \times n}$ we can discuss the existence of the unique matrix \mathbf{A}^{-1} , called the *inverse* of \mathbf{A} , verifying the equality that for all $\mathbf{v} \in \mathbb{R}^n$, $\mathbf{A}^{-1}\mathbf{A}\mathbf{v} = \mathbf{A}\mathbf{A}^{-1}\mathbf{v} = \mathbf{v}$ holds. If the inverse of a given matrix exists, we call the latter *nonsingular*. The inverse matrix exists if and only if the columns of \mathbf{A} are linearly independent; if and only if the columns of \mathbf{A}^T are linearly independent; if and only if the system of linear equations $\mathbf{A}\mathbf{x} = \mathbf{v}$ has a unique solution for every $\mathbf{v} \in \mathbb{R}^n$; if and only if the homogeneous system of equations $\mathbf{A}\mathbf{x} = \mathbf{0}^n$ has $\mathbf{x} = \mathbf{0}^n$ as its unique solution. From this definition it follows that \mathbf{A} is nonsingular if and only if \mathbf{A}^T is nonsingular, and, furthermore, $(\mathbf{A}^{-1})^T = (\mathbf{A}^T)^{-1}$ and therefore will be denoted simply by \mathbf{A}^{-T} . At last, if \mathbf{A} and \mathbf{B} are two nonsingular square matrices of the same size, then \mathbf{AB} is nonsingular (why?) and $(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$.

If, for some vector $\mathbf{v} \in \mathbb{R}^n$ and scalar $\alpha \in \mathbb{R}$ it holds that $\mathbf{A}\mathbf{v} = \alpha\mathbf{v}$, then we call \mathbf{v} an *eigenvector* of \mathbf{A} , corresponding to the *eigenvalue* α of \mathbf{A} . Eigenvectors, corresponding to a given eigenvalue, form a linear *subspace* of \mathbb{R}^n ; two nonzero eigenvectors, corresponding to two distinct eigenvalues are linearly independent. In general, every matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ has n eigenvalues (counted with multiplicity), maybe complex, which are furthermore roots of the *characteristic equation* $\det(\mathbf{A} - \lambda\mathbf{I}^n) = 0$, where $\mathbf{I}^n \in \mathbb{R}^{n \times n}$ is the *identity matrix*, characterized by the fact that for all $\mathbf{v} \in \mathbb{R}^n$ it holds that $\mathbf{I}^n\mathbf{v} = \mathbf{v}$. The norm of the matrix \mathbf{A} is in fact equal to the largest absolute value of its eigenvalues. The matrix \mathbf{A} is nonsingular if and only if none of its eigenvalues are equal to zero, and in this case the eigenvalues of \mathbf{A}^{-1} are equal to the inverted eigenvalues of \mathbf{A} . The eigenvalues of \mathbf{A}^T are equal to the eigenvalues of \mathbf{A} .

We call \mathbf{A} *symmetric* if and only if $\mathbf{A}^T = \mathbf{A}$. All eigenvalues of symmetric matrices are real, and eigenvectors corresponding to distinct eigenvalues are orthogonal.

Even if \mathbf{A} is not square, $\mathbf{A}^T \mathbf{A}$ as well as $\mathbf{A} \mathbf{A}^T$ are square and symmetric. If the columns of \mathbf{A} are linearly independent, then $\mathbf{A}^T \mathbf{A}$ is nonsingular. (Similarly, if the columns of \mathbf{A}^T are linearly independent, then $\mathbf{A} \mathbf{A}^T$ is nonsingular.)

Sometimes, we will use the following simple fact: for every $\mathbf{A} \in \mathbb{R}^{k \times n}$ with elements a_{ij} , $i = 1, \dots, k$, $j = 1, \dots, n$, it holds that $a_{ij} = (\tilde{\mathbf{e}}_i, \mathbf{A} \mathbf{e}_j)$, where $(\tilde{\mathbf{e}}_1, \dots, \tilde{\mathbf{e}}_k)$ is the standard basis in \mathbb{R}^k , and $(\mathbf{e}_1, \dots, \mathbf{e}_n)$ is the standard basis in \mathbb{R}^n , $i = 1, \dots, k$, $j = 1, \dots, n$.

We will say that $\mathbf{A} \in \mathbb{R}^{n \times n}$ is *positive semidefinite* (respectively, *positive definite*), and denote this by $\mathbf{A} \succeq \mathbf{0}$ (respectively, $\mathbf{A} \succ \mathbf{0}$) if and only if for all $\mathbf{v} \in \mathbb{R}^n$ it holds that $(\mathbf{v}, \mathbf{A} \mathbf{v}) \geq 0$ (respectively, for all $\mathbf{v} \in \mathbb{R}^n$, $\mathbf{v} \neq \mathbf{0}^n$, it holds that $(\mathbf{v}, \mathbf{A} \mathbf{v}) > 0$). The matrix \mathbf{A} is positive semidefinite (respectively, positive definite) if and only if its eigenvalues are nonnegative (respectively, positive).

For two symmetric matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$ we will write $\mathbf{A} \succeq \mathbf{B}$ (respectively, $\mathbf{A} \succ \mathbf{B}$) if and only if $\mathbf{A} - \mathbf{B} \succeq \mathbf{0}$ (respectively, $\mathbf{A} - \mathbf{B} \succ \mathbf{0}$).

2.3 Analysis

Consider a sequence $\{\mathbf{x}_k\} \subset \mathbb{R}^n$. We will write $\lim_{k \rightarrow \infty} \mathbf{x}_k = \mathbf{x}$, for some $\mathbf{x} \in \mathbb{R}^n$, or just $\mathbf{x}_k \rightarrow \mathbf{x}$, if and only if $\lim_{k \rightarrow \infty} \|\mathbf{x}_k - \mathbf{x}\| = 0$. We will say in this case that $\{\mathbf{x}_k\}$ *converges to \mathbf{x}* , or, equivalently, that \mathbf{x} is the *limit* of $\{\mathbf{x}_k\}$. Owing to the triangle inequality, every sequence might have at most one limit. (Why?) At the same time, there are sequences that do not converge. Moreover, an arbitrary non-converging sequence might contain a converging subsequence (or even several subsequences). We will refer to the limits of such converging subsequences as *limit points* of a given sequence $\{\mathbf{x}_k\}$.

A subset $S \subset \mathbb{R}^n$ is called *bounded* if there exists a constant $C > 0$ such that for all $\mathbf{x} \in S$: $\|\mathbf{x}\| \leq C$; otherwise, the set will be called *unbounded*. Now, let $S \subset \mathbb{R}^n$ be bounded. An interesting and very important fact about bounded subsets $S \subset \mathbb{R}^n$ is that every sequence $\{\mathbf{x}_k\} \subset S$ contains a convergent subsequence.

The set $B_\varepsilon(\mathbf{x}) = \{\mathbf{y} \in \mathbb{R}^n \mid \|\mathbf{x} - \mathbf{y}\| < \varepsilon\}$ is called an *open ball* of radius $\varepsilon > 0$ with center $\mathbf{x} \in \mathbb{R}^n$. A set $S \subseteq \mathbb{R}^n$ is called *open* if and only if for all $\mathbf{x} \in S$ there exists an $\varepsilon > 0$ such that $B_\varepsilon(\mathbf{x}) \subset S$. A set S is *closed* if and only if its complement $\mathbb{R}^n \setminus S$ is open. An equivalent definition of closedness in terms of sequences is that a set $S \subseteq \mathbb{R}^n$ is closed if and only if all the limit points of any sequence $\{\mathbf{x}_k\} \subset S$ belong to S . There exist sets which are neither closed nor open. The set \mathbb{R}^n is both open and closed. (Why?)

The *closure* of a set $S \subseteq \mathbb{R}^n$ (notation: $\text{cl } S$) is the smallest closed

set containing S ; equivalently, it can be defined as the intersection of all closed sets in \mathbb{R}^n containing S . More constructively, the closure $\text{cl } S$ can be obtained by considering all limit points of all sequences in S . The closure is a closed set, and, quite naturally, the closure of a closed set equals the set itself.

The *interior* of a set $S \subseteq \mathbb{R}^n$ (notation: $\text{int } S$) is the largest open set contained in S . The interior of an open set equals the set itself.

Finally, the *boundary* of a set $S \subseteq \mathbb{R}^n$ (notation: $\text{bd } S$, or ∂S) is the set difference $\text{cl } S \setminus \text{int } S$.

A *neighbourhood* of a point $\mathbf{x} \in \mathbb{R}^n$ is an arbitrary open set containing \mathbf{x} .

Consider a function $f : S \rightarrow \mathbb{R}$, where $S \subseteq \mathbb{R}^n$. We say that f is *continuous* at $\mathbf{x}_0 \in S$ if and only if for every sequence $\{\mathbf{x}_k\} \subset S$ such that $\mathbf{x}_k \rightarrow \mathbf{x}_0$ it holds that $\lim_{k \rightarrow \infty} f(\mathbf{x}_k) = f(\mathbf{x}_0)$. We say that f is continuous on S if and only if f is continuous at every point of S .

Now, let $f : S \rightarrow \mathbb{R}$ be a continuous function defined on some open set S . We say that $f'(\mathbf{x}_0; \mathbf{d}) \in \mathbb{R}$ is a *directional derivative* of f at $\mathbf{x}_0 \in S$ in the direction $\mathbf{d} \in \mathbb{R}^n$ if the following limit exists:

$$f'(\mathbf{x}_0, \mathbf{d}) = \lim_{t \downarrow 0} \frac{f(\mathbf{x}_0 + t\mathbf{d}) - f(\mathbf{x}_0)}{t},$$

and then f will be called *directionally differentiable* at $\mathbf{x}_0 \in S$ in the direction \mathbf{d} . Clearly, if we fix $\mathbf{x}_0 \in S$ and assume that $f'(\mathbf{x}_0; \mathbf{d})$ exists for some \mathbf{d} , then for every $\alpha \geq 0$ we have that $f'(\mathbf{x}_0; \alpha\mathbf{d}) = \alpha f'(\mathbf{x}_0; \mathbf{d})$. If further $f'(\mathbf{x}_0; \mathbf{d})$ is linear in \mathbf{d} , then there exists a vector called the *gradient* of f at $\mathbf{x}_0 \in S$, denoted by $\nabla f(\mathbf{x}_0) \in \mathbb{R}^n$, such that $f'(\mathbf{x}_0; \mathbf{d}) = (\nabla f(\mathbf{x}_0), \mathbf{d})$ and f is then called *differentiable* at $\mathbf{x}_0 \in S$. Naturally, we say that f is differentiable on S if it is differentiable at every point in S .

Equivalently, the gradient can be defined as follows: $\nabla f(\mathbf{x}_0) \in \mathbb{R}^n$ is the gradient of f at \mathbf{x}_0 if and only if there exists a function $o : \mathbb{R} \rightarrow \mathbb{R}$ such that

$$f(\mathbf{x}) = f(\mathbf{x}_0) + (\nabla f(\mathbf{x}_0), \mathbf{x} - \mathbf{x}_0) + o(\|\mathbf{x} - \mathbf{x}_0\|), \quad (2.1)$$

and moreover

$$\lim_{t \downarrow 0} \frac{o(t)}{t} = 0. \quad (2.2)$$

For a differentiable function $f : S \rightarrow \mathbb{R}$ we can go one step further and define second derivatives of f . Namely, a differentiable function f will be called *twice differentiable* at $\mathbf{x}_0 \in S$ if and only if there exists a symmetric matrix denoted by $\nabla^2 f(\mathbf{x}_0)$, and referred to as the *Hessian*

matrix, and a function $o : \mathbb{R} \rightarrow \mathbb{R}$ verifying (2.2), such that

$$\begin{aligned} f(\mathbf{x}) &= f(\mathbf{x}_0) + (\nabla f(\mathbf{x}_0), \mathbf{x} - \mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0, \nabla^2 f(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0)) \\ &\quad + o(\|\mathbf{x} - \mathbf{x}_0\|^2). \end{aligned} \quad (2.3)$$

Sometimes it will be convenient to discuss *vector-valued* functions $\mathbf{f} : S \rightarrow \mathbb{R}^k$. We say that $\mathbf{f} = (f_1, \dots, f_k)^\top$ is continuous if every f_i , $i = 1, \dots, k$ is; similarly we define differentiability. In the latter case, by $\nabla \mathbf{f} \in \mathbb{R}^{n \times k}$ we denote a matrix with columns $(\nabla f_1, \dots, \nabla f_k)$. Its transpose is often referred to as the *Jacobian* of \mathbf{f} .

We call a continuous function $f : S \rightarrow \mathbb{R}$ *continuously differentiable* [notation: $f \in C^1(S)$] if it is differentiable on S and the gradient $\nabla f : S \rightarrow \mathbb{R}^n$ is continuous on S . We call $f : S \rightarrow \mathbb{R}$ *twice continuously differentiable* [notation: $f \in C^2(S)$], if it is continuously differentiable and in addition every component of $\nabla f : S \rightarrow \mathbb{R}^n$ is continuously differentiable.

The following alternative forms of (2.1) and (2.3) will be useful some times. If $f : S \rightarrow \mathbb{R}$ is continuously differentiable on S , and $\mathbf{x}_0 \in S$, then for every \mathbf{x} in some neighbourhood of \mathbf{x}_0 we have

$$f(\mathbf{x}) = f(\mathbf{x}_0) + (\nabla f(\boldsymbol{\xi}), \mathbf{x} - \mathbf{x}_0), \quad (2.4)$$

where $\boldsymbol{\xi} = \lambda \mathbf{x}_0 + (1 - \lambda)\mathbf{x}$, for some $0 \leq \lambda \leq 1$, is a point between \mathbf{x} and \mathbf{x}_0 . (This result is also known as the *mean-value theorem*.) Similarly, for twice differentiable functions we have

$$f(\mathbf{x}) = f(\mathbf{x}_0) + (\nabla f(\mathbf{x}_0), \mathbf{x} - \mathbf{x}_0) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_0, \nabla^2 f(\boldsymbol{\xi})(\mathbf{x} - \mathbf{x}_0)), \quad (2.5)$$

with the same notation.

If $f, g : \mathbb{R}^n \rightarrow \mathbb{R}$ are both differentiable, then $f + g$ and fg are, and $\nabla(f + g) = \nabla f + \nabla g$, $\nabla(fg) = g\nabla f + f\nabla g$. Moreover, if g is never zero, then f/g is differentiable and $\nabla(f/g) = (g\nabla f - f\nabla g)/g^2$.

If both $\mathbf{F} : \mathbb{R}^n \rightarrow \mathbb{R}^k$ and $h : \mathbb{R}^k \rightarrow \mathbb{R}$ are differentiable, then $h(\mathbf{F})$ is, and $(\nabla h(\mathbf{F}))(\mathbf{x}) = (\nabla \mathbf{F})(\mathbf{x}) \cdot (\nabla h)(\mathbf{F}(\mathbf{x}))$.

Finally, consider a vector-valued function $\mathbf{F} : \mathbb{R}^k \times \mathbb{R}^n \rightarrow \mathbb{R}^k$. Assume that \mathbf{F} is continuously differentiable in some neighbourhood $\mathcal{N}_u \times \mathcal{N}_x$ of the point $(\mathbf{u}_0, \mathbf{x}_0) \in \mathbb{R}^k \times \mathbb{R}^n$, and that $\mathbf{F}(\mathbf{u}_0, \mathbf{x}_0) = \mathbf{0}^k$. If the square matrix $\nabla_u \mathbf{F}(\mathbf{u}_0, \mathbf{x}_0)$ is nonsingular, then there exists a unique function $\boldsymbol{\varphi} : \mathcal{N}'_x \rightarrow \mathcal{N}'_u$ such that $\mathbf{F}(\boldsymbol{\varphi}(\mathbf{x}), \mathbf{x}) \equiv \mathbf{0}^k$ in \mathcal{N}'_x , where $\mathcal{N}'_u \times \mathcal{N}'_x \subset \mathcal{N}_u \times \mathcal{N}_x$ is another neighbourhood of $(\mathbf{u}_0, \mathbf{x}_0)$. Furthermore, $\boldsymbol{\varphi}$ is differentiable at \mathbf{x}_0 , and

$$\nabla \boldsymbol{\varphi}(\mathbf{x}_0) = -(\nabla_u \mathbf{F}(\mathbf{u}_0, \mathbf{x}_0))^{-1} \nabla_x \mathbf{F}(\mathbf{u}_0, \mathbf{x}_0).$$

The function φ is known as the *implicit function* defined by the system of equations $\mathbf{F}(\mathbf{u}, \mathbf{x}) = \mathbf{0}^k$.

Now we consider two special but very important cases.

Let for some $\mathbf{a} \in \mathbb{R}^n$ define a *linear function* $f : \mathbb{R}^n \rightarrow \mathbb{R}$ via $f(\mathbf{x}) := (\mathbf{a}, \mathbf{x})$. By the Cauchy–Bunyakowski–Schwarz inequality this function is continuous, and writing $f(\mathbf{x}) - f(\mathbf{x}_0) = (\mathbf{a}, \mathbf{x} - \mathbf{x}_0)$ for every $\mathbf{x}_0 \in \mathbb{R}^n$ we immediately identify from the definitions of the gradient and the Hessian that $\nabla f = \mathbf{a}$, $\nabla^2 f = \mathbf{0}^{n \times n}$.

Similarly, for some $\mathbf{A} \in \mathbb{R}^{n \times n}$ define a *quadratic function* $f(\mathbf{x}) = (\mathbf{x}, \mathbf{A}\mathbf{x})$. This function is also continuous, and since $f(\mathbf{x}) - f(\mathbf{x}_0) = (\mathbf{A}\mathbf{x}_0, \mathbf{x} - \mathbf{x}_0) + (\mathbf{x}_0, \mathbf{A}(\mathbf{x} - \mathbf{x}_0)) + (\mathbf{x} - \mathbf{x}_0, \mathbf{A}(\mathbf{x} - \mathbf{x}_0)) = ((\mathbf{A} + \mathbf{A}^T)\mathbf{x}_0, \mathbf{x} - \mathbf{x}_0) + 0.5(\mathbf{x} - \mathbf{x}_0, (\mathbf{A} + \mathbf{A}^T)(\mathbf{x} - \mathbf{x}_0))$, we identify $\nabla f(\mathbf{x}_0) = (\mathbf{A} + \mathbf{A}^T)\mathbf{x}_0$, and $\nabla^2 f(\mathbf{x}_0) = \mathbf{A} + \mathbf{A}^T$. If the matrix \mathbf{A} is symmetric, then these expressions reduce to $\nabla f(\mathbf{x}_0) = 2\mathbf{A}\mathbf{x}_0$, and $\nabla^2 f(\mathbf{x}_0) = 2\mathbf{A}$.

Convex analysis

III

3.1 Convexity of sets

Definition 3.1 (convex set) Let $S \subseteq \mathbb{R}^n$. The set S is convex if

$$\left. \begin{array}{l} \mathbf{x}^1, \mathbf{x}^2 \in S \\ \lambda \in (0, 1) \end{array} \right\} \implies \lambda \mathbf{x}^1 + (1 - \lambda) \mathbf{x}^2 \in S$$

holds. ■

A set S is convex if, from everywhere in S , all other points of S are “visible.”

Figure 3.1 illustrates a convex set.

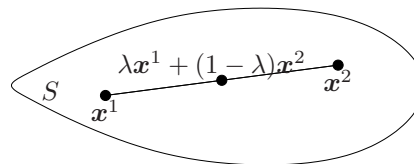


Figure 3.1: A convex set. (For the intermediate vector shown, the value of λ is $\approx 1/2$.)

Two non-convex sets are shown in Figure 3.2.

Example 3.2 (convex and non-convex sets) By using the definition of a convex set, the following can be established:

- (a) The set \mathbb{R}^n is a convex set.
- (b) The empty set is a convex set.

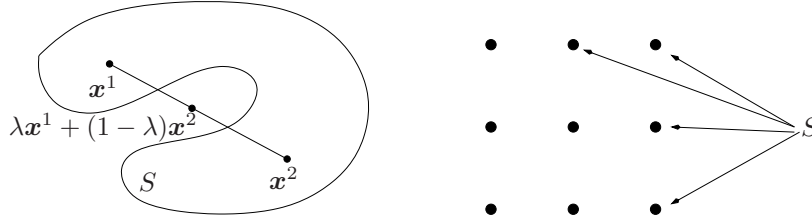


Figure 3.2: Two non-convex sets.

- (c) The set $\{\mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x}\| \leq a\}$ is convex for every value of $a \in \mathbb{R}$.
 [Note: $\|\cdot\|$ here denotes any vector norm, but we will almost always use the 2-norm,

$$\|\mathbf{x}\|_2 := \sqrt{\sum_{j=1}^n x_j^2}.$$

We will not write the index ₂, but instead use the 2-norm implicitly whenever writing $\|\cdot\|$.

- (d) The set $\{\mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x}\| = a\}$ is non-convex for every $a > 0$.
 (e) The set $\{0, 1, 2\}$ is non-convex. (The second illustration in Figure 3.2 is such a case of a set of integral points in \mathbb{R}^2 .) ■

Proposition 3.3 (convex intersection) *Suppose that S_k , $k \in \mathcal{K}$, is any collection of convex sets. Then, the intersection*

$$S := \bigcap_{k \in \mathcal{K}} S_k$$

is a convex set.

Proof. Let both \mathbf{x}^1 and \mathbf{x}^2 belong to S . (If two such points cannot be found, then the result holds vacuously.) Then, $\mathbf{x}^1 \in S_k$ and $\mathbf{x}^2 \in S_k$ for all $k \in \mathcal{K}$. Take $\lambda \in (0, 1)$. Then, $\lambda \mathbf{x}^1 + (1 - \lambda) \mathbf{x}^2 \in S_k, k \in \mathcal{K}$, by the convexity of the sets S_k . So, $\lambda \mathbf{x}^1 + (1 - \lambda) \mathbf{x}^2 \in \bigcap_{k \in \mathcal{K}} S_k = S$. ■

3.2 Polyhedral theory

3.2.1 Convex hulls

Consider the set $V := \{\mathbf{v}^1, \mathbf{v}^2\}$, where $\mathbf{v}^1, \mathbf{v}^2 \in \mathbb{R}^n$ and $\mathbf{v}^1 \neq \mathbf{v}^2$. A set naturally related to V is the line in \mathbb{R}^n through \mathbf{v}^1 and \mathbf{v}^2 [see Figure

3.3(b)], that is, $\{\lambda \mathbf{v}^1 + (1 - \lambda) \mathbf{v}^2 \mid \lambda \in \mathbb{R}\} = \{\lambda_1 \mathbf{v}^1 + \lambda_2 \mathbf{v}^2 \mid \lambda_1, \lambda_2 \in \mathbb{R}; \lambda_1 + \lambda_2 = 1\}$. Another set naturally related to V is the line segment between \mathbf{v}^1 and \mathbf{v}^2 [see Figure 3.3(c)], that is, $\{\lambda \mathbf{v}^1 + (1 - \lambda) \mathbf{v}^2 \mid \lambda \in [0, 1]\} = \{\lambda_1 \mathbf{v}^1 + \lambda_2 \mathbf{v}^2 \mid \lambda_1, \lambda_2 \geq 0; \lambda_1 + \lambda_2 = 1\}$. Motivated by these examples we define the *affine hull* and the *convex hull* of a set in \mathbb{R}^n .

Definition 3.4 (affine hull) Let $V := \{\mathbf{v}^1, \dots, \mathbf{v}^k\} \subset \mathbb{R}^n$. The affine hull of V is the set

$$\text{aff } V := \left\{ \lambda_1 \mathbf{v}^1 + \dots + \lambda_k \mathbf{v}^k \mid \lambda_1, \dots, \lambda_k \in \mathbb{R}; \sum_{i=1}^k \lambda_i = 1 \right\}.$$

The affine hull of an arbitrary set $V \subseteq \mathbb{R}^n$ is the smallest affine subspace that includes V .

A point $\lambda_1 \mathbf{v}^1 + \dots + \lambda_k \mathbf{v}^k$, where $\mathbf{v}^1, \dots, \mathbf{v}^k \in V$ and $\lambda_1, \dots, \lambda_k \in \mathbb{R}$ such that $\sum_{i=1}^k \lambda_i = 1$, is called an *affine combination* of the points $\mathbf{v}^1, \dots, \mathbf{v}^k$ (the number k of points in the sum must be finite). ■

Definition 3.5 (convex hull) Let $V := \{\mathbf{v}^1, \dots, \mathbf{v}^k\} \subset \mathbb{R}^n$. The convex hull of V is the set

$$\text{conv } V := \left\{ \lambda_1 \mathbf{v}^1 + \dots + \lambda_k \mathbf{v}^k \mid \lambda_1, \dots, \lambda_k \geq 0; \sum_{i=1}^k \lambda_i = 1 \right\}.$$

The convex hull of an arbitrary set $V \subseteq \mathbb{R}^n$ is the smallest convex set that includes V .

A point $\lambda_1 \mathbf{v}^1 + \dots + \lambda_k \mathbf{v}^k$, where $\mathbf{v}^1, \dots, \mathbf{v}^k \in V$ and $\lambda_1, \dots, \lambda_k \geq 0$ such that $\sum_{i=1}^k \lambda_i = 1$, is called a *convex combination* of the points $\mathbf{v}^1, \dots, \mathbf{v}^k$ (the number k of points in the sum must be finite). ■

Example 3.6 (affine hull, convex hull) (a) The affine hull of three or more points in \mathbb{R}^2 not all lying on the same line is \mathbb{R}^2 itself. The convex hull of five points in \mathbb{R}^2 is shown in Figure 3.4 (observe that the “corners” of the convex hull of the points are some of the points themselves).

(b) The affine hull of three points not all lying on the same line in \mathbb{R}^3 is the plane through the points.

(c) The affine hull of an affine space is the space itself and the convex hull of a convex set is the set itself. ■

From the definition of convex hull of a finite set it follows that the convex hull equals the set of all convex combinations of points in the set. It turns out that this also holds for arbitrary sets.

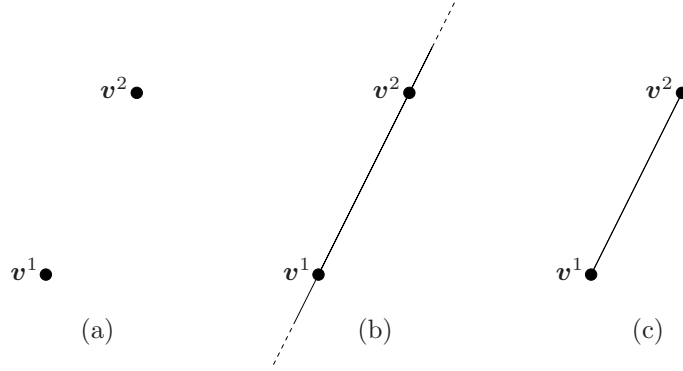


Figure 3.3: (a) The set V . (b) The set $\text{aff } V$. (c) The set $\text{conv } V$.

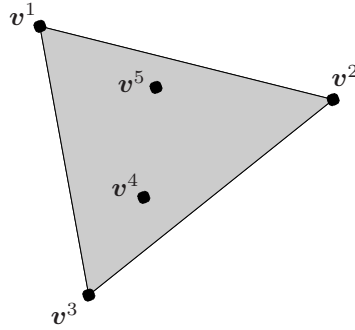


Figure 3.4: The convex hull of five points in \mathbb{R}^2 .

Proposition 3.7 *Let $V \subseteq \mathbb{R}^n$. Then, $\text{conv } V$ is the set of all convex combinations of points of V .*

Proof. Let Q be the set of all convex combinations of points of V . The inclusion $Q \subseteq \text{conv } V$ follows from the definition of a convex set (since $\text{conv } V$ is a convex set). We next show that Q is a convex set. If $\mathbf{x}^1, \mathbf{x}^2 \in Q$, then $\mathbf{x}^1 = \alpha_1 \mathbf{a}^1 + \cdots + \alpha_k \mathbf{a}^k$ and $\mathbf{x}^2 = \beta_1 \mathbf{b}^1 + \cdots + \beta_m \mathbf{b}^m$ for some $\mathbf{a}^1, \dots, \mathbf{a}^k, \mathbf{b}^1, \dots, \mathbf{b}^m \in V$ and $\alpha_1, \dots, \alpha_k, \beta_1, \dots, \beta_m \geq 0$ such that $\sum_{i=1}^k \alpha_i = \sum_{i=1}^m \beta_i = 1$. Let $\lambda \in (0, 1)$. Then

$$\begin{aligned} \lambda \mathbf{x}^1 + (1 - \lambda) \mathbf{x}^2 &= \lambda \alpha_1 \mathbf{a}^1 + \cdots + \lambda \alpha_k \mathbf{a}^k \\ &\quad + (1 - \lambda) \beta_1 \mathbf{b}^1 + \cdots + (1 - \lambda) \beta_m \mathbf{b}^m, \end{aligned}$$

and since $\lambda\alpha_1 + \cdots + \lambda\alpha_k + (1-\lambda)\beta_1 + \cdots + (1-\lambda)\beta_m = 1$, we have that $\lambda\mathbf{x}^1 + (1-\lambda)\mathbf{x}^2 \in Q$, so Q is convex. Since Q is convex and $V \subseteq Q$ it follows that $\text{conv } V \subseteq Q$ (from the definition of convex hull of an arbitrary set in \mathbb{R}^n it follows that $\text{conv } V$ is the smallest convex set that contains V). Therefore $Q = \text{conv } V$. ■

Proposition 3.7 shows that every point of the convex hull of a set can be written as a convex combination of points from the set. It tells, however, nothing about how many points that are required. This is the content of Carathéodory's Theorem.

Theorem 3.8 (Carathéodory's Theorem) *Let $\mathbf{x} \in \text{conv } V$, where $V \subseteq \mathbb{R}^n$. Then, \mathbf{x} can be expressed as a convex combination of $n+1$ or fewer points of V .*

Proof. From Proposition 3.7 it follows that $\mathbf{x} = \lambda_1\mathbf{a}^1 + \cdots + \lambda_m\mathbf{a}^m$ for some $\mathbf{a}^1, \dots, \mathbf{a}^m \in V$ and $\lambda_1, \dots, \lambda_m \geq 0$ such that $\sum_{i=1}^m \lambda_i = 1$. We assume that this representation of \mathbf{x} is chosen so that \mathbf{x} cannot be expressed as a convex combination of fewer than m points of V . It follows that no two of the points $\mathbf{a}^1, \dots, \mathbf{a}^m$ are equal and that $\lambda_1, \dots, \lambda_m > 0$. We prove the theorem by showing that $m \leq n+1$. Assume that $m > n+1$. Then the set $\{\mathbf{a}^1, \dots, \mathbf{a}^m\}$ must be affinely dependent, so there exist $\alpha_1, \dots, \alpha_m \in \mathbb{R}$, not all zero, such that $\sum_{i=1}^m \alpha_i \mathbf{a}^i = \mathbf{0}^n$ and $\sum_{i=1}^m \alpha_i = 0$. Let $\varepsilon > 0$ be such that $\lambda_1 + \varepsilon\alpha_1, \dots, \lambda_m + \varepsilon\alpha_m$ are non-negative with at least one of them zero (such an ε exists since the λ 's are all positive and at least one of the α 's must be negative). Then, $\mathbf{x} = \sum_{i=1}^m (\lambda_i + \varepsilon\alpha_i)\mathbf{a}^i$, and if terms with zero coefficients are omitted this is a representation of \mathbf{x} with fewer than m points; this is a contradiction. ■

3.2.2 Polytopes

We are now ready to define the geometrical object *polytope*.

Definition 3.9 (polytope) *A subset P of \mathbb{R}^n is a polytope if it is the convex hull of finitely many points in \mathbb{R}^n .* ■

Example 3.10 (polytopes) (a) The set shown in Figure 3.4 is a polytope.

(b) A cube and a tetrahedron are polytopes in \mathbb{R}^3 . ■

We next show how to characterize a polytope as the convex hull of its *extreme points*.

Definition 3.11 (extreme point) A point \mathbf{v} of a convex set P is called an extreme point if whenever $\mathbf{v} = \lambda \mathbf{x}^1 + (1 - \lambda) \mathbf{x}^2$, where $\mathbf{x}^1, \mathbf{x}^2 \in P$ and $\lambda \in (0, 1)$, then $\mathbf{v} = \mathbf{x}^1 = \mathbf{x}^2$. ■

Example 3.12 (extreme points) The set shown in Figure 3.3(c) has the extreme points \mathbf{v}^1 and \mathbf{v}^2 . The set shown in Figure 3.4 has the extreme points \mathbf{v}^1 , \mathbf{v}^2 , and \mathbf{v}^3 . The set shown in Figure 3.3(b) does not have any extreme points. ■

Lemma 3.13 Let $V := \{\mathbf{v}^1, \dots, \mathbf{v}^k\} \subset \mathbb{R}^n$ and let P be the polytope $\text{conv } V$. Then, each extreme point of P lies in V .

Proof. Assume that $\mathbf{w} \notin V$ is an extreme point of P . We have that $\mathbf{w} = \sum_{i=1}^k \lambda_i \mathbf{v}^i$, for some $\lambda_i \geq 0$ such that $\sum_{i=1}^k \lambda_i = 1$. At least one of the λ_i 's must be nonzero, say λ_1 . If $\lambda_1 = 1$ then $\mathbf{w} = \mathbf{v}^1$, a contradiction, so $\lambda_1 \in (0, 1)$. We have that

$$\mathbf{w} = \lambda_1 \mathbf{v}^1 + (1 - \lambda_1) \sum_{i=2}^k \frac{\lambda_i}{1 - \lambda_1} \mathbf{v}^i.$$

Since $\sum_{i=2}^k \lambda_i / (1 - \lambda_1) = 1$ we have that $\sum_{i=2}^k \lambda_i / (1 - \lambda_1) \mathbf{v}^i \in P$, but \mathbf{w} is an extreme point of P so $\mathbf{w} = \mathbf{v}^1$, a contradiction. ■

Proposition 3.14 Let $V := \{\mathbf{v}^1, \dots, \mathbf{v}^k\} \subset \mathbb{R}^n$ and let P be the polytope $\text{conv } V$. Then P is equal to the convex hull of its extreme points.

Proof. Let Q be the set of extreme points of P . If $\mathbf{v}^i \in Q$ for all $i = 1, \dots, k$ we are done, so assume that $\mathbf{v}^1 \notin Q$. Then $\mathbf{v}^1 = \lambda \mathbf{u} + (1 - \lambda) \mathbf{w}$ for some $\lambda \in (0, 1)$ and $\mathbf{u}, \mathbf{w} \in P$, $\mathbf{u} \neq \mathbf{w}$. Further, $\mathbf{u} = \sum_{i=1}^k \alpha_i \mathbf{v}^i$ and $\mathbf{w} = \sum_{i=1}^k \beta_i \mathbf{v}^i$, for some $\alpha_1, \dots, \alpha_k, \beta_1, \dots, \beta_k \geq 0$ such that $\sum_{i=1}^k \alpha_i = \sum_{i=1}^k \beta_i = 1$. Hence,

$$\mathbf{v}^1 = \lambda \sum_{i=1}^k \alpha_i \mathbf{v}^i + (1 - \lambda) \sum_{i=1}^k \beta_i \mathbf{v}^i = \sum_{i=1}^k (\lambda \alpha_i + (1 - \lambda) \beta_i) \mathbf{v}^i.$$

It must hold that $\alpha_1, \beta_1 \neq 1$, since otherwise $\mathbf{u} = \mathbf{w} = \mathbf{v}^1$, a contradiction. Therefore,

$$\mathbf{v}^1 = \sum_{i=2}^k \frac{\lambda \alpha_i + (1 - \lambda) \beta_i}{1 - (\lambda \alpha_1 + (1 - \lambda) \beta_1)} \mathbf{v}^i,$$

and since $\sum_{i=2}^k (\lambda\alpha_i + (1-\lambda)\beta_i) / (1-\lambda\alpha_1 - (1-\lambda)\beta_1) = 1$ it follows that $\text{conv } V = \text{conv } (V \setminus \{\mathbf{v}^1\})$. Similarly, every $\mathbf{v}^i \notin Q$ can be removed, and we end up with a set $T \subseteq V$ such that $\text{conv } T = \text{conv } V$ and $T \subseteq Q$. On the other hand, from Lemma 3.13 we have that every extreme point of the set $\text{conv } T$ lies in T and since $\text{conv } T = \text{conv } V$ it follows that Q is the set of extreme points of $\text{conv } T$, so $Q \subseteq T$. Hence, $T = Q$ and we are done. ■

3.2.3 Polyhedra

Closely related to the polytope is the polyhedron. We will show that every polyhedron is the sum of a polytope and a polyhedral cone. In the next subsection we show that a set is a polytope if and only if it is a bounded polyhedron.

Definition 3.15 (polyhedron) *A subset P of \mathbb{R}^n is a polyhedron if there exist a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ and a vector $\mathbf{b} \in \mathbb{R}^m$ such that*

$$P = \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{Ax} \leq \mathbf{b}\}$$

holds. ■

The importance of polyhedra is obvious, since the set of feasible solutions of every linear programming problem is a polyhedron.

Example 3.16 (polyhedra) (a) Figure 3.5 shows the bounded polyhedron $P := \{\mathbf{x} \in \mathbb{R}^2 \mid x_1 \geq 2; x_1 + x_2 \leq 6; 2x_1 - x_2 \leq 4\}$.

(b) The unbounded polyhedron $P := \{\mathbf{x} \in \mathbb{R}^2 \mid x_1 + x_2 \geq 2; x_1 - x_2 \leq 2; 3x_1 - x_2 \geq 0\}$ is shown in Figure 3.6. ■

Often it is hard to decide whether a point in a convex set is an extreme point or not. This is not the case for the polyhedron since there is an algebraic characterization of the extreme points of such a set. Given an $\tilde{\mathbf{x}} \in \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{Ax} \leq \mathbf{b}\}$ we refer to the rows of $\mathbf{A}\tilde{\mathbf{x}} \leq \mathbf{b}$ that are fulfilled with equality as the *equality subsystem* of $\mathbf{A}\tilde{\mathbf{x}} \leq \mathbf{b}$, and denote it by $\tilde{\mathbf{A}}\tilde{\mathbf{x}} = \tilde{\mathbf{b}}$, that is, $(\tilde{\mathbf{A}}, \tilde{\mathbf{b}})$ consists of the rows $i \in \{1, \dots, m\}$ of (\mathbf{A}, \mathbf{b}) such that $\mathbf{A}_i \tilde{\mathbf{x}} = b_i$, where \mathbf{A}_i is the i^{th} row of \mathbf{A} . The number of rows in $(\tilde{\mathbf{A}}, \tilde{\mathbf{b}})$ is denoted by \tilde{m} .

Theorem 3.17 (algebraic characterization of extreme points) *Let $\tilde{\mathbf{x}} \in P = \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{Ax} \leq \mathbf{b}\}$, where $\mathbf{A} \in \mathbb{R}^{m \times n}$ has $\text{rank } \mathbf{A} = n$ and $\mathbf{b} \in \mathbb{R}^m$. Further, let $\tilde{\mathbf{A}}\tilde{\mathbf{x}} = \tilde{\mathbf{b}}$ be the equality subsystem of $\mathbf{A}\tilde{\mathbf{x}} \leq \mathbf{b}$. Then $\tilde{\mathbf{x}}$ is an extreme point of P if and only if $\text{rank } \tilde{\mathbf{A}} = n$.*

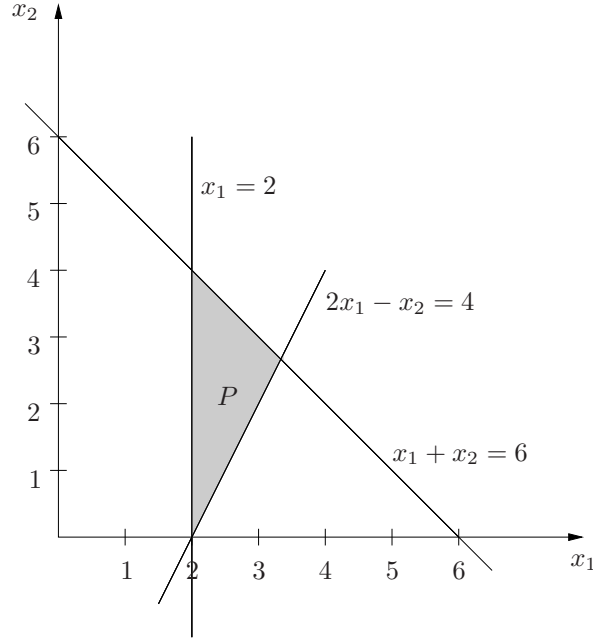


Figure 3.5: Illustration of the bounded polyhedron $P := \{ \mathbf{x} \in \mathbb{R}^2 \mid x_1 \geq 2; x_1 + x_2 \leq 6; 2x_1 - x_2 \leq 4 \}$.

Proof. $[\Rightarrow]$ Suppose that $\tilde{\mathbf{x}}$ is an extreme point of P . If $\mathbf{A}\tilde{\mathbf{x}} < \mathbf{b}$ then $\tilde{\mathbf{x}} + \varepsilon \mathbf{1}^n, \tilde{\mathbf{x}} - \varepsilon \mathbf{1}^n \in P$ if $\varepsilon > 0$ is sufficiently small. But $\tilde{\mathbf{x}} = 1/2(\tilde{\mathbf{x}} + \varepsilon \mathbf{1}^n) + 1/2(\tilde{\mathbf{x}} - \varepsilon \mathbf{1}^n)$ which contradicts that $\tilde{\mathbf{x}}$ is an extreme point, so assume that at least one of the rows in $\mathbf{A}\tilde{\mathbf{x}} \leq \mathbf{b}$ is fulfilled with equality. If $\tilde{\mathbf{A}}\tilde{\mathbf{x}} = \tilde{\mathbf{b}}$ is the equality subsystem of $\mathbf{A}\tilde{\mathbf{x}} \leq \mathbf{b}$ and $\text{rank } \tilde{\mathbf{A}} \leq n - 1$, then there exists a $\mathbf{w} \neq \mathbf{0}^n$ such that $\tilde{\mathbf{A}}\mathbf{w} = \mathbf{0}^{\tilde{m}}$, so $\tilde{\mathbf{x}} + \varepsilon \mathbf{w}, \tilde{\mathbf{x}} - \varepsilon \mathbf{w} \in P$ if $\varepsilon > 0$ is sufficiently small. But $\tilde{\mathbf{x}} = 1/2(\tilde{\mathbf{x}} + \varepsilon \mathbf{w}) + 1/2(\tilde{\mathbf{x}} - \varepsilon \mathbf{w})$, which contradicts that $\tilde{\mathbf{x}}$ is an extreme point. Hence, $\text{rank } \tilde{\mathbf{A}} = n$.

$[\Leftarrow]$ Assume that $\text{rank } \tilde{\mathbf{A}} = n$. Then, $\tilde{\mathbf{x}}$ is the unique solution to $\tilde{\mathbf{A}}\mathbf{x} = \tilde{\mathbf{b}}$. If $\tilde{\mathbf{x}}$ is not an extreme point of P it follows that $\tilde{\mathbf{x}} = \lambda \mathbf{u} + (1 - \lambda)\mathbf{v}$ for some $\lambda \in (0, 1)$ and $\mathbf{u}, \mathbf{v} \in P, \mathbf{u} \neq \mathbf{v}$. This yields that $\lambda \tilde{\mathbf{A}}\mathbf{u} + (1 - \lambda)\tilde{\mathbf{A}}\mathbf{v} = \tilde{\mathbf{b}}$, and since $\mathbf{A}\mathbf{u} \leq \mathbf{b}$ and $\mathbf{A}\mathbf{v} \leq \mathbf{b}$ it follows that $\tilde{\mathbf{A}}\mathbf{u} = \tilde{\mathbf{A}}\mathbf{v} = \tilde{\mathbf{b}}$, which contradicts that $\tilde{\mathbf{x}}$ is the unique solution to $\tilde{\mathbf{A}}\mathbf{x} = \tilde{\mathbf{b}}$. Therefore $\tilde{\mathbf{x}}$ must be an extreme point. ■

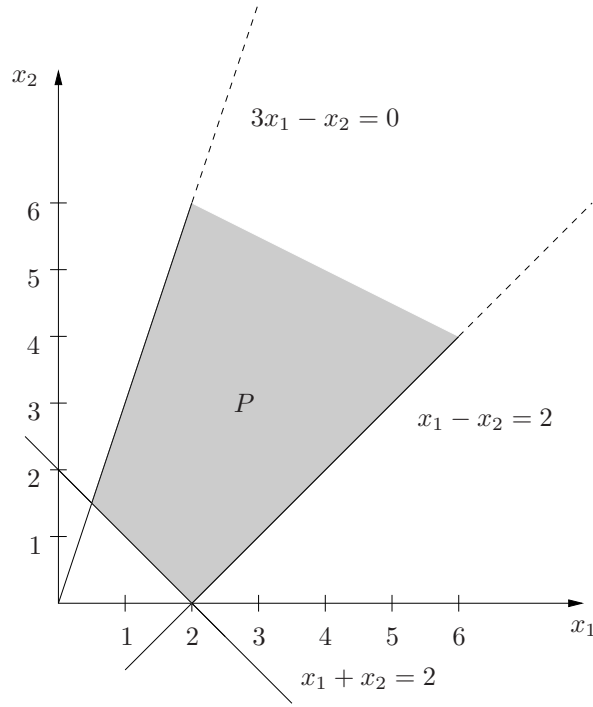


Figure 3.6: Illustration of the unbounded polyhedron $P := \{ \mathbf{x} \in \mathbb{R}^2 \mid x_1 + x_2 \geq 2; x_1 - x_2 \leq 2; 3x_1 - x_2 \geq 0 \}$.

Corollary 3.18 Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$. The number of extreme points of the polyhedron $P := \{ \mathbf{x} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{x} \leq \mathbf{b} \}$ is finite.

Proof. The theorem implies that the number of extreme points of P never exceeds the number of ways in which n objects can be chosen from a set of m objects, that is, the number of extreme points is less than or equal to

$$\binom{m}{n} = \frac{m!}{n!(m-n)!}.$$

We are done. ■

Remark 3.19 Since the number of extreme points is finite, the convex hull of the extreme points of a polyhedron is a polytope. ■

Definition 3.20 (cone) A subset C of \mathbb{R}^n is a cone if $\lambda \mathbf{x} \in C$ whenever $\mathbf{x} \in C$ and $\lambda > 0$. ■

Example 3.21 (cone) (a) The set $\{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{x} \leq \mathbf{0}^m\}$, where $\mathbf{A} \in \mathbb{R}^{m \times n}$, is a cone. Since this set is a polyhedron, this type of cone is usually called a *polyhedral cone*.

(b) Figure 3.7(a) illustrates a convex cone and Figure 3.7(b) illustrates a non-convex cone in \mathbb{R}^2 . ■

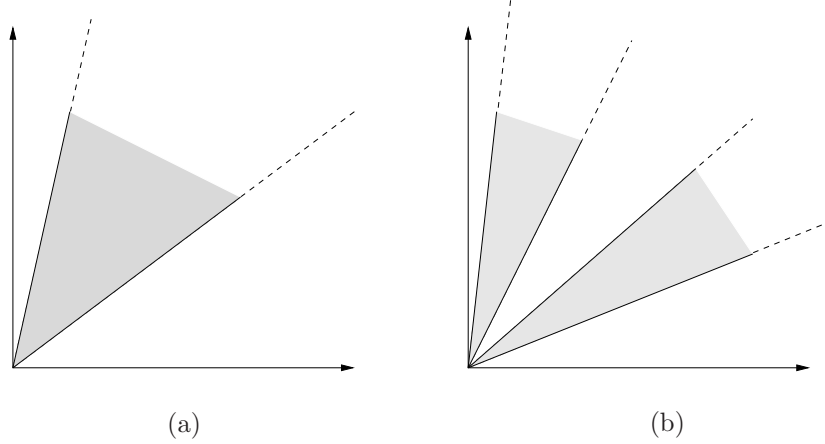


Figure 3.7: (a) A convex cone in \mathbb{R}^2 . (b) A non-convex cone in \mathbb{R}^2 .

We have arrived at the most important theorem of this section, namely the Representation Theorem, which tells that every polyhedron is the sum of a polytope and a polyhedral cone. The Representation Theorem will have great importance in the linear programming theory in Chapter 8.

Theorem 3.22 (Representation Theorem) Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$. Let $Q := \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{x} \leq \mathbf{b}\}$, P denote the convex hull of the extreme points of Q , and $C := \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{x} \leq \mathbf{0}^m\}$. If $\text{rank } \mathbf{A} = n$ then $Q = P + C := \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{x} = \mathbf{u} + \mathbf{v} \text{ for some } \mathbf{u} \in P \text{ and } \mathbf{v} \in C\}$. In other words, every polyhedron (that has at least one extreme point) is the sum of a polytope and a polyhedral cone.

Proof. Let $\tilde{\mathbf{x}} \in Q$ and $\tilde{\mathbf{A}}\tilde{\mathbf{x}} = \tilde{\mathbf{b}}$ be the corresponding equality subsystem of $\mathbf{A}\mathbf{x} \leq \mathbf{b}$. We prove the theorem by induction on the rank of $\tilde{\mathbf{A}}$.

If $\text{rank } \tilde{\mathbf{A}} = n$ it follows from Theorem 3.17 that $\tilde{\mathbf{x}}$ is an extreme point of Q , so $\tilde{\mathbf{x}} \in P + C$, since $\mathbf{0}^n \in C$. Now assume that $\tilde{\mathbf{x}} \in P + C$ for

all $\tilde{\mathbf{x}} \in Q$ with $k \leq \text{rank } \tilde{\mathbf{A}} \leq n$, and choose $\tilde{\mathbf{x}} \in Q$ with $\text{rank } \tilde{\mathbf{A}} = k - 1$. Then there exists a $\mathbf{w} \neq \mathbf{0}^n$ such that $\tilde{\mathbf{A}}\mathbf{w} = \mathbf{0}^m$. If $|\lambda|$ is sufficiently small it follows that $\tilde{\mathbf{x}} + \lambda\mathbf{w} \in Q$. (Why?) If $\tilde{\mathbf{x}} + \lambda\mathbf{w} \in Q$ for all $\lambda \in \mathbb{R}$ we must have $\mathbf{A}\mathbf{w} = \mathbf{0}^n$ which implies $\text{rank } \mathbf{A} \leq n - 1$, a contradiction. Suppose that there exists a largest λ^+ such that $\tilde{\mathbf{x}} + \lambda^+\mathbf{w} \in Q$. Then if $\tilde{\mathbf{A}}(\tilde{\mathbf{x}} + \lambda^+\mathbf{w}) = \tilde{\mathbf{b}}$ is the equality subsystem of $\mathbf{A}(\tilde{\mathbf{x}} + \lambda^+\mathbf{w}) \leq \mathbf{b}$ we must have $\text{rank } \tilde{\mathbf{A}} \geq k$. (Why?) By the induction hypothesis it then follows that $\tilde{\mathbf{x}} + \lambda^+\mathbf{w} \in P + C$. On the other hand, if $\tilde{\mathbf{x}} + \lambda\mathbf{w} \in Q$ for all $\lambda \geq 0$ then $\mathbf{A}\mathbf{w} \leq \mathbf{0}^m$, so $\mathbf{w} \in C$. Similarly, if $\tilde{\mathbf{x}} + \lambda(-\mathbf{w}) \in Q$ for all $\lambda \geq 0$ then $-\mathbf{w} \in C$, and if there exists a largest λ^- such that $\tilde{\mathbf{x}} + \lambda^-(-\mathbf{w}) \in Q$ then $\tilde{\mathbf{x}} + \lambda^-(-\mathbf{w}) \in P + C$.

Above we got a contradiction if none of λ^+ or λ^- existed. If only one of them exists, say λ^+ , then $\tilde{\mathbf{x}} + \lambda^+\mathbf{w} \in P + C$ and $-\mathbf{w} \in C$, and it follows that $\tilde{\mathbf{x}} \in P + C$. Otherwise, if both λ^+ and λ^- exist then $\tilde{\mathbf{x}} + \lambda^+\mathbf{w} \in P + C$ and $\tilde{\mathbf{x}} + \lambda^-(-\mathbf{w}) \in P + C$, and $\tilde{\mathbf{x}}$ can be written as a convex combination of these points, which gives $\tilde{\mathbf{x}} \in P + C$. We have shown that $\tilde{\mathbf{x}} \in P + C$ for all $\tilde{\mathbf{x}} \in Q$ with $k - 1 \leq \text{rank } \tilde{\mathbf{A}} \leq n$ and the theorem follows by induction. \blacksquare

Example 3.23 (illustration of the Representation Theorem) Figure 3.8(a) shows a bounded polyhedron. The interior point $\tilde{\mathbf{x}}$ can be written as a convex combination of the extreme point \mathbf{x}^5 and the point \mathbf{v} on the boundary, that is, there is a $\lambda \in (0, 1)$ such that

$$\tilde{\mathbf{x}} = \lambda\mathbf{x}^5 + (1 - \lambda)\mathbf{v}.$$

Further, the point \mathbf{v} can be written as a convex combination of the extreme points \mathbf{x}^2 and \mathbf{x}^3 , that is, there exists a $\mu \in (0, 1)$ such that

$$\mathbf{v} = \mu\mathbf{x}^2 + (1 - \mu)\mathbf{x}^3.$$

This gives that

$$\tilde{\mathbf{x}} = \lambda\mathbf{x}^5 + (1 - \lambda)\mu\mathbf{x}^2 + (1 - \lambda)(1 - \mu)\mathbf{x}^3,$$

and since $\lambda, (1 - \lambda)\mu, (1 - \lambda)(1 - \mu) \geq 0$ and

$$\lambda + (1 - \lambda)\mu + (1 - \lambda)(1 - \mu) = 1$$

holds we have that $\tilde{\mathbf{x}}$ lies in the convex hull of the extreme points \mathbf{x}^2 , \mathbf{x}^3 , and \mathbf{x}^5 .

Figure 3.8(b) shows an unbounded polyhedron. The interior point $\tilde{\mathbf{x}}$ can be written as a convex combination of the extreme point \mathbf{x}^3 and the point \mathbf{v} on the boundary, that is, there exists a $\lambda \in (0, 1)$ such that

$$\tilde{\mathbf{x}} = \lambda\mathbf{x}^3 + (1 - \lambda)\mathbf{v}.$$

The point \mathbf{v} lies on the halfline $\{\mathbf{x} \in \mathbb{R}^2 \mid \mathbf{x} = \mathbf{x}^2 + \mu(\mathbf{x}^1 - \mathbf{x}^2), \mu \geq 0\}$. All the points on this halfline are feasible, which gives that if the polyhedron is given by $\{\mathbf{x} \in \mathbb{R}^2 \mid \mathbf{A}\mathbf{x} \leq \mathbf{b}\}$ then

$$\mathbf{A}(\mathbf{x}^2 + \mu(\mathbf{x}^1 - \mathbf{x}^2)) = \mathbf{A}\mathbf{x}^2 + \mu\mathbf{A}(\mathbf{x}^1 - \mathbf{x}^2) \leq \mathbf{b}, \quad \mu \geq 0.$$

But then we must have that $\mathbf{A}(\mathbf{x}^1 - \mathbf{x}^2) \leq \mathbf{0}^2$ since otherwise some component of $\mu\mathbf{A}(\mathbf{x}^1 - \mathbf{x}^2)$ tends to infinity as μ tends to infinity. Therefore $\mathbf{x}^1 - \mathbf{x}^2$ lies in the cone $C := \{\mathbf{x} \in \mathbb{R}^2 \mid \mathbf{A}\mathbf{x} \leq \mathbf{0}^2\}$. Now there exists a $\mu \geq 0$ such that

$$\mathbf{v} = \mathbf{x}^2 + \mu(\mathbf{x}^1 - \mathbf{x}^2),$$

and it follows that

$$\tilde{\mathbf{x}} = \lambda\mathbf{x}^3 + (1 - \lambda)\mathbf{x}^2 + (1 - \lambda)\mu(\mathbf{x}^1 - \mathbf{x}^2),$$

so since $(1 - \lambda)\mu \geq 0$ and $\mathbf{x}^1 - \mathbf{x}^2 \in C$, $\tilde{\mathbf{x}}$ is the sum of a point in the convex hull of the extreme points and a point in the polyhedral cone C .

Note that the representation of a vector $\tilde{\mathbf{x}}$ in a polyhedron is normally *not* uniquely determined; in the case of Figure 3.8(a), for example, we can also represent $\tilde{\mathbf{x}}$ as a convex combination of \mathbf{x}^1 , \mathbf{x}^4 , and \mathbf{x}^5 . ■

3.2.4 The Separation Theorem and Farkas' Lemma

We introduce the important concept of separation and use it to show that every polytope is a polyhedron.

Theorem 3.24 (Separation Theorem) *Suppose that the set $C \subseteq \mathbb{R}^n$ is closed and convex, and that the point \mathbf{y} does not lie in C . Then there exist a vector $\boldsymbol{\pi} \neq \mathbf{0}^n$ and $\alpha \in \mathbb{R}$ such that $\boldsymbol{\pi}^T \mathbf{y} > \alpha$ and $\boldsymbol{\pi}^T \mathbf{x} \leq \alpha$ for all $\mathbf{x} \in C$.* ■

We postpone the proof of this theorem since it requires the Weierstrass Theorem 4.7 and the first order necessary optimality condition given in Proposition 4.23(b). Instead the proof is presented in Section 4.4.

The Separation Theorem is easy to describe geometrically: If a point \mathbf{y} does not lie in a closed and convex set C , then there exists a hyperplane that separates \mathbf{y} from C .

Example 3.25 (illustration of the Separation Theorem) Consider the convex and closed set $C := \{\mathbf{x} \in \mathbb{R}^2 \mid \|\mathbf{x}\| \leq 1\}$ (i.e., C is the unit disc in

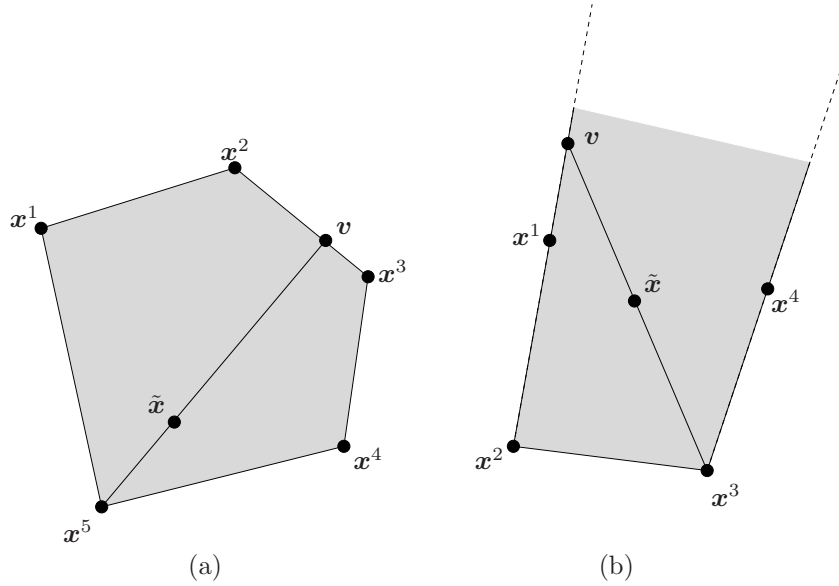


Figure 3.8: Illustration of the Representation Theorem (a) in the bounded case, and (b) in the unbounded case.

\mathbb{R}^2), and the point $\mathbf{y} := (1.5, 1.5)^T$. Since $\mathbf{y} \notin C$ the Separation Theorem implies that there exists a line in \mathbb{R}^2 that separates \mathbf{y} from C . (This line is however not unique: in Figure 3.9 we see that the line given by $\boldsymbol{\pi} = (1, 1)^T$ and $\alpha = 2$ is a separating line, while the one constructed in the proof of Theorem 3.24 is actually a tangent plane to C .) ■

Theorem 3.26 *A set P is a polytope if and only if it is a bounded polyhedron.*

Proof. [\Leftarrow] From the Representation Theorem 3.22 we get that a bounded polyhedron is the convex hull of its extreme points and hence by Remark 3.19 a polytope.

[\Rightarrow] Let $V := \{\mathbf{v}^1, \dots, \mathbf{v}^k\} \subset \mathbb{R}^n$ and let P be the polytope $\text{conv } V$. In order to prove that P is a polyhedron we must show that P is the solution set of some finite system of linear inequalities. The idea of the proof is to define a bounded polyhedron consisting of the coefficients and right-hand sides of all valid inequalities for P and then apply the Representation Theorem to select a finite subset of those valid inequalities.

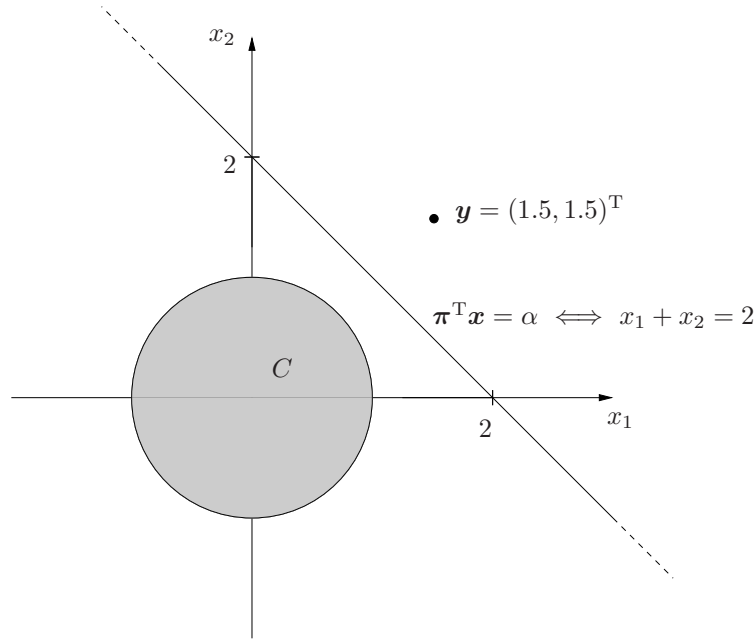


Figure 3.9: Illustration of the Separation Theorem: the unit disk is separated from \mathbf{y} by the line $\{\mathbf{x} \in \mathbb{R}^2 \mid x_1 + x_2 = 2\}$.

To carry this out, consider the set $Q \subset \mathbb{R}^{n+1}$ defined as

$$\left\{ \begin{pmatrix} \mathbf{a} \\ b \end{pmatrix} \mid \mathbf{a} \in \mathbb{R}^n; b \in \mathbb{R}; -\mathbf{1}^n \leq \mathbf{a} \leq \mathbf{1}^n; -1 \leq b \leq 1; \mathbf{a}^T \mathbf{v} \leq b, \mathbf{v} \in V \right\}.$$

Since V is a finite set, Q is a polyhedron. Further, Q is bounded, so by the Representation Theorem we know that Q is the convex hull of its extreme points given by

$$\begin{pmatrix} \mathbf{a}^1 \\ b_1 \end{pmatrix}, \dots, \begin{pmatrix} \mathbf{a}^m \\ b_m \end{pmatrix}.$$

We will prove that the linear system

$$(\mathbf{a}^1)^T \mathbf{x} \leq b_1, \dots, (\mathbf{a}^m)^T \mathbf{x} \leq b_m, \quad (3.1)$$

defines P . We first show that P is contained in the solution set of (3.1). Suppose that $\tilde{\mathbf{x}} \in P$. Then $\tilde{\mathbf{x}} = \lambda_1 \mathbf{v}^1 + \dots + \lambda_k \mathbf{v}^k$ for some $\lambda_1, \dots, \lambda_k \geq 0$

such that $\sum_{i=1}^k \lambda_i = 1$. Thus, for each $i = 1, \dots, m$, we have

$$\begin{aligned} (\mathbf{a}^i)^T \tilde{\mathbf{x}} &= (\mathbf{a}^i)^T (\lambda_1 \mathbf{v}^1 + \dots + \lambda_k \mathbf{v}^k) = \lambda_1 (\mathbf{a}^i)^T \mathbf{v}^1 + \dots + \lambda_k (\mathbf{a}^i)^T \mathbf{v}^k \\ &\leq \lambda_1 b_i + \dots + \lambda_k b_i = b_i, \end{aligned}$$

so $\tilde{\mathbf{x}}$ satisfies all inequalities in (3.1).

In order to show that the solution set of (3.1) is contained in P , let $\tilde{\mathbf{x}}$ be a solution to (3.1) and suppose that $\tilde{\mathbf{x}} \notin P$. Then, by the Separation Theorem 3.24 there exist a vector $\boldsymbol{\pi} \neq \mathbf{0}^n$ and $\alpha \in \mathbb{R}$ such that $\boldsymbol{\pi}^T \tilde{\mathbf{x}} > \alpha$ and $\boldsymbol{\pi}^T \mathbf{x} \leq \alpha$ for all $\mathbf{x} \in P$. By scaling $\boldsymbol{\pi}^T \mathbf{x} \leq \alpha$ by a positive constant if necessary, we may assume that $-\mathbf{1}^n \leq \boldsymbol{\pi} \leq \mathbf{1}^n$ and $-1 \leq \alpha \leq 1$. That is, we may assume that $\begin{pmatrix} \boldsymbol{\pi} \\ \alpha \end{pmatrix} \in Q$. So we may write

$$\begin{pmatrix} \boldsymbol{\pi} \\ \alpha \end{pmatrix} = \lambda_1 \begin{pmatrix} \mathbf{a}^1 \\ b_1 \end{pmatrix} + \dots + \lambda_m \begin{pmatrix} \mathbf{a}^m \\ b_m \end{pmatrix},$$

for some $\lambda_1, \dots, \lambda_m \geq 0$ such that $\sum_{i=1}^m \lambda_i = 1$. Therefore,

$$\boldsymbol{\pi}^T \tilde{\mathbf{x}} = \lambda_1 (\mathbf{a}^1)^T \tilde{\mathbf{x}} + \dots + \lambda_m (\mathbf{a}^m)^T \tilde{\mathbf{x}} \leq \lambda_1 b_1 + \dots + \lambda_m b_m = \alpha.$$

But this is a contradiction, since $\boldsymbol{\pi}^T \tilde{\mathbf{x}} > \alpha$. So $\tilde{\mathbf{x}} \in P$, which completes the proof. \blacksquare

We introduce the concept of finitely generated cones. In the proof of Farkas' Lemma below we will use that finitely generated cones are convex and closed, and in order to show this fact we prove that finitely generated cones are polyhedral sets.

Definition 3.27 (finitely generated cone) A finitely generated cone is one that is generated by a finite set, that is, a cone of the form

$$\text{cone}\{\mathbf{v}^1, \dots, \mathbf{v}^m\} := \{\lambda_1 \mathbf{v}^1 + \dots + \lambda_m \mathbf{v}^m \mid \lambda_1, \dots, \lambda_m \geq 0\},$$

where $\mathbf{v}^1, \dots, \mathbf{v}^m \in \mathbb{R}^n$. Note that if $\mathbf{A} \in \mathbb{R}^{m \times n}$, then the set $\{\mathbf{y} \in \mathbb{R}^m \mid \mathbf{y} = \mathbf{A}\mathbf{x}; \mathbf{x} \geq \mathbf{0}^n\}$ is a finitely generated cone. \blacksquare

Recall that a cone that is a polyhedron is called a *polyhedral cone*. We show that a finitely generated cone is always a polyhedral cone and vice versa.

Theorem 3.28 A convex cone in \mathbb{R}^n is finitely generated if and only if it is polyhedral.

Proof. \Rightarrow Assume that C is the finitely generated cone

$$\text{cone}\{\mathbf{v}^1, \dots, \mathbf{v}^m\},$$

where $\mathbf{v}^1, \dots, \mathbf{v}^m \in \mathbb{R}^n$. From Theorem 3.26 we know that polytopes are polyhedral sets, so $\text{conv}\{\mathbf{0}^n, \mathbf{v}^1, \dots, \mathbf{v}^m\}$ is the solution set of some linear inequalities

$$(\mathbf{a}^1)^T \mathbf{x} \leq b_1, \dots, (\mathbf{a}^k)^T \mathbf{x} \leq b_k. \quad (3.2)$$

Since the solution set of these inequalities contains $\mathbf{0}^n$ we must have $b_1, \dots, b_k \geq 0$. We show that C is the polyhedral cone A that equals the solution set of the inequalities of (3.2) for which $b_i = 0$. Since $\mathbf{v}^1, \dots, \mathbf{v}^m \in A$ we have $C \subseteq A$. In order to show that $A \subseteq C$, assume that $\mathbf{w} \in A$. Then $\lambda \mathbf{w}$ is in the solution set of (3.2) if $\lambda > 0$ is sufficiently small. Hence there exists a $\lambda > 0$ such that

$$\begin{aligned} \lambda \mathbf{w} &\in \{\mathbf{x} \in \mathbb{R}^n \mid (\mathbf{a}^1)^T \mathbf{x} \leq b_1, \dots, (\mathbf{a}^k)^T \mathbf{x} \leq b_k\} \\ &= \text{conv}\{\mathbf{0}^n, \mathbf{v}^1, \dots, \mathbf{v}^m\} \subseteq C, \end{aligned}$$

so $\mathbf{w} \in (1/\lambda)C = C$. Hence $A \subseteq C$, and $C = A$.

\Leftarrow Suppose that C is a polyhedral cone in \mathbb{R}^n . Let P be a polytope in \mathbb{R}^n such that $\mathbf{0}^n \in \text{int } P$ (that is, $\mathbf{0}^n$ lies in the interior of P). Then $C \cap P$ is a bounded polyhedron and hence the Representation Theorem gives that $C \cap P = \text{conv}\{\mathbf{v}^1, \dots, \mathbf{v}^m\}$, where $\mathbf{v}^1, \dots, \mathbf{v}^m$ is the set of extreme points of $C \cap P$. We show that C is the finitely generated cone $\text{cone}\{\mathbf{v}^1, \dots, \mathbf{v}^m\}$. Since $\mathbf{v}^1, \dots, \mathbf{v}^m \in C$ and C is a polyhedral cone we get that $\text{cone}\{\mathbf{v}^1, \dots, \mathbf{v}^m\} \subseteq C$. If $\mathbf{c} \in C$, then, since $\mathbf{0}^n \in \text{int } P$, there exists a $\lambda > 0$ such that $\lambda \mathbf{c} \in P$. Thus,

$$\lambda \mathbf{c} \in C \cap P = \text{conv}\{\mathbf{v}^1, \dots, \mathbf{v}^m\} \subseteq \text{cone}\{\mathbf{v}^1, \dots, \mathbf{v}^m\},$$

and so $\mathbf{c} \in (1/\lambda)\text{cone}\{\mathbf{v}^1, \dots, \mathbf{v}^m\} = \text{cone}\{\mathbf{v}^1, \dots, \mathbf{v}^m\}$. Hence it follows that $C \subseteq \text{cone}\{\mathbf{v}^1, \dots, \mathbf{v}^m\}$, and $C = \text{cone}\{\mathbf{v}^1, \dots, \mathbf{v}^m\}$. ■

Corollary 3.29 *Finitely generated cones in \mathbb{R}^n are convex and closed.*

Proof. Halfspaces, that is, sets of the form $\{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{a}^T \mathbf{x} \leq b\}$ for some vector $\mathbf{a} \in \mathbb{R}^n$ and $b \in \mathbb{R}$, are convex and closed. (Why?) By the theorem a finitely generated cone is the intersection of finitely many halfspaces and thus the corollary follows from Proposition 3.3 and the fact that intersections of closed sets are closed. ■

We close this section by proving the famous Farkas' Lemma by using the Separation Theorem 3.24 and the fact that finitely generated cones are convex and closed.

Theorem 3.30 (Farkas' Lemma) *Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$. Then, exactly one of the systems*

$$\begin{aligned} \mathbf{Ax} &= \mathbf{b}, \\ \mathbf{x} &\geq \mathbf{0}^n, \end{aligned} \tag{I}$$

and

$$\begin{aligned} \mathbf{A}^T \boldsymbol{\pi} &\leq \mathbf{0}^n, \\ \mathbf{b}^T \boldsymbol{\pi} &> 0, \end{aligned} \tag{II}$$

has a feasible solution, and the other system is inconsistent.

Proof. Let $C := \{ \mathbf{y} \in \mathbb{R}^m \mid \mathbf{y} = \mathbf{Ax}; \mathbf{x} \geq \mathbf{0}^n \}$. If (I) is infeasible then $\mathbf{b} \notin C$. The set C is a finitely generated cone. Hence, by Corollary 3.29, it follows that C is convex and closed so by the Separation Theorem 3.24 there exist a vector $\boldsymbol{\pi} \neq \mathbf{0}^m$ and $\alpha \in \mathbb{R}$ such that $\mathbf{b}^T \boldsymbol{\pi} > \alpha$ and $\mathbf{y}^T \boldsymbol{\pi} \leq \alpha$ for all $\mathbf{y} \in C$, that is,

$$\mathbf{x}^T \mathbf{A}^T \boldsymbol{\pi} \leq \alpha, \quad \mathbf{x} \geq \mathbf{0}^n. \tag{3.3}$$

Since $\mathbf{0}^m \in C$ it follows that $\alpha \geq 0$, so $\mathbf{b}^T \boldsymbol{\pi} > 0$, and if there exists an $\tilde{\mathbf{x}} \geq \mathbf{0}^n$ such that $\tilde{\mathbf{x}}^T \mathbf{A}^T \boldsymbol{\pi} > 0$, then (3.3) cannot hold for any α (if $\lambda \geq 0$ then $\lambda \tilde{\mathbf{x}} \geq \mathbf{0}^n$ and $(\lambda \tilde{\mathbf{x}})^T \mathbf{A}^T \boldsymbol{\pi} = \lambda \tilde{\mathbf{x}}^T \mathbf{A}^T \boldsymbol{\pi}$ tends to infinity as λ tends to infinity). Therefore we must have that $\mathbf{x}^T \mathbf{A}^T \boldsymbol{\pi} \leq 0$ for all $\mathbf{x} \geq \mathbf{0}^n$, and this holds if and only if $\mathbf{A}^T \boldsymbol{\pi} \leq \mathbf{0}^n$, which means that (II) is feasible.

On the other hand, if (I) has a feasible solution, say $\tilde{\mathbf{x}} \geq \mathbf{0}^n$, then $\mathbf{A}\tilde{\mathbf{x}} = \mathbf{b}$, so if there is a solution to (II), say $\tilde{\boldsymbol{\pi}}$, then $\tilde{\mathbf{x}}^T \mathbf{A}^T \tilde{\boldsymbol{\pi}} = \mathbf{b}^T \tilde{\boldsymbol{\pi}} > 0$. But then $\mathbf{A}^T \tilde{\boldsymbol{\pi}} > \mathbf{0}^n$ (since $\tilde{\mathbf{x}} \geq \mathbf{0}^n$), a contradiction. Hence (II) is infeasible. ■

3.3 Convex functions

Definition 3.31 (convex function) *Suppose that $S \subseteq \mathbb{R}^n$. A function $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is convex at $\bar{\mathbf{x}} \in S$ if*

$$\left. \begin{aligned} \mathbf{x} &\in S \\ \lambda &\in (0, 1) \\ \lambda \bar{\mathbf{x}} + (1 - \lambda) \mathbf{x} &\in S \end{aligned} \right\} \implies f(\lambda \bar{\mathbf{x}} + (1 - \lambda) \mathbf{x}) \leq \lambda f(\bar{\mathbf{x}}) + (1 - \lambda) f(\mathbf{x}).$$

The function f is convex on S if it is convex at every $\bar{\mathbf{x}} \in S$. ■

In other words, a convex function is such that a linear interpolation never is lower than the function itself.¹

From the definition follows that a function $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is convex on a convex set $S \subseteq \mathbb{R}^n$ if and only if

$$\left. \begin{array}{l} \mathbf{x}^1, \mathbf{x}^2 \in S \\ \lambda \in (0, 1) \end{array} \right\} \implies f(\lambda \mathbf{x}^1 + (1 - \lambda) \mathbf{x}^2) \leq \lambda f(\mathbf{x}^1) + (1 - \lambda) f(\mathbf{x}^2).$$

Definition 3.32 (concave function) Suppose that $S \subseteq \mathbb{R}^n$. A function $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is concave at $\bar{\mathbf{x}} \in S$ if $-f$ is convex at $\bar{\mathbf{x}}$.

The function f is concave on S if it is concave at every $\bar{\mathbf{x}} \in S$. ■

Definition 3.33 (strictly convex/concave function) A function $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is strictly convex at $\bar{\mathbf{x}} \in S$ if

$$\left. \begin{array}{l} \mathbf{x} \in S, \mathbf{x} \neq \bar{\mathbf{x}} \\ \lambda \in (0, 1) \\ \lambda \bar{\mathbf{x}} + (1 - \lambda) \mathbf{x} \in S \end{array} \right\} \implies f(\lambda \bar{\mathbf{x}} + (1 - \lambda) \mathbf{x}) < \lambda f(\bar{\mathbf{x}}) + (1 - \lambda) f(\mathbf{x}).$$

The function f is strictly convex (concave) on S if it is strictly convex (concave) at every $\bar{\mathbf{x}} \in S$. ■

In other words, a strictly convex function is such that a linear interpolation is strictly above the function itself.

Figure 3.10 illustrates a strictly convex function.

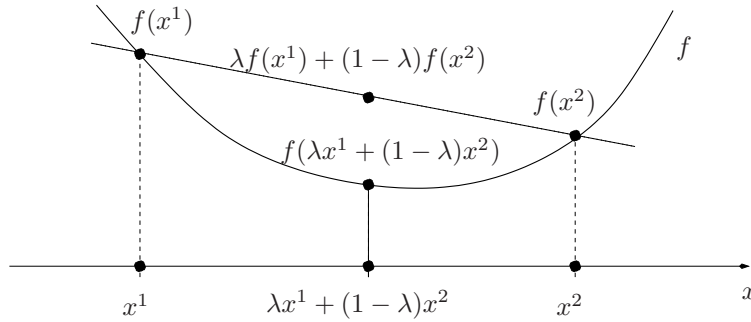


Figure 3.10: A strictly convex function.

¹Words like “lower” and “above” should be understood in the sense of the comparison between the y -coordinates of the respective function at the same coordinates in x .

Example 3.34 (convex functions) By using the definition of a convex function, the following can be established:

(a) The function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ defined by $f(\mathbf{x}) := \|\mathbf{x}\|$ is convex on \mathbb{R}^n .

(b) Let $\mathbf{c} \in \mathbb{R}^n$, $a \in \mathbb{R}$. The affine function $\mathbf{x} \mapsto f(\mathbf{x}) := \mathbf{c}^T \mathbf{x} + a = \sum_{j=1}^n c_j x_j + a$ is both convex and concave on \mathbb{R}^n . The affine functions are also the only finite functions that are both convex and concave. ■

Figure 3.11 illustrates a non-convex function.

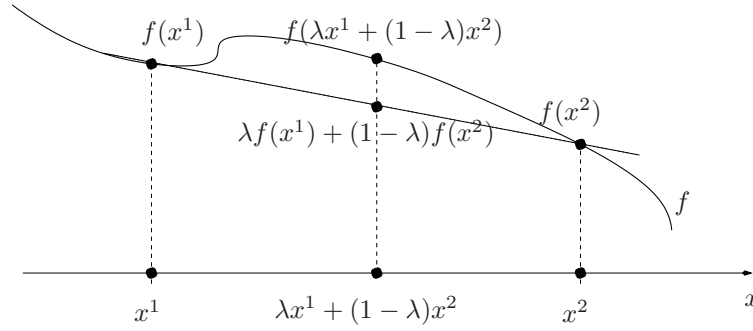


Figure 3.11: A non-convex function.

Proposition 3.35 (sums of convex functions) Suppose that $S \subseteq \mathbb{R}^n$. Let $f_k, k \in \mathcal{K}$, with \mathcal{K} finite, be a collection of functions $f_k : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$. Let $\alpha_k \geq 0, k \in \mathcal{K}$. If each function $f_k, k \in \mathcal{K}$, is convex at $\bar{\mathbf{x}} \in S$, then so is the function $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ defined by $f(\mathbf{x}) := \sum_{k \in \mathcal{K}} \alpha_k f_k(\mathbf{x})$.

In particular, suppose that S is convex and that f_k is convex on S for each $k \in \mathcal{K}$. Then,

$$\alpha_k \geq 0, \quad k \in \mathcal{K} \implies \sum_{k \in \mathcal{K}} \alpha_k f_k \text{ is convex on } S$$

holds.

Proof. The proof is left as an exercise. ■

Proposition 3.36 (convexity of composite functions) Suppose that $S \subseteq \mathbb{R}^n$ and $P \subseteq \mathbb{R}$. Let further $g : S \rightarrow \mathbb{R}$ be a function which is convex on S , and $f : P \rightarrow \mathbb{R}$ be convex and non-decreasing [$y \geq x \implies f(y) \geq f(x)$] on P . Then, the composite function $f(g)$ is convex on the set $\{\mathbf{x} \in S \mid g(\mathbf{x}) \in P\}$.

Proof. Let $\mathbf{x}^1, \mathbf{x}^2 \in S \cap \{\mathbf{x} \in \mathbb{R}^n \mid g(\mathbf{x}) \in P\}$, and $\lambda \in (0, 1)$. Then,

$$\begin{aligned} f(g(\lambda \mathbf{x}^1 + (1 - \lambda)\mathbf{x}^2)) &\leq f(\lambda g(\mathbf{x}^1) + (1 - \lambda)g(\mathbf{x}^2)) \\ &\leq \lambda f(g(\mathbf{x}^1)) + (1 - \lambda)f(g(\mathbf{x}^2)), \end{aligned}$$

where the first inequality follows from the convexity of g and the property of f being increasing, and the second inequality from the convexity of f . ■

The following example functions are important in the development of penalty methods in linear and nonlinear optimization; their convexity is crucial when developing a convergence theory for such algorithms.

Example 3.37 (convex composite functions) Suppose that the function $g : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex.

(a) The function $\mathbf{x} \mapsto -\log(-g(\mathbf{x}))$ is convex on the set $\{\mathbf{x} \in \mathbb{R}^n \mid g(\mathbf{x}) < 0\}$. (This function will be of interest in the analysis of interior point methods; see Section 13.1.)

(b) The function $\mathbf{x} \mapsto -1/g(\mathbf{x})$ is convex on the set $\{\mathbf{x} \in \mathbb{R}^n \mid g(\mathbf{x}) < 0\}$.

[Note: This function is convex, but the above rule for composite functions cannot be used. Utilize the definition of a convex function instead. The domain of the function must here be limited, because $x \mapsto 1/x$ is convex only for positive x .]

(c) The function $\mathbf{x} \mapsto 1/\log(-g(\mathbf{x}))$ is convex on the set $\{\mathbf{x} \in \mathbb{R}^n \mid g(\mathbf{x}) < -1\}$.

[Note: This function is convex, but the above rule for composite functions cannot be used. Utilize the definition of a convex function instead. The domain of the function must here be limited, because $x \mapsto 1/x$ is convex only for positive x .] ■

We next characterize the convexity of a function on \mathbb{R}^n by the convexity of its *epigraph* in \mathbb{R}^{n+1} .

[Note: the *graph* of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is the boundary of $\text{epi } f$, which still resides in \mathbb{R}^{n+1} . See Figure 3.12 for an example, corresponding to the convex function in Figure 3.10.]

Definition 3.38 (epigraph) The epigraph of a function $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is the set

$$\text{epi } f := \{(\mathbf{x}, \alpha) \in \mathbb{R}^n \times \mathbb{R} \mid f(\mathbf{x}) \leq \alpha\}. \quad (3.4)$$

The epigraph of the function f restricted to the set $S \subseteq \mathbb{R}^n$ is

$$\text{epi}_S f := \{(\mathbf{x}, \alpha) \in S \times \mathbb{R} \mid f(\mathbf{x}) \leq \alpha\}. \quad (3.5)$$

■

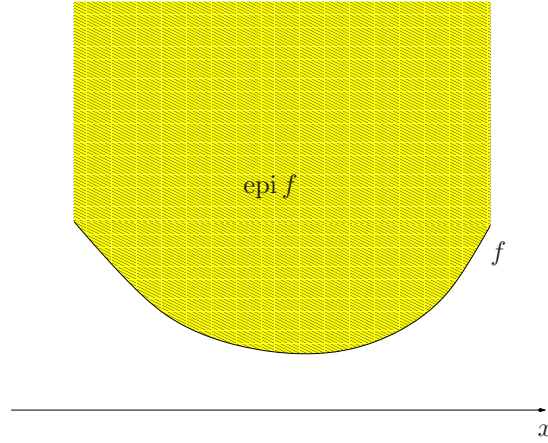


Figure 3.12: A convex function and its epigraph.

Theorem 3.39 Suppose that $S \subseteq \mathbb{R}^n$ is a convex set. Then, the function $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is convex on S if, and only if, its epigraph restricted to S is a convex set in \mathbb{R}^{n+1} .

Proof. $[\implies]$ Suppose that f is convex on S . Let $(\mathbf{x}^1, \alpha_1), (\mathbf{x}^2, \alpha_2) \in \text{epi}_S f$. Let $\lambda \in (0, 1)$. By the convexity of f on S ,

$$\begin{aligned} f(\lambda \mathbf{x}^1 + (1 - \lambda) \mathbf{x}^2) &\leq \lambda f(\mathbf{x}^1) + (1 - \lambda) f(\mathbf{x}^2) \\ &\leq \lambda \alpha_1 + (1 - \lambda) \alpha_2. \end{aligned}$$

Hence, $[\lambda \mathbf{x}^1 + (1 - \lambda) \mathbf{x}^2, \lambda \alpha_1 + (1 - \lambda) \alpha_2] \in \text{epi}_S f$, so $\text{epi}_S f$ is a convex set in \mathbb{R}^{n+1} .

$[\impliedby]$ Suppose that $\text{epi}_S f$ is convex. Let $\mathbf{x}^1, \mathbf{x}^2 \in S$, whence

$$(\mathbf{x}^1, f(\mathbf{x}^1)), (\mathbf{x}^2, f(\mathbf{x}^2)) \in \text{epi}_S f.$$

Let $\lambda \in (0, 1)$. By the convexity of $\text{epi}_S f$, it follows that

$$[\lambda \mathbf{x}^1 + (1 - \lambda) \mathbf{x}^2, \lambda f(\mathbf{x}^1) + (1 - \lambda) f(\mathbf{x}^2)] \in \text{epi}_S f,$$

that is, $f(\lambda \mathbf{x}^1 + (1 - \lambda) \mathbf{x}^2) \leq \lambda f(\mathbf{x}^1) + (1 - \lambda) f(\mathbf{x}^2)$. Hence, f is convex on S . ■

When f is in C^1 (once differentiable, with continuous partial derivatives) or C^2 (twice differentiable, with continuous partial second derivatives), then convexity can be characterized also in terms of these derivatives. The results show how with stronger differentiability properties the characterizations become more and more useful in practice.

Theorem 3.40 (convexity characterizations in C^1) *Let $f \in C^1$ on an open convex set S .*

(a) *f is convex on $S \iff f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x})$, for all $\mathbf{x}, \mathbf{y} \in S$.*

(b) *f is convex on $S \iff [\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})]^T(\mathbf{x} - \mathbf{y}) \geq 0$, for all $\mathbf{x}, \mathbf{y} \in S$.*

The result in (a) states, in words, that “every tangent plane to the function surface in \mathbb{R}^{n+1} lies on, or below, the epigraph of f ”, or, that “a first-order approximation is below f .”

The result in (b) states that ∇f is “monotone on S .”

[Note: when $n = 1$, the result in (b) states that f is convex if and only if its derivative f' is non-decreasing, that is, that it is monotonically increasing.]

Proof. (a) \implies Take $\mathbf{x}^1, \mathbf{x}^2 \in S$ and $\lambda \in (0, 1)$. Then,

$$\begin{aligned} \lambda f(\mathbf{x}^1) + (1 - \lambda)f(\mathbf{x}^2) &\geq f(\lambda\mathbf{x}^1 + (1 - \lambda)\mathbf{x}^2) \\ &\iff [\lambda > 0] \\ f(\mathbf{x}^1) - f(\mathbf{x}^2) &\geq (1/\lambda)[f(\lambda\mathbf{x}^1 + (1 - \lambda)\mathbf{x}^2) - f(\mathbf{x}^2)]. \end{aligned}$$

Let $\lambda \downarrow 0$. Then, the right-hand side of the above inequality tends to the directional derivative of f at \mathbf{x}^2 in the direction of $(\mathbf{x}^1 - \mathbf{x}^2)$, so that in the limit it becomes

$$f(\mathbf{x}^1) - f(\mathbf{x}^2) \geq \nabla f(\mathbf{x}^2)^T(\mathbf{x}^1 - \mathbf{x}^2).$$

The result follows.

\impliedby We have that

$$\begin{aligned} f(\mathbf{x}^1) &\geq f(\lambda\mathbf{x}^1 + (1 - \lambda)\mathbf{x}^2) + (1 - \lambda)\nabla f(\lambda\mathbf{x}^1 + (1 - \lambda)\mathbf{x}^2)^T(\mathbf{x}^1 - \mathbf{x}^2), \\ f(\mathbf{x}^2) &\geq f(\lambda\mathbf{x}^1 + (1 - \lambda)\mathbf{x}^2) + \lambda\nabla f(\lambda\mathbf{x}^1 + (1 - \lambda)\mathbf{x}^2)^T(\mathbf{x}^2 - \mathbf{x}^1). \end{aligned}$$

Multiply the inequalities by λ and $(1 - \lambda)$, respectively, and add them together to get the result sought.

(b) \implies Using (a), and the two inequalities

$$\begin{aligned} f(\mathbf{y}) &\geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}), & \mathbf{x}, \mathbf{y} \in S, \\ f(\mathbf{x}) &\geq f(\mathbf{y}) + \nabla f(\mathbf{y})^T(\mathbf{x} - \mathbf{y}), & \mathbf{x}, \mathbf{y} \in S, \end{aligned}$$

added together, yields that $[\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})]^T(\mathbf{x} - \mathbf{y}) \geq 0$, for all $\mathbf{x}, \mathbf{y} \in S$.

\impliedby The mean-value theorem states that

$$f(\mathbf{x}^2) - f(\mathbf{x}^1) = \nabla f(\mathbf{x})^T(\mathbf{x}^2 - \mathbf{x}^1), \quad (3.6)$$

where $\mathbf{x} = \lambda \mathbf{x}^1 + (1 - \lambda) \mathbf{x}^2$ for some $\lambda \in (0, 1)$. By assumption, $[\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}^1)]^T (\mathbf{x} - \mathbf{x}^1) \geq 0$, so $(1 - \lambda)[\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}^1)]^T (\mathbf{x}^2 - \mathbf{x}^1) \geq 0$. From this follows that $\nabla f(\mathbf{x})^T (\mathbf{x}^2 - \mathbf{x}^1) \geq \nabla f(\mathbf{x}^1)^T (\mathbf{x}^2 - \mathbf{x}^1)$. By using this inequality and (3.6), we get that $f(\mathbf{x}^2) \geq f(\mathbf{x}^1) + \nabla f(\mathbf{x}^1)^T (\mathbf{x}^2 - \mathbf{x}^1)$. We are done. ■

Figure 3.13 illustrates part (a) of Theorem 3.40.

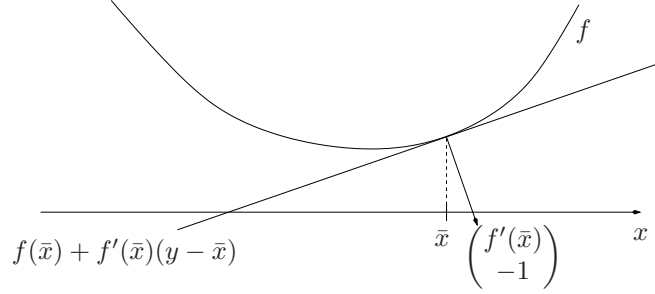


Figure 3.13: A tangent plane to the graph of a convex function.

By replacing the inequalities in (a) and (b) in the theorem by *strict* inequalities, and adding the requirement that $\mathbf{x} \neq \mathbf{y}$ holds in the statements, we can establish a characterization also of *strictly* convex functions. The statement in (a) then says that the tangential hyperplane lies strictly below the function except at the tangent point, and (b) states that the gradient mapping is *strictly* monotone.

Still more can be said in C^2 :

Theorem 3.41 (convexity characterizations in C^2 , I) *Let f be in C^2 on an open, convex set $S \subseteq \mathbb{R}^n$.*

- (a) *f is convex on $S \iff \nabla^2 f(\mathbf{x})$ is positive semidefinite for all $\mathbf{x} \in S$.*
- (b) *$\nabla^2 f(\mathbf{x})$ is positive definite for all $\mathbf{x} \in S \implies f$ is strictly convex on S .*

[Note: When $n = 1$ and S is an open interval, the above reduce to the following familiar results: (a) f is convex on S if and only if $f''(x) \geq 0$ for every $x \in S$; (b) f is strictly convex on S if $f''(x) > 0$ for every $x \in S$.]

Proof. (a) [\implies] Suppose that f is convex and let $\bar{\mathbf{x}} \in S$. We must show that $\mathbf{p}^T \nabla^2 f(\bar{\mathbf{x}}) \mathbf{p} \geq 0$ for all $\mathbf{p} \in \mathbb{R}^n$ holds.

Since S open, for any given $\mathbf{p} \in \mathbb{R}^n$, $\bar{\mathbf{x}} + \alpha \mathbf{p} \in S$ whenever $|\alpha| \neq 0$ is small enough. We utilize Theorem 3.40(a) as follows: by the twice

differentiability of f ,

$$f(\bar{\mathbf{x}} + \alpha \mathbf{p}) \geq f(\bar{\mathbf{x}}) + \alpha \nabla f(\bar{\mathbf{x}})^T \mathbf{p}, \quad (3.7)$$

$$f(\bar{\mathbf{x}} + \alpha \mathbf{p}) = f(\bar{\mathbf{x}}) + \alpha \nabla f(\bar{\mathbf{x}})^T \mathbf{p} + \frac{1}{2} \alpha^2 \mathbf{p}^T \nabla^2 f(\bar{\mathbf{x}}) \mathbf{p} + o(\alpha^2). \quad (3.8)$$

Subtracting (3.8) from (3.7), we get

$$\frac{1}{2} \alpha^2 \mathbf{p}^T \nabla^2 f(\bar{\mathbf{x}}) \mathbf{p} + o(\alpha^2) \geq 0.$$

Dividing by α^2 and letting $\alpha \rightarrow 0$ it follows that $\mathbf{p}^T \nabla^2 f(\bar{\mathbf{x}}) \mathbf{p} \geq 0$.

[\Leftarrow] Suppose that the Hessian matrix is positive semidefinite at each point in S . The proof depends on the following second-order mean-value theorem: for every $\mathbf{x}, \mathbf{y} \in S$, there exists $\ell \in [0, 1]$ such that

$$f(\mathbf{y}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{1}{2} (\mathbf{y} - \mathbf{x})^T \nabla^2 f[\mathbf{x} + \ell(\mathbf{y} - \mathbf{x})] (\mathbf{y} - \mathbf{x}). \quad (3.9)$$

By assumption, the last term in (3.9) is non-negative, whence we obtain the convexity characterization in Theorem 3.40(a).

(b) [\Rightarrow] By the assumptions, the last term in (3.9) is always positive when $\mathbf{y} \neq \mathbf{x}$, whence we obtain the strict convexity characterization in C^1 . ■

It is important to note that the opposite direction in the result (b) is false. A simple example that establishes this fact is the function defined by $f(x) := x^4$, $S := \mathbb{R}$; f is strictly convex on \mathbb{R} (why?), but its second derivative at zero is $f''(0) = 0$.

The case of quadratic functions is interesting to mention in particular. For quadratic functions, that is, functions of the form

$$f(\mathbf{x}) := (1/2) \mathbf{x}^T \mathbf{Q} \mathbf{x} - \mathbf{q}^T \mathbf{x} + a, \quad (3.10)$$

for some symmetric matrix $\mathbf{Q} \in \mathbb{R}^{n \times n}$, vector $\mathbf{q} \in \mathbb{R}^n$ and scalar $a \in \mathbb{R}$, it holds that $\nabla^2 f(\mathbf{x}) \equiv \mathbf{Q}$ for every \mathbf{x} where f is defined, so the value $\nabla^2 f(\mathbf{x})$ does not depend on \mathbf{x} . In this case, we can state a stronger result than in Theorem 3.41: *the quadratic function f is convex on the open, convex set $S \subseteq \mathbb{R}^n$ if and only if \mathbf{Q} is positive semidefinite; f is strictly convex on S if and only if \mathbf{Q} is positive definite*. To prove this result is simple from the above result for general C^2 functions, and is left as an exercise.

What happens when S is not full-dimensional (which is often the case)? Take, for example, $f(\mathbf{x}) := x_1^2 - x_2^2$ and $S := \{\mathbf{x} \in \mathbb{R}^2 \mid x_1 \in \mathbb{R}; x_2 = 0\}$. Then, f is convex on S but $\nabla^2 f(\mathbf{x})$ is not positive semidefinite anywhere on S . The below result covers this type of case. Its proof is left as an exercise.

Theorem 3.42 (convexity characterizations in C^2 , II) *Let $S \subseteq \mathbb{R}^n$ be a nonempty convex set and $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be in C^2 on \mathbb{R}^n . Let C be the subspace parallel to the affine hull of S . Then,*

$$f \text{ is convex on } S \iff \mathbf{p}^T \nabla^2 f(\mathbf{x}) \mathbf{p} \geq 0 \text{ for every } \mathbf{x} \in S \text{ and } \mathbf{p} \in C.$$

In particular, when S has a nonempty interior, f is convex if and only if $\nabla^2 f(\mathbf{x})$ is positive semidefinite for every $\mathbf{x} \in S$. ■

We have already seen that the convexity of a function is intimately connected to the convexity of a certain set, namely the epigraph of the function. The following result shows that a particular type of set, defined by those vectors that bound a convex function from above, is a convex set. Later, we will utilize this result to establish the convexity of feasible sets in some optimization problems.

Definition 3.43 (level set) *Let $g : \mathbb{R}^n \rightarrow \mathbb{R}$ be a function. The level set of g with respect to the value $b \in \mathbb{R}$ is the set*

$$\text{lev}_g(b) := \{ \mathbf{x} \in \mathbb{R}^n \mid g(\mathbf{x}) \leq b \}. \quad (3.11)$$

■

Figure 3.14 illustrates a level set of a convex function.

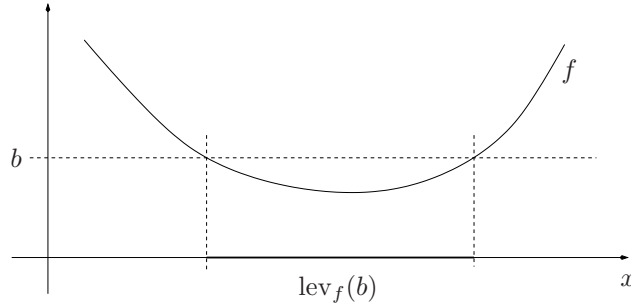


Figure 3.14: A level set of a convex function.

Proposition 3.44 (convex level sets from convex functions) *Suppose that the function $g : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex. Then, for every value of $b \in \mathbb{R}$, the level set $\text{lev}_g(b)$ is a convex set. It is moreover closed.*

Proof. The result follows immediately from the definitions of a convex set and a convex function. Let $\mathbf{x}^1, \mathbf{x}^2$ both satisfy the constraint that

$g(\mathbf{x}) \leq b$ holds, and let $\lambda \in (0, 1)$. (If not two such points $\mathbf{x}^1, \mathbf{x}^2$ can be found, then the result holds vacuously.) Then, by the convexity of g , $g(\lambda\mathbf{x}^1 + (1 - \lambda)\mathbf{x}^2) \leq \lambda b + (1 - \lambda)b = b$, so the set $\text{lev}_g(b)$ is convex.

The fact that a convex function which is defined on \mathbb{R}^n is continuous establishes that the set $\text{lev}_g(b)$ is always closed.² (Why?) ■

Definition 3.45 (convex problem) Suppose that the set $X \subseteq \mathbb{R}^n$ is closed and convex. Suppose further that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex and that the functions $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$, $i \in \mathcal{I}$, are convex. Suppose, finally, that the functions $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$, $i \in \mathcal{E}$, are affine. Then, the problem to

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}), \\ & \text{subject to} && g_i(\mathbf{x}) \leq 0, \quad i \in \mathcal{I}, \\ & && g_i(\mathbf{x}) = 0, \quad i \in \mathcal{E}, \\ & && \mathbf{x} \in X, \end{aligned}$$

is called a convex problem. ■

The name is natural, because the objective function is a convex one, and the feasible set is closed and convex as well. In order to establish the latter, we refer first to Proposition 3.44 to establish that the inequality constraints define convex sets [note that in the similar problem (1.1) the inequalities are given as \geq -constraints, and then we require g_i , $i \in \mathcal{I}$, to be concave functions in order to have a convex problem], and ask the reader to prove that a constraint of the form $\mathbf{a}_i^T \mathbf{x} = b_i$ defines a convex set as well. Finally, we refer to Proposition 3.3 to establish that the intersection of all the convex sets defined by X , \mathcal{I} , and \mathcal{E} is convex.

3.4 Application: the projection of a vector onto a convex set

In Figure 3.15 we illustrate the Euclidean projection of some vectors onto a convex set.

We see that the Euclidean projection of $\mathbf{w} \in \mathbb{R}^n$ is the vector in S which is nearest (in the Euclidean norm) to \mathbf{w} : the vector $\text{Proj}_S(\mathbf{w})$ is the unique optimum in the problem of finding

$$\text{minimum}_{\mathbf{x} \in S} \|\mathbf{x} - \mathbf{w}\|.$$

²That convex functions are continuous will be established in Theorem 4.27.

Application: the projection of a vector onto a convex set

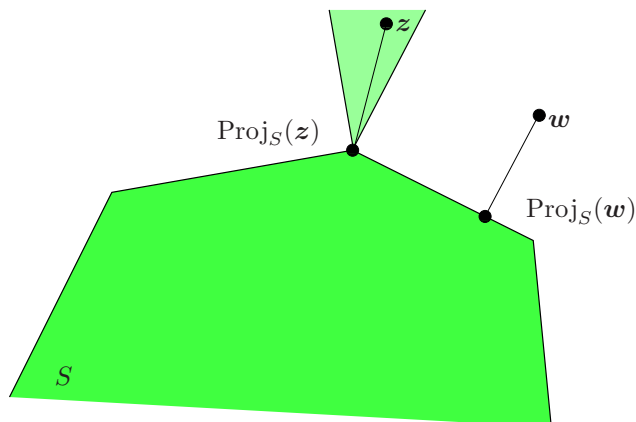


Figure 3.15: The projection of two vectors onto a convex set.

The vector $w - \text{Proj}_S(w)$ clearly is *normal* to the set S . The point z has the Euclidean projection $\text{Proj}_S(z)$, but there are also several other vectors with the same projection; the figure shows in a special shading the set of vectors z which all have that same projection onto S . This set is a cone, which we refer to as the *normal cone* to S at $x = \text{Proj}_S(z)$. In the case of the point $\text{Proj}_S(w)$ the normal cone reduces to a ray—which of course is also a cone. (The difference between these two sets is largely the consequence of the fact that there is only one constraint active at $\text{Proj}_S(w)$, while there are two constraints active at $\text{Proj}_S(z)$; when developing the KKT conditions in Chapter 5 we shall see how strongly the active constraints influence the appearance of the optimality conditions.)

We will also return to this image already in Section 4.6.3, because it contains the building blocks of the optimality conditions for an optimization problem with an objective function in C^1 over a closed convex set. For now, we will establish only one property of the projection operation Proj_S , namely that the *distance function*, dist_S , defined by

$$\text{dist}_S(x) := \|x - \text{Proj}_S(x)\|, \quad x \in \mathbb{R}^n, \quad (3.12)$$

is a convex function on \mathbb{R}^n . In particular, then, this function is continuous. (Later, we will establish also that the projection operation Proj_S is a well-defined operation whenever S is nonempty, closed and convex, and that the operation has particularly nice continuity properties. Before we can do so, however, we need to establish some results on the existence

of optimal solutions.)

Let $\mathbf{x}^1, \mathbf{x}^2 \in \mathbb{R}^n$, and $\lambda \in (0, 1)$. Then,

$$\begin{aligned}
 \text{dist}_S(\lambda \mathbf{x}^1 + (1 - \lambda) \mathbf{x}^2) &= \|(\lambda \mathbf{x}^1 + (1 - \lambda) \mathbf{x}^2) \\
 &\quad - \text{Proj}_S(\lambda \mathbf{x}^1 + (1 - \lambda) \mathbf{x}^2)\| \\
 &\leq \|(\lambda \mathbf{x}^1 + (1 - \lambda) \mathbf{x}^2) \\
 &\quad - (\lambda \text{Proj}_S(\mathbf{x}^1) + (1 - \lambda) \text{Proj}_S(\mathbf{x}^2))\| \\
 &\leq \lambda \|\mathbf{x}^1 - \text{Proj}_S(\mathbf{x}^1)\| \\
 &\quad + (1 - \lambda) \|\mathbf{x}^2 - \text{Proj}_S(\mathbf{x}^2)\| \\
 &= \lambda \text{dist}_S(\mathbf{x}^1) + (1 - \lambda) \text{dist}_S(\mathbf{x}^2),
 \end{aligned}$$

where the first inequality comes from the fact that $\lambda \text{Proj}_S(\mathbf{x}^1) + (1 - \lambda) \text{Proj}_S(\mathbf{x}^2) \in S$, but it does not necessarily define $\text{Proj}_S(\lambda \mathbf{x}^1 + (1 - \lambda) \mathbf{x}^2)$ (it may have a longer distance), and the second is the triangle inequality.

The proof is illustrated in Figure 3.16.

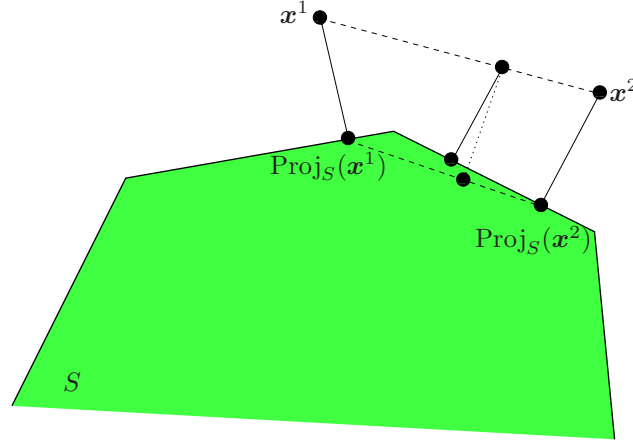


Figure 3.16: The distance function is convex. From the intermediate vector $\lambda \mathbf{x}^1 + (1 - \lambda) \mathbf{x}^2$, the distance to the vector $\lambda \text{Proj}_S(\mathbf{x}^1) + (1 - \lambda) \text{Proj}_S(\mathbf{x}^2)$ [the dotted line segment] clearly is longer than the distance to its projection on S [shown as a solid line].

3.5 Notes and further reading

The subject of this chapter—convex analysis—has a long history, going back about a century. Much of the early work on convex sets and functions, for example, the theory of separation of convex sets, go back to the work of Minkowski [Min10, Min11]. Other expositions are found in [Fen51, Roc70, StW70], which all are classical in the field. More easily accessible are the modern books [BoL00, BNO03]. Lighter introductions are also found in [BSS93, HiL93]. The most influential of all of these is *Convex Analysis* by R. T. Rockafellar [Roc70].

Carathéodory's Theorem 3.8 is found in [Car07, Car11]. Farkas' Lemma in Theorem 3.30 is found in [Far1902]. Theorem 3.42 is given as Exercise 1.8 in [BNO03].

The early history of polyhedral convexity is found in [Mot36].

3.6 Exercises

Exercise 3.1 (convexity of polyhedra) Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$. Show that the polyhedron

$$P := \{ \mathbf{x} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{x} \leq \mathbf{b} \},$$

is a convex set.

Exercise 3.2 (polyhedra) Which of the following sets are polyhedra?

(a) $S := \{ y_1 \mathbf{a} + y_2 \mathbf{b} \mid -1 \leq y_1 \leq 1; -1 \leq y_2 \leq 1 \}$, where $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ are fixed.

(b) $S := \{ \mathbf{x} \in \mathbb{R}^n \mid \mathbf{x} \geq \mathbf{0}^n; \mathbf{x}^T \mathbf{1}^n = 1; \sum_{i=1}^n x_i a_i = b_1; \sum_{i=1}^n x_i a_i^2 = b_2 \}$, where $a_i \in \mathbb{R}$ for $i = 1, \dots, n$, and $b_1, b_2 \in \mathbb{R}$ are fixed.

(c) $S := \{ \mathbf{x} \in \mathbb{R}^n \mid \mathbf{x} \geq \mathbf{0}^n; \mathbf{x}^T \mathbf{y} \leq 1 \text{ for all } \mathbf{y} \text{ such that } \|\mathbf{y}\|_2 = 1 \}$.

(d) $S := \{ \mathbf{x} \in \mathbb{R}^n \mid \mathbf{x} \geq \mathbf{0}^n; \mathbf{x}^T \mathbf{y} \leq 1 \text{ for all } \mathbf{y} \text{ such that } \sum_{i=1}^n |y_i| = 1 \}$.

(e) $S := \{ \mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x} - \mathbf{x}^0\|_2 \leq \|\mathbf{x} - \mathbf{x}^1\|_2 \}$, where $\mathbf{x}^0, \mathbf{x}^1 \in \mathbb{R}^n$ are fixed.

(f) $S := \{ \mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x} - \mathbf{x}^0\|_2 \leq \|\mathbf{x} - \mathbf{x}^i\|_2, i = 1, \dots, k \}$, where $\mathbf{x}^0, \dots, \mathbf{x}^k \in \mathbb{R}^n$ are fixed.

Exercise 3.3 (extreme points) Consider the polyhedron P defined by

$$x_1 + x_2 \leq 2,$$

$$x_2 \leq 1,$$

$$x_3 \leq 2,$$

$$x_2 + x_3 \leq 2.$$

(a) Is $\mathbf{x}^1 := (1, 1, 0)^T$ an extreme point to P ?

(b) Is $\mathbf{x}^2 := (1, 1, 1)^T$ an extreme point to P ?

Exercise 3.4 (existence of extreme points in LPs) Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ be such that $\text{rank } \mathbf{A} = m$, and let $\mathbf{b} \in \mathbb{R}^m$. Show that if the polyhedron

$$P := \{ \mathbf{x} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{x} = \mathbf{b}; \quad \mathbf{x} \geq \mathbf{0}^n \}$$

has a feasible solution, then it has an extreme point.

Exercise 3.5 (illustration of the Representation Theorem) Let

$$Q := \{ \mathbf{x} \in \mathbb{R}^2 \mid -2x_1 + x_2 \leq 1; \quad x_1 - x_2 \leq 1; \quad -x_1 - x_2 \leq -1 \},$$

$$C := \{ \mathbf{x} \in \mathbb{R}^2 \mid -2x_1 + x_2 \leq 0; \quad x_1 - x_2 \leq 0; \quad -x_1 - x_2 \leq 0 \},$$

and P be the convex hull of the extreme points of Q . Show that the feasible point $\tilde{\mathbf{x}} = (1, 1)^T$ can be written as

$$\tilde{\mathbf{x}} = \mathbf{p} + \mathbf{c},$$

where $\mathbf{p} \in P$ and $\mathbf{c} \in C$.

Exercise 3.6 (separation) Show that there is only one hyperplane in \mathbb{R}^3 which separates the disjoint closed convex sets A and B defined by

$$A := \{ (0, x_2, 1)^T \mid x_2 \in \mathbb{R} \}, \quad B := \{ \mathbf{x} \in \mathbb{R}^3 \mid \mathbf{x} \geq \mathbf{0}^3; \quad x_1 x_2 \geq x_3^2 \},$$

and that this hyperplane meets both A and B .

Exercise 3.7 (separation) Show that each closed convex set A in \mathbb{R}^n is the intersection of all the closed halfspaces in \mathbb{R}^n containing A .

Exercise 3.8 (application of Farkas' Lemma) In a paper submitted for publication in an operations research journal, the author considered the set

$$P := \left\{ \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \in \mathbb{R}^{n+m} \mid \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} \geq \mathbf{c}; \quad \mathbf{x} \geq \mathbf{0}^n; \quad \mathbf{y} \geq \mathbf{0}^m \right\},$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{B} \in \mathbb{R}^{m \times m}$ is positive semidefinite, and $\mathbf{c} \in \mathbb{R}^m$. The author explicitly assumed that the set P is compact in \mathbb{R}^{n+m} . A reviewer of the paper pointed out that the only compact set of the above form is the empty set. Prove the reviewer's assertion.

Exercise 3.9 (convex sets) Let $S_1 := \{ \mathbf{x} \in \mathbb{R}^2 \mid x_1 + x_2 \leq 1; \quad x_1 \geq 0 \}$, $S_2 := \{ \mathbf{x} \in \mathbb{R}^2 \mid x_1 - x_2 \geq 0; \quad x_1 \leq 1 \}$, and $S := S_1 \cup S_2$. Prove that S_1 and S_2 are convex sets and that S is not convex. Hence, the union of convex sets is not necessarily a convex set.

Exercise 3.10 (convex functions) Determine if the function $f(\mathbf{x}) := 2x_1^2 - 3x_1x_2 + 5x_2^2 - 2x_1 + 6x_2$ is convex, concave, or neither, on \mathbb{R}^2 .

Exercise 3.11 (convex functions) Let $a > 0$. Consider the following functions in one variable:

- (a) $f(x) := \ln x$, for $x > 0$;
- (b) $f(x) := -\ln x$, for $x > 0$;
- (c) $f(x) := -\ln(1 - e^{-ax})$, for $x > 0$;
- (d) $f(x) := \ln(1 + e^{ax})$;
- (e) $f(x) := e^{ax}$;
- (f) $f(x) := x \ln x$, for $x > 0$.

Which of these functions are convex (or, strictly convex)?

Exercise 3.12 (convex functions) Consider the following functions:

- (a) $f(\mathbf{x}) := \ln(e^{x_1} + e^{x_2})$;
- (b) $f(\mathbf{x}) := \ln \sum_{j=1}^n e^{a_j x_j}$, where a_j , $j = 1, \dots, n$, are constants;
- (c) $f(\mathbf{x}) := \sqrt{\sum_{j=1}^n x_j^2}$;
- (d) $f(\mathbf{x}) := x_1^2/x_2$, for $x_2 > 0$;
- (e) $f(\mathbf{x}) := -\sqrt{x_1 x_2}$, for $x_1, x_2 > 0$;
- (f) $f(\mathbf{x}) := -\left(\prod_{j=1}^n x_j\right)^{1/n}$, for $x_j > 0$, $j = 1, \dots, n$.

Which of these functions are convex (or, strictly convex)?

Exercise 3.13 (convex functions) Consider the following function:

$$f(x, y) := 2x^2 - 2xy + \frac{1}{2}y^2 + 3x - y.$$

- (a) Express the function in matrix-vector form.
- (b) Is the Hessian singular?
- (c) Is f a convex function?

Exercise 3.14 (convex sets) Consider the following sets:

- (a) $\{\mathbf{x} \in \mathbb{R}^2 \mid x_1^2 + x_2^2 \leq 1; x_1^2 + x_2^2 \geq 1/4\}$;
- (b) $\{\mathbf{x} \in \mathbb{R}^n \mid x_j \geq 0, j = 1, \dots, n\}$;
- (c) $\{\mathbf{x} \in \mathbb{R}^n \mid x_1^2 + x_2^2 + \dots + x_n^2 = 1\}$;
- (d) $\{\mathbf{x} \in \mathbb{R}^2 \mid x_1 + x_2 \leq 5; x_1 - x_2 \leq 10; x_1 \geq 0; x_2 \geq 0\}$;
- (e) $\{\mathbf{x} \in \mathbb{R}^2 \mid x_1 - x_2 \geq 1; x_1^3 + x_2^2 \leq 10; 2x_1 + x_2 \leq 8; x_1 \geq 1; x_2 \geq 0\}$.

Investigate, in each case, whether the set defined is convex or not. In the latter case, provide a counter-example.

Exercise 3.15 (convex sets) Is the set defined by

$$S := \{\mathbf{x} \in \mathbb{R}^2 \mid 2e^{-x_1+x_2^2} \leq 4; -x_1^2 + 3x_1x_2 - 3x_2^2 \geq -1\}$$

a convex set?

Exercise 3.16 (convex sets) Is the set defined by

$$S := \{ \mathbf{x} \in \mathbb{R}^2 \mid x_1 - x_2^2 \geq 1; x_1^3 + x_2^2 \leq 10; 2x_1 + x_2 \leq 8; x_1 \geq 1; x_2 \geq 0 \}$$

a convex set?

Exercise 3.17 (convex problem) Suppose that the function $g : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex on \mathbb{R}^n and that $\mathbf{d} \in \mathbb{R}^n$. Is the problem to

$$\begin{aligned} \text{maximize} \quad & -\sum_{j=1}^n x_j^2, \\ \text{subject to} \quad & -\frac{1}{\ln(-g(\mathbf{x}))} \geq 0, \\ & \mathbf{d}^T \mathbf{x} = 2, \\ & g(\mathbf{x}) \leq -2, \\ & \mathbf{x} \geq \mathbf{0}^n \end{aligned}$$

a convex problem?

Exercise 3.18 (convex problem) Is the problem to

$$\begin{aligned} \text{maximize} \quad & x_1 \ln x_1, \\ \text{subject to} \quad & x_1^2 + x_2^2 \geq 0, \\ & \mathbf{x} \geq \mathbf{0}^2 \end{aligned}$$

a convex problem?

Part III

Optimality Conditions

An introduction to optimality conditions

IV

4.1 Local and global optimality

Consider the problem to

$$\text{minimize } f(\mathbf{x}), \quad (4.1a)$$

$$\text{subject to } \mathbf{x} \in S, \quad (4.1b)$$

where $S \subseteq \mathbb{R}^n$ is a nonempty set and $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is a given function.

Consider the function given in Figure 4.1.

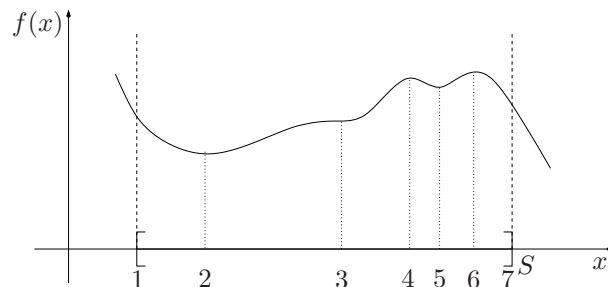


Figure 4.1: A one-dimensional function and its possible extremal points.

For a minimization problem of f in one variable over a closed interval S , the interesting points are:

- (i) boundary points of S ;
- (ii) stationary points, that is, where $f'(x) = 0$;
- (iii) discontinuities in f or f' .

In the case of the function in Figure 4.1 the points 1 and 7 are of category (i); 2, 3, 4, 5, and 6 of category (ii); and none of category (iii).

Definition 4.1 (global minimum) *Consider the problem (4.1). Let $\mathbf{x}^* \in S$. We say that \mathbf{x}^* is a global minimum of f over S if f attains its lowest value over S at \mathbf{x}^* .*

In other words, $\mathbf{x}^ \in S$ is a global minimum of f over S if*

$$f(\mathbf{x}^*) \leq f(\mathbf{x}), \quad \mathbf{x} \in S \quad (4.2)$$

holds. ■

Let $B_\varepsilon(\mathbf{x}^*) := \{\mathbf{y} \in \mathbb{R}^n \mid \|\mathbf{y} - \mathbf{x}^*\| < \varepsilon\}$ be the open Euclidean ball with radius ε centered at \mathbf{x}^* .

Definition 4.2 (local minimum) *Consider the problem (4.1). Let $\mathbf{x}^* \in S$.*

(a) *We say that \mathbf{x}^* is a local minimum of f over S if there exists a small enough ball intersected with S around \mathbf{x}^* such that it is a globally optimal solution in that smaller set.*

In other words, $\mathbf{x}^ \in S$ is a local minimum of f over S if*

$$\exists \varepsilon > 0 \text{ such that } f(\mathbf{x}^*) \leq f(\mathbf{x}), \quad \mathbf{x} \in S \cap B_\varepsilon(\mathbf{x}^*). \quad (4.3)$$

(b) *We say that $\mathbf{x}^* \in S$ is a strict local minimum of f over S if, in (4.3), the inequality holds strictly for $\mathbf{x} \neq \mathbf{x}^*$.* ■

Note that a global minimum in particular is a local minimum. When is a local minimum a global one? This question is resolved in the case of convex problems, as the following fundamental theorem shows.

Theorem 4.3 (Fundamental Theorem of global optimality) *Consider the problem (4.1), where S is a convex set and f is convex on S . Then, every local minimum of f over S is also a global minimum.*

Proof. Suppose that \mathbf{x}^* is a local minimum but not a global one, while $\bar{\mathbf{x}}$ is a global minimum. Then, $f(\bar{\mathbf{x}}) < f(\mathbf{x}^*)$. Let $\lambda \in (0, 1)$. By the convexity of S and f , $\lambda\bar{\mathbf{x}} + (1 - \lambda)\mathbf{x}^* \in S$, and $f(\lambda\bar{\mathbf{x}} + (1 - \lambda)\mathbf{x}^*) \leq \lambda f(\bar{\mathbf{x}}) + (1 - \lambda)f(\mathbf{x}^*) < f(\mathbf{x}^*)$. Choosing $\lambda > 0$ small enough then leads to a contradiction to the local optimality of \mathbf{x}^* . ■

There is an intuitive image that can be seen from the proof design: If \mathbf{x}^* is a local minimum, then f cannot “go down-hill” from \mathbf{x}^* in any direction, but if $\bar{\mathbf{x}}$ has a lower value, then f has to go down-hill sooner or later. No convex function can have this shape.

The example in Figure 4.2 shows a case where, without convexity, a vector \mathbf{x}^* may be a local minimum of a function $f \in C^1$ with respect to every *line segment* that passes through \mathbf{x}^* , and yet it is not even a local minimum of f over \mathbb{R}^n .

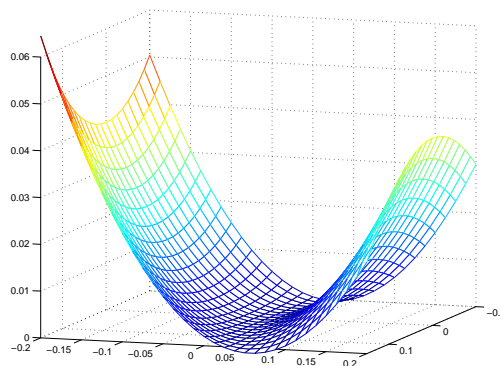


Figure 4.2: A three-dimensional graph of the function $f(x, y) := (y - x^2)(y - 4x^2)$. The origin is a local minimum with respect to every line that passes through it, but it is not a local minimum of f .

In fact, this situation may also occur in the convex case when $f \notin C^1$.

In the simple one-dimensional example in Figure 4.1, finding and checking the different points of the form (i)–(iii) was easy; however there are, of course, examples even in \mathbb{R} which makes this “algorithm” impossible to use, and when considering the multi-dimensional case (that is, $n > 1$) this is a completely absurd “method” for solving an optimization problem.

Remark 4.4 (checking for local optimality is hard) It cannot be over-emphasized that it is hard to check, for a general constrained problem, whether a given feasible solution actually is locally or globally optimal. We refer to one such negative result: Pardalos and Schnitger [PaS88] consider a constrained quadratic optimization problem, that is, where the objective function has the form (3.10), see also (4.6) below; citing from their abstract, *the problem of checking local optimality for a feasible point and the problem of checking if a local minimum is strict, are NP-hard problems. As a consequence the problem of checking whether a function is locally strictly convex is also NP-hard.* Since finding a globally optimal solution to problems in this class is also NP-hard (cf. [PaV91]), it follows that checking for local optimality is, in terms of

worst-case complexity, just as difficult as that to solve the entire problem. ■

In the following we will develop necessary and sufficient conditions for \mathbf{x}^* to be a local or a global optimal solution to the problem (4.1) for any dimension $n \geq 1$, and which are useful and possible to check. Before we do that, however, we will establish when a globally optimal solution to the problem (4.1) exists.

4.2 Existence of optimal solutions

4.2.1 A classic result

We first pave the way for a classic result from calculus: Weierstrass' Theorem.

Definition 4.5 (weakly coercive, coercive functions) *Let $S \subseteq \mathbb{R}^n$ be a nonempty and closed set, and $f : S \rightarrow \mathbb{R}$ be a given function.*

(a) *We say that f is weakly coercive with respect to the set S if S is bounded or for every $N > 0$ there exists an $M > 0$ such that $f(\mathbf{x}) \geq N$ whenever $\|\mathbf{x}\| \geq M$.*

In other words, f is weakly coercive if either S is bounded or

$$\lim_{\substack{\|\mathbf{x}\| \rightarrow \infty \\ \mathbf{x} \in S}} f(\mathbf{x}) = \infty$$

holds.

(b) *We say that f is coercive with respect to the set S if S is bounded or for every $N > 0$ there exists an $M > 0$ such that $f(\mathbf{x})/\|\mathbf{x}\| \geq N$ whenever $\|\mathbf{x}\| \geq M$.*

In other words, f is coercive if either S is bounded or

$$\lim_{\substack{\|\mathbf{x}\| \rightarrow \infty \\ \mathbf{x} \in S}} f(\mathbf{x})/\|\mathbf{x}\| = \infty$$

holds. ■

The weak coercivity of $f : S \rightarrow \mathbb{R}$ is (for nonempty closed sets S) equivalent to the property that f has bounded level sets restricted to S (cf. Definition 3.43). (Why?)

A coercive function grows faster than any linear function. In fact, for convex functions f , f being coercive is equivalent to $\mathbf{x} \mapsto f(\mathbf{x}) - \mathbf{a}^T \mathbf{x}$

being weakly coercive for every vector $\mathbf{a} \in \mathbb{R}^n$. This property is a very useful one for certain analyses in the context of Lagrangian duality.¹

We next introduce two extended notions of continuity.

Definition 4.6 (semi-continuity) Consider a function $f : S \rightarrow \mathbb{R}$, where $S \subseteq \mathbb{R}^n$ is nonempty.

(a) The function f is said to be lower semi-continuous at $\bar{\mathbf{x}} \in S$ if the value $f(\bar{\mathbf{x}})$ is less than or equal to every limit of f as $\mathbf{x}_k \rightarrow \bar{\mathbf{x}}$, that is, $\lim_{k \rightarrow \infty} \mathbf{x}_k = \bar{\mathbf{x}}$.

In other words, f is lower semi-continuous at $\bar{\mathbf{x}} \in S$ if

$$\mathbf{x}_k \rightarrow \bar{\mathbf{x}} \quad \implies \quad f(\bar{\mathbf{x}}) \leq \liminf_{k \rightarrow \infty} f(\mathbf{x}_k).$$

(b) The function f is said to be upper semi-continuous at $\bar{\mathbf{x}} \in S$ if the value $f(\bar{\mathbf{x}})$ is greater than or equal to every limit of f as $\mathbf{x}_k \rightarrow \bar{\mathbf{x}}$.

In other words, f is upper semi-continuous at $\bar{\mathbf{x}} \in S$ if

$$\mathbf{x}_k \rightarrow \bar{\mathbf{x}} \quad \implies \quad f(\bar{\mathbf{x}}) \geq \limsup_{k \rightarrow \infty} f(\mathbf{x}_k).$$

We say that f is lower semi-continuous on S (respectively, upper semi-continuous on S) if it is lower semi-continuous (respectively, upper semi-continuous) at every $\bar{\mathbf{x}} \in S$. ■

Lower semi-continuous functions in one variable have the appearance shown in Figure 4.3.

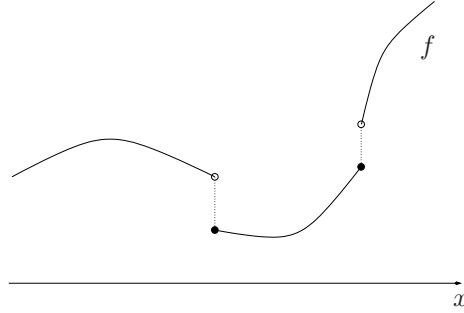


Figure 4.3: A lower semi-continuous function in one variable.

Establish the following important relations:

¹For example, in Section 6.3.2 we suppose that the ground set X is compact in order for the Lagrangian dual function q to be finite. It is possible to replace the boundedness condition on X with a coercivity condition on f .

(a) The function f mentioned in Definition 4.6 is *continuous* at $\bar{\mathbf{x}} \in S$ if and only if it is *both* lower and upper semi-continuous at $\bar{\mathbf{x}}$.

(b) The lower semi-continuity of f is equivalent to the closedness of all its level sets $\text{lev}_f(b)$, $b \in \mathbb{R}$ (cf. Definition 3.43), as well as the closedness of its epigraph (cf. Definition 3.38). (Why?)

Next follows the famous existence theorem credited to Karl Weierstrass (see, however, Section 4.7).

Theorem 4.7 (Weierstrass' Theorem) *Let $S \subseteq \mathbb{R}^n$ be a nonempty and closed set, and $f : S \rightarrow \mathbb{R}$ be a lower semi-continuous function on S . If f is weakly coercive with respect to S , then there exists a nonempty, closed and bounded (thus compact) set of globally optimal solutions to the problem (4.1).*

Proof. We first assume that S is bounded, and proceed by choosing a sequence $\{\mathbf{x}_k\}$ in S such that

$$\lim_{k \rightarrow \infty} f(\mathbf{x}_k) = \inf_{\mathbf{x} \in S} f(\mathbf{x}).$$

(The infimum of f over S is the lowest limit of all sequences of the form $\{f(\mathbf{x}_k)\}$ with $\{\mathbf{x}_k\} \subset S$, so such a sequence of vectors \mathbf{x}_k is what we here are choosing.)

Due to the boundedness of S , the sequence $\{\mathbf{x}_k\}$ must have limit points, all of which lie in S because of the closedness of S . Let $\bar{\mathbf{x}}$ be an arbitrary limit point of $\{\mathbf{x}_k\}$, corresponding to the subsequence $\mathcal{K} \subseteq \mathbb{Z}_+$. Then, by the lower semi-continuity of f ,

$$f(\bar{\mathbf{x}}) \leq \liminf_{k \in \mathcal{K}} f(\mathbf{x}_k) = \lim_{k \in \mathcal{K}} f(\mathbf{x}_k) = \inf_{\mathbf{x} \in S} f(\mathbf{x}).$$

Since $\bar{\mathbf{x}}$ attains the infimum of f over S , $\bar{\mathbf{x}}$ is a global minimum of f over S . This limit point of $\{\mathbf{x}_k\}$ was arbitrarily chosen; any other choice (provided more than one exists) has the same (optimal) objective value.

Suppose next that f is weakly coercive, and consider the same sequence $\{\mathbf{x}_k\}$ in S . Then, by the weak coercivity assumption, either $\{\mathbf{x}_k\}$ is bounded or the elements of the sequence $\{f(\mathbf{x}_k)\}$ tend to infinity. The non-emptiness of S implies that $\inf_{\mathbf{x} \in S} f(\mathbf{x}) < \infty$ holds, and hence we conclude that $\{\mathbf{x}_k\}$ is bounded. We can then utilize the same arguments as in the previous paragraph and conclude that also in this case there exists a globally optimal solution. We are done. ■

Before moving on we take a closer look at the proof of this result, because it is instrumental in order to understand the importance of some of the assumptions that we make about the optimization models that we

pose. We notice that the closedness of S is really crucial; if S is not closed then a sequence generated in S may converge to a point outside of S , which means that we would converge to an infeasible and of course also non-optimal solution. This is the reason why the generic optimization model (1.1) stated in Chapter 1 does not contain any constraints of the form

$$g_i(\mathbf{x}) < 0, \quad i \in \mathcal{SI},$$

where \mathcal{SI} denotes *strict inequality*. The reason is that such constraints in general may describe non-closed sets.

4.2.2 *Non-standard results

Weierstrass' Theorem 4.7 is next improved for certain convex instances of the problem (4.1). The main purpose of presenting these results is to show the role of convexity and to illustrate the special properties of convex quadratic programs and linear programs. The proofs are complex and are left out; see the references in Section 4.7.

Theorem 4.8 (existence of optimal solutions, convex polynomials) *Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a convex polynomial function. Suppose further that the set S can be described by inequality constraints of the form $g_i(\mathbf{x}) \leq 0$, $i = 1, \dots, m$, where each function g_i is convex and polynomial. The problem (4.1) then has a nonempty (as well as closed and convex) set of globally optimal solutions if and only if f is lower bounded on S .* ■

In the following result, we let S be a nonempty polyhedron, and suppose that it is possible to describe it as the following finite (cf. Definition 3.15) set of linear constraints:

$$S := \{ \mathbf{x} \in \mathbb{R}^n \mid \mathbf{Ax} \leq \mathbf{b}; \quad \mathbf{Ex} = \mathbf{d} \}, \quad (4.4)$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{E} \in \mathbb{R}^{\ell \times n}$, $\mathbf{b} \in \mathbb{R}^m$, and $\mathbf{d} \in \mathbb{R}^\ell$. The *recession cone* to S then is the following set, defining the set of directions that are feasible at every point in S :²

$$\text{rec}_S := \{ \mathbf{p} \in \mathbb{R}^n \mid \mathbf{Ap} \leq \mathbf{0}^m; \quad \mathbf{Ep} = \mathbf{0}^\ell \}. \quad (4.5)$$

(For the definition of the set of feasible directions at a given vector \mathbf{x} , see Definition 4.20.)

We also suppose that

$$f(\mathbf{x}) := \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{q}^T \mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^n, \quad (4.6)$$

²Recall the cone C in the Representation Theorem 3.22.

where $\mathbf{Q} \in \mathbb{R}^{n \times n}$ is a symmetric and positive semidefinite matrix and $\mathbf{q} \in \mathbb{R}^n$. We define the recession cone to any convex function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ as follows: the *recession cone* to f is the recession cone to the level set of f (cf. Definition 3.43), defined for any value of b for which the corresponding level set of f is nonempty.³ In the special case of the convex quadratic function given in (4.6),

$$\text{rec}_f = \{ \mathbf{p} \in \mathbb{R}^n \mid \mathbf{Q}\mathbf{p} = \mathbf{0}^n; \quad \mathbf{q}^T \mathbf{p} \leq 0 \}.$$

This is the set of directions that nowhere are ascent directions to f .

Corollary 4.9 (Frank–Wolfe Theorem) *Suppose that S is the polyhedron described by (4.4) and f is the convex quadratic function given by (4.6), so that the problem (4.1) is a convex quadratic programming problem. Then, the following three statements are equivalent.*

- (a) *The problem (4.1) has a nonempty (as well as a closed and convex) set of globally optimal solutions.*
- (b) *f is lower bounded on S .*
- (c) *For every vector \mathbf{p} in the intersection of the recession cone rec_S to S and the null space $N(\mathbf{Q})$ of the matrix \mathbf{Q} , it holds that $\mathbf{q}^T \mathbf{p} \geq 0$. In other words,*

$$\mathbf{p} \in \text{rec}_S \cap N(\mathbf{Q}) \quad \implies \quad \mathbf{q}^T \mathbf{p} \geq 0$$

holds. ■

The statement in (c) shows that the conditions for the existence of an optimal solution in the case of convex quadratic programs are milder than in the general convex case. In the latter case, we can state a slight improvement over the Weierstrass Theorem 4.7 that if, in the problem (4.1), f is convex on S where the latter is nonempty, closed and convex, then the problem has a nonempty, convex and compact set of globally optimal solutions if and only if $\text{rec}_S \cap \text{rec}_f = \{\mathbf{0}^n\}$. The improvements in the above results for polyhedral, in particular quadratic, programs stem from the fact that convex polynomial functions cannot be lower bounded and yet not have a global minimum.

[Note: Consider the special case of the problem (4.1) where $f(x) := 1/x$ and $S := [1, +\infty)$. It is clear that f is bounded from below on S , in fact by the value zero which is the infimum of f over S , but it never attains the value zero on S , and therefore this problem has no optimal solution. Of course, f is not a polynomial function.]

³Check that this cone actually is independent of the value of b under this only requirement. Also confirm that if the level set $\text{lev}_f(b)$ is (nonempty and) bounded for some $b \in \mathbb{R}$ then it is bounded for every $b \in \mathbb{R}$, thanks to the convexity of f .

Corollary 4.10 (a fundamental theorem in linear programming) *Suppose, in the Frank–Wolfe Theorem, that f is linear, that is, that $\mathbf{Q} = \mathbf{0}^{n \times n}$. Then, the problem (4.1) is identical to a linear programming (LP) problem. Then, the following three statements are equivalent.*

(a) *The problem (4.1) has a nonempty (as well as a closed and convex polyhedral) set of globally optimal solutions.*

(b) *f is lower bounded on S .*

(c) *For every vector \mathbf{p} in the recession cone rec_S to S , it holds that $\mathbf{q}^T \mathbf{p} \geq 0$. In other words,*

$$\mathbf{p} \in \text{rec}_S \quad \implies \quad \mathbf{q}^T \mathbf{p} \geq 0$$

holds. ■

Corollary 4.10 will in fact be established later on in Theorem 8.10, by the use of polyhedral convexity, when we specialize our treatment of non-linear optimization to that of linear optimization. Since we have already established the Representation Theorem 3.22, proving Corollary 4.10 for the case of LP will be easy: since the objective function is linear, every feasible direction $\mathbf{p} \in \text{rec}_S$ with $\mathbf{q}^T \mathbf{p} < 0$ leads to an unbounded solution from any vector $\mathbf{x} \in S$.

4.2.3 Special optimal solution sets

Under strict convexity, we can establish the following result.

Proposition 4.11 (unique optimal solution under strict convexity) *Suppose that in the problem (4.1) f is strictly convex on S and the set S is convex. Then, there can be at most one globally optimal solution.*

Proof. Suppose, by means of contradiction, that \mathbf{x}^* and \mathbf{x}^{**} are two different globally optimal solutions. Then, for every $\lambda \in (0, 1)$, we have that

$$f(\lambda \mathbf{x}^* + (1 - \lambda) \mathbf{x}^{**}) < \lambda f(\mathbf{x}^*) + (1 - \lambda) f(\mathbf{x}^{**}) = f(\mathbf{x}^*) [= f(\mathbf{x}^{**})].$$

Since $\lambda \mathbf{x}^* + (1 - \lambda) \mathbf{x}^{**} \in S$, we have found an entire interval of points which are strictly better than \mathbf{x}^* or \mathbf{x}^{**} . This is impossible. ■

We finally characterize a class of optimization problems over polytopes whose optimal solution set, if nonempty, includes an extreme point.

Consider the maximization problem to

$$\begin{aligned} & \text{maximize} && f(\mathbf{x}), \\ & \text{subject to} && \mathbf{x} \in P, \end{aligned} \tag{4.7}$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a convex function and $P \subset \mathbb{R}^n$ is a nonempty, bounded polyhedron (that is, a polytope). Then, from the Representation Theorem 3.22 it follows below that an optimal solution can be found among the extreme points of P . Theorem 8.10 establishes a corresponding result for linear programs that does not rely on Weierstrass' Theorem.

Theorem 4.12 (optimal extreme point) *An optimal solution to (4.7) can be found among the extreme points of P .*

Proof. The function f is continuous (since it is convex, cf. Theorem 4.27 below); further, P is a nonempty and compact set. Hence, there exists an optimal solution $\tilde{\mathbf{x}}$ to (4.7) by Weierstrass' Theorem 4.7. The Representation Theorem 3.22 implies that $\tilde{\mathbf{x}} = \lambda_1 \mathbf{v}^1 + \cdots + \lambda_k \mathbf{v}^k$ for some extreme points $\mathbf{v}^1, \dots, \mathbf{v}^k$ of P and $\lambda_1, \dots, \lambda_k \geq 0$ such that $\sum_{i=1}^k \lambda_i = 1$. But then (from the convexity of f)

$$\begin{aligned} f(\tilde{\mathbf{x}}) &= f(\lambda_1 \mathbf{v}^1 + \cdots + \lambda_k \mathbf{v}^k) \leq \lambda_1 f(\mathbf{v}^1) + \cdots + \lambda_k f(\mathbf{v}^k) \\ &\leq \lambda_1 f(\tilde{\mathbf{x}}) + \cdots + \lambda_k f(\tilde{\mathbf{x}}) = f(\tilde{\mathbf{x}}), \end{aligned}$$

which gives that $f(\tilde{\mathbf{x}}) = f(\mathbf{v}^i)$ for some $i = 1, \dots, k$. ■

Remark 4.13 Every linear function is convex, so Theorem 4.12 implies, in particular, that every linear program over a nonempty and bounded polyhedron has an optimal extreme point. ■

4.3 Optimality in unconstrained optimization

In Theorem 4.3 we have established that locally optimal solutions also are global in the convex case. What are the necessary and sufficient conditions for a vector \mathbf{x}^* to be a local optimum? This is an important question, because the algorithms that we will investigate for solving important classes of optimization problems are always devised based on those conditions that we would like to fulfill. This is a statement that seems to be true universally: *efficient, locally or globally convergent iterative algorithms for an optimization problem are directly based on its necessary and/or sufficient local optimality conditions.*

We begin by establishing these conditions for the case of unconstrained optimization, where the objective function is in C^1 . Every proof is based on the Taylor expansion of the objective function up to order one or two.

Our problem here is the following:

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \ f(\mathbf{x}), \quad (4.8)$$

where f is in C^1 on \mathbb{R}^n [for short we say: in C^1 or $C^1(\mathbb{R}^n)$].

Theorem 4.14 (necessary optimality conditions, C^1 case) *Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is in C^1 on \mathbb{R}^n . Then,*

$$\mathbf{x}^* \text{ is a local minimum of } f \text{ over } \mathbb{R}^n \implies \nabla f(\mathbf{x}^*) = \mathbf{0}^n.$$

Note that

$$\nabla f(\mathbf{x}) = \left(\frac{\partial f(\mathbf{x})}{\partial x_j} \right)_{j=1}^n,$$

so the requirement thus is that $\frac{\partial f(\mathbf{x}^*)}{\partial x_j} = 0$, $j = 1, \dots, n$.

Just as for the case $n = 1$, we refer to this condition as \mathbf{x}^* being a *stationary point* of f .

[Note: For $n = 1$, Theorem 4.14 reduces to: $x^* \in \mathbb{R}$ is a *local minimum* $\implies f'(x^*) = 0$.]

Proof. (By contradiction.) Suppose that \mathbf{x}^* is a local minimum, but that $\nabla f(\mathbf{x}^*) \neq \mathbf{0}^n$. Let $\mathbf{p} := -\nabla f(\mathbf{x}^*)$, and study the Taylor expansion around $\mathbf{x} = \mathbf{x}^*$ in the direction of \mathbf{p} :

$$f(\mathbf{x}^* + \alpha \mathbf{p}) = f(\mathbf{x}^*) + \alpha \nabla f(\mathbf{x}^*)^T \mathbf{p} + o(\alpha),$$

where $o : \mathbb{R} \rightarrow \mathbb{R}$ is such that $o(s)/s \rightarrow 0$ when $s \rightarrow 0$. We get that

$$\begin{aligned} f(\mathbf{x}^* + \alpha \mathbf{p}) &= f(\mathbf{x}^*) - \alpha \|\nabla f(\mathbf{x}^*)\|^2 + o(\alpha) \\ &< f(\mathbf{x}^*) \end{aligned}$$

for all small enough $\alpha > 0$, since $\|\nabla f(\mathbf{x}^*)\| \neq 0$. This completes the proof. ■

The opposite direction is false: take $f(x) = x^3$; then, $\bar{x} = 0$ is stationary, but it is neither a local minimum nor a local maximum.

The proof is instrumental in that it provides a sufficient condition for a vector \mathbf{p} to define a *descent direction*, that is, a direction such that a small step along it yields a lower objective value. We first define this notion properly.

Definition 4.15 (descent direction) *Let the function $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ be given. Let $\mathbf{x} \in \mathbb{R}^n$ be a vector such that $f(\mathbf{x})$ is finite. Let $\mathbf{p} \in \mathbb{R}^n$. We say that the vector $\mathbf{p} \in \mathbb{R}^n$ is a descent direction with respect to f at \mathbf{x} if*

$$\exists \delta > 0 \text{ such that } f(\mathbf{x} + \alpha \mathbf{p}) < f(\mathbf{x}) \text{ for every } \alpha \in (0, \delta]$$

holds. ■

Proposition 4.16 (sufficient condition for descent) *Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is in C^1 around a point \mathbf{x} for which $f(\mathbf{x}) < +\infty$, and that $\mathbf{p} \in \mathbb{R}^n$. Then,*

$$\nabla f(\mathbf{x})^T \mathbf{p} < 0 \implies \mathbf{p} \text{ is a descent direction with respect to } f \text{ at } \mathbf{x}$$

holds.

Proof. Since f is in C^1 around \mathbf{x} , we can construct a Taylor expansion of f , as above:

$$f(\mathbf{x} + \alpha \mathbf{p}) = f(\mathbf{x}) + \alpha \nabla f(\mathbf{x})^T \mathbf{p} + o(\alpha).$$

Since $\nabla f(\mathbf{x})^T \mathbf{p} < 0$, we obtain that $f(\mathbf{x} + \alpha \mathbf{p}) < f(\mathbf{x})$ for all sufficiently small values of $\alpha > 0$. ■

Notice that at a point $\mathbf{x} \in \mathbb{R}^n$ there may be other descent directions $\mathbf{p} \in \mathbb{R}^n$ beside those satisfying that $\nabla f(\mathbf{x})^T \mathbf{p} < 0$; in Example 11.2(b) we show how directions of *negative curvature* stemming from eigenvectors corresponding to negative eigenvalues of the Hessian matrix $\nabla^2 f(\mathbf{x})$ can be utilized.

If f in addition is convex then the opposite implication in the above proposition is true, thus making the descent property equivalent to the property that the directional derivative is negative. Since this result can be stated also for *non-differentiable* functions f (in which case we must of course replace the expression “ $\nabla f(\mathbf{x})^T \mathbf{p}$ ” with the classic expression for the directional derivative, $f'(\mathbf{x}; \mathbf{p}) := \lim_{\alpha \rightarrow 0+} \frac{1}{\alpha} [f(\mathbf{x} + \alpha \mathbf{p}) - f(\mathbf{x})]$), we shall relegate the proof of this equivalence to our presentation of the subdifferentiability analysis of convex functions in Section 6.3.1, in particular to Proposition 6.18.

If f has stronger differentiability properties, then we can say even more what a local optimum must be like.

Theorem 4.17 (necessary optimality conditions, C^2 case) *Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is in C^2 on \mathbb{R}^n . Then,*

$$\mathbf{x}^* \text{ is a local minimum of } f \implies \begin{cases} \nabla f(\mathbf{x}^*) = \mathbf{0}^n \\ \nabla^2 f(\mathbf{x}^*) \text{ is positive semidefinite.} \end{cases}$$

[Note: For $n = 1$, Theorem 4.17 reduces to: $x^* \in \mathbb{R}$ is a local minimum of f over $\mathbb{R} \implies f'(x^*) = 0$ and $f''(x^*) \geq 0$.]

Proof. Consider the Taylor expansion of f up to order two around \mathbf{x}^* and in the direction of a vector $\mathbf{p} \in \mathbb{R}^n$:

$$f(\mathbf{x}^* + \alpha \mathbf{p}) = f(\mathbf{x}^*) + \alpha \nabla f(\mathbf{x}^*)^T \mathbf{p} + \frac{\alpha^2}{2} \mathbf{p}^T \nabla^2 f(\mathbf{x}^*) \mathbf{p} + o(\alpha^2).$$

Suppose that \mathbf{x}^* satisfies $\nabla f(\mathbf{x}^*) = \mathbf{0}^n$, but that there is a vector $\mathbf{p} \neq \mathbf{0}^n$ with $\mathbf{p}^T \nabla^2 f(\mathbf{x}^*) \mathbf{p} < 0$, that is, $\nabla^2 f(\mathbf{x}^*)$ is not positive semidefinite. Then the above yields that $f(\mathbf{x}^* + \alpha \mathbf{p}) < f(\mathbf{x}^*)$ for all small enough $\alpha > 0$, whence \mathbf{x}^* cannot be a local minimum. ■

Also in this case, the opposite direction is false; the same counterexample as that after Theorem 4.14 applies.

In Example 11.2(b) we provide an example of a descent direction that has the form provided in the above proof; it is based on \mathbf{p} being an eigenvector corresponding to a negative eigenvalue of $\nabla^2 f(\mathbf{x}^*)$.

The next result shows that under some circumstances, we can establish local optimality of a stationary point.

Theorem 4.18 (sufficient optimality conditions, C^2 case) *Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is in C^2 on \mathbb{R}^n . Then,*

$$\left. \begin{array}{l} \nabla f(\mathbf{x}^*) = \mathbf{0}^n \\ \nabla^2 f(\mathbf{x}^*) \text{ is positive definite} \end{array} \right\} \implies \mathbf{x}^* \text{ strict local min of } f \text{ over } \mathbb{R}^n.$$

[Note: For $n = 1$, Theorem 4.18 reduces to: $f'(x^*) = 0$ and $f''(x^*) > 0 \implies x^* \in \mathbb{R}$ is a strict local minimum of f over \mathbb{R} .]

Proof. Suppose that $\nabla f(\mathbf{x}^*) = \mathbf{0}^n$ and $\nabla^2 f(\mathbf{x}^*)$ is positive definite. Take an arbitrary vector $\mathbf{p} \in \mathbb{R}^n$, $\mathbf{p} \neq \mathbf{0}^n$. Then,

$$\begin{aligned} f(\mathbf{x}^* + \alpha \mathbf{p}) &= f(\mathbf{x}^*) + \alpha \underbrace{\nabla f(\mathbf{x}^*)^T \mathbf{p}}_{=0} + \frac{\alpha^2}{2} \underbrace{\mathbf{p}^T \nabla^2 f(\mathbf{x}^*) \mathbf{p}}_{>0} + o(\alpha^2) \\ &> f(\mathbf{x}^*) \end{aligned}$$

for all small enough $\alpha > 0$. As \mathbf{p} was arbitrary, the above implies that \mathbf{x}^* is a strict local minimum of f over \mathbb{R}^n . ■

We naturally face the following question: When is a stationary point a global minimum? The answer is given next. (Investigate the connection between this result and the Fundamental Theorem 4.3.)

Theorem 4.19 (necessary and sufficient global optimality conditions) Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex and in C^1 on \mathbb{R}^n . Then,

$$\mathbf{x}^* \text{ is a global minimum of } f \text{ over } \mathbb{R}^n \iff \nabla f(\mathbf{x}^*) = \mathbf{0}^n.$$

Proof. $[\implies]$ This has already been shown in Theorem 4.14, since a global minimum is a local minimum.

$[\impliedby]$ The convexity of f yields that for every $\mathbf{y} \in \mathbb{R}^n$,

$$f(\mathbf{y}) \geq f(\mathbf{x}^*) + \nabla f(\mathbf{x}^*)^T (\mathbf{y} - \mathbf{x}^*) = f(\mathbf{x}^*),$$

where the equality stems from the property that $\nabla f(\mathbf{x}^*) = \mathbf{0}^n$. ■

4.4 Optimality for optimization over convex sets

We consider a quite general optimization problem of the form:

$$\text{minimize } f(\mathbf{x}), \tag{4.9a}$$

$$\text{subject to } \mathbf{x} \in S, \tag{4.9b}$$

where $S \subseteq \mathbb{R}^n$ is nonempty, closed and convex, and $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is in C^1 on S .

A noticeable difference to unconstrained optimization is the fact that whether a vector $\mathbf{p} \in \mathbb{R}^n$ can be used as a direction of movement from a point $\mathbf{x} \in S$ depends on the constraints defining S ; if \mathbf{x} is an interior point of S then every $\mathbf{p} \in \mathbb{R}^n$ is a *feasible direction*, otherwise only certain directions will be feasible. That is, it all depends on whether there are any *active constraints* of S at \mathbf{x} or not. We will define these terms in detail next, and then develop necessary and sufficient optimality conditions based on them. The optimality conditions are natural extensions of those for the case of unconstrained optimization, and reduce to them when $S = \mathbb{R}^n$. Further, we will develop a way of measuring the distance to an optimal solution in terms of the value of the objective function f , which is valid for convex problems. As a result of this development, we will also be able to finally establish the Separation Theorem 3.24, whose proof has been postponed until now. (See Section 4.6.2 for the proof.)

Definition 4.20 (feasible direction) Suppose that $\mathbf{x} \in S$, where $S \subseteq \mathbb{R}^n$, and that $\mathbf{p} \in \mathbb{R}^n$. Then, the vector \mathbf{p} defines a feasible direction at \mathbf{x} if a small step in the direction of \mathbf{p} does not lead outside of the set S .

In other words, the vector \mathbf{p} defines a feasible direction at $\mathbf{x} \in S$ if

$$\exists \delta > 0 \text{ such that } \mathbf{x} + \alpha \mathbf{p} \in S \text{ for all } \alpha \in [0, \delta]$$

holds. ■

Recall that in the discussion following Theorem 4.8 we defined the set of feasible directions of a polyhedral set, that is, the set of directions that are feasible at every feasible point. For a general set S it would hence be the set

$$\{ \mathbf{p} \in \mathbb{R}^n \mid \forall \mathbf{x} \in S \exists \delta > 0 \text{ such that } \mathbf{x} + \alpha \mathbf{p} \in S \text{ for all } \alpha \in [0, \delta] \}.$$

For nonempty, closed and convex sets S , this set is nonempty if and only if the set S also is unbounded. (Why?)

Definition 4.21 (active constraints) Suppose that the set $S \subset \mathbb{R}^n$ is defined by a finite collection of equality and inequality constraints:

$$S = \{ \mathbf{x} \in \mathbb{R}^n \mid g_i(\mathbf{x}) = 0, \quad i \in \mathcal{E}; \quad g_i(\mathbf{x}) \leq 0, \quad i \in \mathcal{I} \},$$

where $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$ ($i \in \mathcal{E} \cup \mathcal{I}$) are given functions. Suppose that $\mathbf{x} \in S$. The set of active constraints at \mathbf{x} is the union of the equality constraints and the set of inequality constraints that are satisfied with equality at \mathbf{x} , that is, the set $\mathcal{E} \cup \mathcal{I}(\mathbf{x})$, where $\mathcal{I}(\mathbf{x}) := \{ i \in \mathcal{I} \mid g_i(\mathbf{x}) = 0 \}$. ■

Example 4.22 (feasible directions for linear constraints) Suppose, as a special case, that the constraints are all linear, that is, that for every $i \in \mathcal{E}$, $g_i(\mathbf{x}) := \mathbf{e}_i^T \mathbf{x} - d_i$ ($\mathbf{e}_i \in \mathbb{R}^n$; $d_i \in \mathbb{R}$), and for every $i \in \mathcal{I}$, $g_i(\mathbf{x}) := \mathbf{a}_i^T \mathbf{x} - b_i$ ($\mathbf{a}_i \in \mathbb{R}^n$; $b_i \in \mathbb{R}$). In other words, in matrix notation, $S := \{ \mathbf{x} \in \mathbb{R}^n \mid \mathbf{E}\mathbf{x} = \mathbf{d}; \quad \mathbf{A}\mathbf{x} \leq \mathbf{b} \}$.

Suppose further that $\mathbf{x} \in S$. Then, the set of feasible directions at \mathbf{x} is the set

$$\{ \mathbf{p} \in \mathbb{R}^n \mid \mathbf{E}\mathbf{p} = \mathbf{0}^\ell; \quad \mathbf{a}_i^T \mathbf{p} \leq 0, \quad i \in \mathcal{I}(\mathbf{x}) \}.$$

Just as S , this is a polyhedron. Moreover, it is a polyhedral cone. ■

Clearly, the set of feasible directions of the polyhedral set S (or, the recession cone of S) is

$$\text{rec}_S := \{ \mathbf{p} \in \mathbb{R}^n \mid \mathbf{E}\mathbf{p} = \mathbf{0}^\ell; \quad \mathbf{A}\mathbf{p} \leq \mathbf{0}^m \},$$

as stated in (4.5). Note moreover that the above set rec_S represents the cone C in the Representation Theorem 3.22.⁴

⁴While that theorem was stated for sets defined only by linear inequalities, we can always rewrite the equalities $\mathbf{E}\mathbf{x} = \mathbf{d}$ as $\mathbf{E}\mathbf{x} \leq \mathbf{d}$, $-\mathbf{E}\mathbf{x} \leq -\mathbf{d}$; the corresponding feasible directions are then given by $\mathbf{E}\mathbf{p} \leq \mathbf{0}^\ell$, $-\mathbf{E}\mathbf{p} \leq \mathbf{0}^\ell$, that is, $\mathbf{E}\mathbf{p} = \mathbf{0}^\ell$.

We can now more or less repeat the arguments for the unconstrained case in order to establish a necessary optimality condition for constrained optimization problems over convex sets. This condition will later on in Chapter 5 be given a general statement for the case of explicit constraints in the form of the famous Karush–Kuhn–Tucker conditions in nonlinear programming.

Proposition 4.23 (necessary optimality conditions, C^1 case) *Suppose that $S \subseteq \mathbb{R}^n$ and that $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is in C^1 around a point $\mathbf{x} \in S$ for which $f(\mathbf{x}) < +\infty$.*

(a) *If $\mathbf{x}^* \in S$ is a local minimum of f over S then $\nabla f(\mathbf{x}^*)^T \mathbf{p} \geq 0$ holds for every feasible direction \mathbf{p} at \mathbf{x}^* .*

(b) *Suppose that S is convex and that f is in C^1 on S . If $\mathbf{x}^* \in S$ is a local minimum of f over S then*

$$\nabla f(\mathbf{x}^*)^T (\mathbf{x} - \mathbf{x}^*) \geq 0, \quad \mathbf{x} \in S \quad (4.10)$$

holds.

Proof. (a) We again utilize the Taylor expansion of f around \mathbf{x}^* :

$$f(\mathbf{x}^* + \alpha \mathbf{p}) = f(\mathbf{x}^*) + \alpha \nabla f(\mathbf{x}^*)^T \mathbf{p} + o(\alpha).$$

The proof is by contradiction. As was shown in Proposition 4.16, if there is a direction \mathbf{p} for which it holds that $\nabla f(\mathbf{x}^*)^T \mathbf{p} < 0$, then $f(\mathbf{x}^* + \alpha \mathbf{p}) < f(\mathbf{x}^*)$ for all sufficiently small values of $\alpha > 0$. It suffices here to state that \mathbf{p} should also be a feasible direction in order to reach a contradiction to the local optimality of \mathbf{x}^* .

(b) If S is convex then every feasible direction \mathbf{p} can be written as a positive scalar times the vector $\mathbf{x} - \mathbf{x}^*$ for *some* vector $\mathbf{x} \in S$. (Why?) The expression (4.10) then follows from the statement in (a). ■

The inequality (4.10) is sometimes referred to as a *variational inequality*. We will utilize it for several purposes: (i) to derive equivalent optimality conditions involving a linear optimization problem as well as the Euclidean projection operation Proj_S introduced in Section 3.4; (ii) to derive descent algorithms for the problem (4.9) in Section 12.2 and 12.4; (iii) to derive a near-optimality condition for convex optimization problems in Section 4.5; and (iv) we will extend it to non-convex sets in the form of the Karush–Kuhn–Tucker conditions in Chapter 5.

In Theorem 4.14 we established that for unconstrained C^1 optimization the necessary optimality condition is that $\nabla f(\mathbf{x}^*) = \mathbf{0}^n$ holds. Notice that that is exactly what becomes of the variational inequality (4.10) when $S = \mathbb{R}^n$, because the only way in which that inequality can hold

for every $\mathbf{x} \in \mathbb{R}^n$ is that $\nabla f(\mathbf{x}^*) = \mathbf{0}^n$ holds. Just as we did in the case of unconstrained optimization, we will call a vector $\mathbf{x}^* \in S$ satisfying (4.10) a *stationary point*.

We will next provide two statements equivalent to the variational inequality (4.10). First up, though, we will provide the extension to Theorem 4.19 to the convex constrained case. Notice the resemblance of their respective proofs.

Theorem 4.24 (necessary and sufficient global optimality conditions) *Suppose that $S \subseteq \mathbb{R}^n$ is nonempty and convex. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be convex and in C^1 on S . Then,*

$$\mathbf{x}^* \text{ is a global minimum of } f \text{ over } S \iff (4.10) \text{ holds.}$$

Proof. $[\implies]$ This has already been shown in Proposition 4.23(b), since a global minimum is a local minimum.

$[\impliedby]$ The convexity of f yields [cf. Theorem 3.40(a)] that for every $\mathbf{y} \in S$,

$$f(\mathbf{y}) \geq f(\mathbf{x}^*) + \nabla f(\mathbf{x}^*)^\top (\mathbf{y} - \mathbf{x}^*) \geq f(\mathbf{x}^*),$$

where the second inequality stems from (4.10). ■

First, we will provide the connection to the Euclidean projection of a vector onto a convex set, discussed in Section 3.4. We claim that the property (4.10) is equivalent to

$$\mathbf{x}^* = \text{Proj}_S[\mathbf{x}^* - \nabla f(\mathbf{x}^*)], \quad (4.11)$$

or, more generally, $\mathbf{x}^* = \text{Proj}_S[\mathbf{x}^* - \alpha \nabla f(\mathbf{x}^*)]$ for every $\alpha > 0$. In other words, a point is stationary if and only if a step in the direction of the steepest descent followed by a Euclidean projection onto S means that we have not moved at all. To prove this fact, we will utilize Proposition 4.23(b) for the optimization problem corresponding to this projection. We are interested in finding the point $\mathbf{x} \in S$ that minimizes the distance to the vector $\mathbf{z} := \mathbf{x}^* - \nabla f(\mathbf{x}^*)$. We can write this as a strictly convex optimization problem as follows:

$$\underset{\mathbf{x} \in S}{\text{minimize}} \quad h(\mathbf{x}) := \frac{1}{2} \|\mathbf{x} - \mathbf{z}\|^2. \quad (4.12)$$

The necessary optimality conditions for this problem, as stated in Proposition 4.23(b), is that

$$\nabla h(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) \geq 0, \quad \mathbf{y} \in S, \quad (4.13)$$

holds. Here, $\nabla h(\mathbf{x}) = \mathbf{x} - \mathbf{z} = \mathbf{x} - [\mathbf{x}^* - \nabla f(\mathbf{x}^*)]$. Since h is strictly convex, by Theorem 4.24 we know that the variational inequality (4.13) characterizes \mathbf{x} as the unique globally optimal solution to the projection problem. We claimed that $\mathbf{x} = \mathbf{x}^*$ is the solution to this problem if and only if \mathbf{x}^* is stationary in the problem (4.9). But this follows immediately, since the variational inequality (4.13), for the special choice of h and $\mathbf{x} = \mathbf{x}^*$, becomes

$$\nabla f(\mathbf{x}^*)^T(\mathbf{y} - \mathbf{x}^*) \geq 0, \quad \mathbf{y} \in S,$$

that is, a statement identical to (4.10). The characterization (4.11) is interesting in that it states that if \mathbf{x}^* is *not* stationary, then the projection operation defined therein must provide a step away from \mathbf{x}^* ; this step will in fact yield a reduced value of f under some additional conditions on the step length α , and so it defines a descent algorithm for (4.9); see Exercise 4.5, and the text in Section 12.4.

So far, we have two equivalent characterizations of a stationary point of f at \mathbf{x}^* : (4.10) and (4.11). The following one is based on a linear optimization problem.

Notice that (4.10) states that $\nabla f(\mathbf{x}^*)^T \mathbf{x} \geq \nabla f(\mathbf{x}^*)^T \mathbf{x}^*$ for every $\mathbf{x} \in S$. Since we obtain an equality by setting $\mathbf{x} = \mathbf{x}^*$ we see that \mathbf{x}^* in fact is a globally optimal solution to the problem to

$$\underset{\mathbf{x} \in S}{\text{minimize}} \nabla f(\mathbf{x}^*)^T \mathbf{x}.$$

In other words, (4.10) is equivalent to the statement

$$\underset{\mathbf{x} \in S}{\text{minimum}} \nabla f(\mathbf{x}^*)^T(\mathbf{x} - \mathbf{x}^*) = 0. \quad (4.14)$$

It is quite obvious that if at some point $\mathbf{x} \in S$,

$$\underset{\mathbf{y} \in S}{\text{minimum}} \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) < 0,$$

then every direction of the form $\mathbf{p} := \bar{\mathbf{y}} - \mathbf{x}$, with

$$\bar{\mathbf{y}} \in \arg \underset{\mathbf{y} \in S}{\text{minimum}} \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}),$$

is a feasible descent direction with respect to f at \mathbf{x} . Again, we have a building block of a descent algorithm for the problem (4.9). [The algorithms that immediately spring out from this characterization are called the *Frank–Wolfe* and *Simplicial decomposition* algorithms, when S is polyhedral; we notice that in the polyhedral case, the linear minimization problem in (4.14) is an LP problem. Read more about these

algorithms in Sections 12.2 and 12.3.] Now having got three equivalent stationarity conditions—(4.10), (4.11), and (4.14)—, we finally provide a fourth one. This one is intimately associated with the projection operation, and it introduces an important geometric concept into the theory of optimality, namely the *normal cone* to a (convex) set S .

We studied a particular choice of \mathbf{z} above, but let us consider an extension of Figure 3.15 which provided an image of the Euclidean projection.

Notice from the above arguments that if we wish to project the vector $\mathbf{z} \in \mathbb{R}^n$ onto S , then the resulting (unique) projection is the vector \mathbf{x} for which the following holds:

$$[\mathbf{x} - \mathbf{z}]^T(\mathbf{y} - \mathbf{x}) \geq 0, \quad \mathbf{y} \in S.$$

Changing sign for clarity, this is the same as

$$[\mathbf{z} - \mathbf{x}]^T(\mathbf{y} - \mathbf{x}) \leq 0, \quad \mathbf{y} \in S.$$

The interpretation of this inequality is that the angle between the two vectors $\mathbf{z} - \mathbf{x}$ (the vector that points towards the point being projected) and the vector $\mathbf{y} - \mathbf{x}$ (the vector that points towards any vector $\mathbf{y} \in S$) is $\geq 90^\circ$. So, the projection operation has the characterization

$$[\mathbf{z} - \text{Proj}_S(\mathbf{z})]^T(\mathbf{y} - \text{Proj}_S(\mathbf{z})) \leq 0, \quad \mathbf{y} \in S. \quad (4.15)$$

The above is shown in Figure 4.4 for $\mathbf{x} = \mathbf{x}^*$ and $\mathbf{z} = \mathbf{x}^* - \nabla f(\mathbf{x}^*)$.

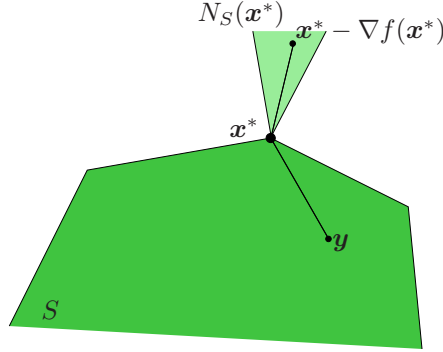


Figure 4.4: Normal cone characterization of a stationary point.

Here, the point being projected is $\mathbf{z} = \mathbf{x}^* - \nabla f(\mathbf{x}^*)$, as used in the characterization of stationarity.

What is left to complete the picture is to define the normal cone, $N_S(\mathbf{x}^*)$, to S at \mathbf{x}^* , depicted in Figure 4.4 in the lighter shade.

Definition 4.25 (normal cone) *Suppose that the set $S \subseteq \mathbb{R}^n$ is closed and convex. Let $\mathbf{x} \in \mathbb{R}^n$. Then, the normal cone to S at \mathbf{x} is the set*

$$N_S(\mathbf{x}) := \begin{cases} \{ \mathbf{v} \in \mathbb{R}^n \mid \mathbf{v}^T(\mathbf{y} - \mathbf{x}) \leq 0, & \mathbf{y} \in S \}, & \text{if } \mathbf{x} \in S, \\ \emptyset, & \text{otherwise.} \end{cases} \quad (4.16)$$

■

According to the definition, we can now define our fourth characterization of a stationary point at \mathbf{x}^* as follows:

$$-\nabla f(\mathbf{x}^*) \in N_S(\mathbf{x}^*). \quad (4.17)$$

What this condition states geometrically is that the angle between the negative gradient and any feasible direction is $\geq 90^\circ$, which, of course, whenever $\nabla f(\mathbf{x}^*) \neq \mathbf{0}^n$, is the same as stating that at \mathbf{x}^* there exist no feasible descent directions. The four conditions (4.10), (4.11), (4.14), and (4.17) are equivalent, and so according to Theorem 4.24 they all are also both necessary and sufficient for the global optimality of \mathbf{x}^* as soon as f is convex.

We remark that in the special case when S is an affine subspace (such as the solution set of a number of linear equations, $S := \{ \mathbf{x} \in \mathbb{R}^n \mid \mathbf{E}\mathbf{x} = \mathbf{d} \}$), the statement (4.17) means that at a stationary point \mathbf{x}^* , $\nabla f(\mathbf{x}^*)$ is parallel to a normal of the subspace.

The normal cone inclusion (4.17) will later be extended to more general sets, where S is described by a finite collection of possibly non-convex constraints. The extension will lead us to the famous Karush–Kuhn–Tucker conditions in Chapter 5. [It turns out to be much more convenient to extend (4.17) than the other three characterizations of stationarity.]

We finish this section by proving a proposition on the behaviour of the gradient of the objective function f on the solution set S^* to convex problems of the form (4.1). The below result shows that ∇f enjoys a stability property, and it also extends the result from the unconstrained case where the value of ∇f always is zero on the solution set.

Proposition 4.26 (invariance of ∇f on the solution set of convex programs) *Suppose that $S \subseteq \mathbb{R}^n$ is convex and that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex and in C^1 on S . Then, the value of $\nabla f(\mathbf{x})$ is constant on the optimal solution set S^* .*

Further, suppose that $\mathbf{x}^ \in S^*$. Then,*

$$S^* = \{ \mathbf{x} \in S \mid \nabla f(\mathbf{x}^*)^T(\mathbf{x} - \mathbf{x}^*) = 0 \text{ and } \nabla f(\mathbf{x}) = \nabla f(\mathbf{x}^*) \}.$$

Proof. Let $\mathbf{x}^* \in S^*$. The definition of the convexity of f shows that

$$f(\mathbf{x}) - f(\mathbf{x}^*) \geq \nabla f(\mathbf{x}^*)^\top (\mathbf{x} - \mathbf{x}^*), \quad \mathbf{x} \in \mathbb{R}^n. \quad (4.18)$$

Let $\bar{\mathbf{x}} \in S^*$. Then, it follows that $\nabla f(\mathbf{x}^*)^\top (\bar{\mathbf{x}} - \mathbf{x}^*) = 0$. By substituting $\nabla f(\mathbf{x}^*)^\top \mathbf{x}^*$ with $\nabla f(\mathbf{x}^*)^\top \bar{\mathbf{x}}$ in (4.18) and using that $f(\mathbf{x}^*) = f(\bar{\mathbf{x}})$, we obtain that

$$f(\mathbf{x}) - f(\bar{\mathbf{x}}) \geq \nabla f(\mathbf{x}^*)^\top (\mathbf{x} - \bar{\mathbf{x}}), \quad \mathbf{x} \in \mathbb{R}^n,$$

which is equivalent to the statement that $\nabla f(\bar{\mathbf{x}}) = \nabla f(\mathbf{x}^*)$. \blacksquare

4.5 Near-optimality in convex optimization

We will here utilize Theorem 4.24 in order to provide a measure of the distance to the optimal solution in terms of the value of f at any feasible point \mathbf{x} .

Let $\mathbf{x} \in S$, and suppose that f is convex on S . Suppose also that $\mathbf{x}^* \in S$ is an arbitrary globally optimal solution, which we suppose exists. From the necessary and sufficient optimality conditions stated in Theorem 4.24, it is clear that unless \mathbf{x} solves (4.9) there exists a $\mathbf{y} \in S$ such that $\nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) < 0$, and hence $\mathbf{p} := \mathbf{y} - \mathbf{x}$ is a feasible descent direction.

Suppose now that

$$\bar{\mathbf{y}} \in \arg \min_{\mathbf{y} \in S} z(\mathbf{y}) := \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}). \quad (4.19)$$

Consider the following string of inequalities and equalities:

$$f(\mathbf{x}) + z(\bar{\mathbf{y}}) = f(\mathbf{x}) + \min_{\mathbf{y} \in S} z(\mathbf{y}) \leq f(\mathbf{x}) + z(\mathbf{x}^*) \leq f(\mathbf{x}^*) \leq f(\mathbf{x}).$$

The equality follows by definition; the first inequality stems from the fact that $\bar{\mathbf{y}}$ solves the linear minimization problem, while the vector \mathbf{x}^* may not; the second inequality follows from the convexity of f on S [cf. Theorem 3.40(a)]; the final inequality follows from the global optimality of \mathbf{x}^* and the feasibility of \mathbf{x} .

From the above, we obtain a closed interval wherein we know that the optimal value of the problem (4.9) lies. Let $f^* := \min_{\mathbf{x} \in S} f(\mathbf{x}) = f(\mathbf{x}^*)$. Then, for every $\mathbf{x} \in S$ and $\bar{\mathbf{y}} \in \arg \min_{\mathbf{y} \in S} z(\mathbf{y})$,

$$f^* \in [f(\mathbf{x}) + z(\bar{\mathbf{y}}), f(\mathbf{x})]. \quad (4.20)$$

Clearly, the length of the interval is defined by how far from zero the value of $z(\bar{\mathbf{y}})$ is. Suppose then that $z(\bar{\mathbf{y}}) \geq -\varepsilon$, for some small value $\varepsilon > 0$. (In an algorithm where a sequence $\{\mathbf{x}_k\}$ is constructed such that it converges to an optimal solution, this will eventually happen for every $\varepsilon > 0$.) Then, from the above we obtain that $f(\mathbf{x}^*) \geq f(\mathbf{x}) + z(\bar{\mathbf{y}}) \geq f(\mathbf{x}) - \varepsilon$; in short,

$$f(\mathbf{x}^*) \geq f(\mathbf{x}) - \varepsilon, \quad \text{or,} \quad f(\mathbf{x}) \leq f^* + \varepsilon. \quad (4.21)$$

We refer to a vector $\mathbf{x} \in S$ satisfying the inequality (4.21) as an ε -*optimal solution*. From the above linear minimization problem we hence have a simple instrument for evaluating the quality of a feasible solution in our problem. Note, again, that convexity is a crucial property enabling this possibility.

4.6 Applications

4.6.1 Continuity of convex functions

A remarkable property of any convex function is that without any additional assumptions it can be shown to be *continuous* relative to any open convex set in the intersection of its effective domain (that is, where the function has a finite value) and its affine hull.⁵ We establish a slightly weaker special case of this result below, in which “relative interior” is replaced by “interior” for simplicity.

Theorem 4.27 (continuity of convex functions) *Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is a convex function, and consider an open convex subset S of its effective domain. The function f is continuous on S .*

Proof. Let $\bar{\mathbf{x}} \in S$. To establish continuity of f at $\bar{\mathbf{x}}$, we must show that given $\varepsilon > 0$, there exists $\delta > 0$ such that $\|\mathbf{x} - \bar{\mathbf{x}}\| \leq \delta$ implies that $|f(\mathbf{x}) - f(\bar{\mathbf{x}})| \leq \varepsilon$. We establish this property in two parts, by showing that f is both lower and upper semi-continuous at $\bar{\mathbf{x}}$ (cf. Definition 4.6).

[upper semi-continuity] By the openness of S , there exists $\delta' > 0$ with $\|\mathbf{x} - \bar{\mathbf{x}}\| \leq \delta'$, implying $\mathbf{x} \in S$. Construct the value of the scalar γ as follows:

$$\gamma := \max_{i \in \{1, 2, \dots, n\}} \{ \max \{ |f(\bar{\mathbf{x}} + \delta' \mathbf{e}_i) - f(\bar{\mathbf{x}})|, |f(\bar{\mathbf{x}} - \delta' \mathbf{e}_i) - f(\bar{\mathbf{x}})| \} \}, \quad (4.22)$$

⁵In other words, f is continuous relative to any relatively open convex subset of its effective domain.

where \mathbf{e}_i is the i^{th} unit vector in \mathbb{R}^n . If $\gamma = 0$ it follows that f is constant in a neighbourhood of $\bar{\mathbf{x}}$ and hence continuous there, so suppose that $\gamma > 0$. Let now

$$\delta := \text{minimum} \left\{ \frac{\delta'}{n}, \frac{\varepsilon \delta'}{\gamma n} \right\}. \quad (4.23)$$

Choose an \mathbf{x} with $\|\mathbf{x} - \bar{\mathbf{x}}\| \leq \delta$. For every $i \in \{1, 2, \dots, n\}$, if $x_i \geq \bar{x}_i$ then define $\mathbf{z}_i := \delta' \mathbf{e}_i$, otherwise $\mathbf{z}_i := -\delta' \mathbf{e}_i$. Then, $\mathbf{x} - \bar{\mathbf{x}} = \sum_{i=1}^n \alpha_i \mathbf{z}_i$, where $\alpha_i \geq 0$ for all i . Moreover,

$$\|\mathbf{x} - \bar{\mathbf{x}}\| = \delta' \|\boldsymbol{\alpha}\|. \quad (4.24)$$

From (4.23), and since $\|\mathbf{x} - \bar{\mathbf{x}}\| \leq \delta$, it follows that $\alpha_i \leq 1/n$ for all i . Hence, by the convexity of f and since $0 \leq \alpha_i n \leq 1$, we get

$$\begin{aligned} f(\mathbf{x}) &= f\left(\bar{\mathbf{x}} + \sum_{i=1}^n \alpha_i \mathbf{z}_i\right) = f\left[\frac{1}{n} \sum_{i=1}^n (\bar{\mathbf{x}} + \alpha_i n \mathbf{z}_i)\right] \\ &\leq \frac{1}{n} \sum_{i=1}^n f(\bar{\mathbf{x}} + \alpha_i n \mathbf{z}_i) \\ &= \frac{1}{n} \sum_{i=1}^n f[(1 - \alpha_i n) \bar{\mathbf{x}} + \alpha_i n (\bar{\mathbf{x}} + \mathbf{z}_i)] \\ &\leq \frac{1}{n} \sum_{i=1}^n [(1 - \alpha_i n) f(\bar{\mathbf{x}}) + \alpha_i n f(\bar{\mathbf{x}} + \mathbf{z}_i)]. \end{aligned}$$

Therefore, $f(\mathbf{x}) - f(\bar{\mathbf{x}}) \leq \sum_{i=1}^n \alpha_i [f(\bar{\mathbf{x}} + \mathbf{z}_i) - f(\bar{\mathbf{x}})]$. From (4.22) it is obvious that $f(\bar{\mathbf{x}} + \mathbf{z}_i) - f(\bar{\mathbf{x}}) \leq \gamma$ for each i ; and since $\alpha_i \geq 0$, it follows that

$$f(\mathbf{x}) - f(\bar{\mathbf{x}}) \leq \gamma \sum_{i=1}^n \alpha_i. \quad (4.25)$$

Noting (4.23), (4.24), it follows that $\alpha_i \leq \varepsilon/n\gamma$, and (4.25) implies that $f(\mathbf{x}) - f(\bar{\mathbf{x}}) \leq \varepsilon$. Hence, we have so far shown that $\|\mathbf{x} - \bar{\mathbf{x}}\| \leq \delta$ implies that $f(\mathbf{x}) - f(\bar{\mathbf{x}}) \leq \varepsilon$. By Definition 4.6(b), f hence is upper semi-continuous at $\bar{\mathbf{x}}$.

[lower semi-continuity] Let $\mathbf{y} := 2\bar{\mathbf{x}} - \mathbf{x}$, and note that $\|\mathbf{y} - \bar{\mathbf{x}}\| \leq \delta$. Therefore, as above,

$$f(\mathbf{y}) - f(\bar{\mathbf{x}}) \leq \varepsilon. \quad (4.26)$$

But $\bar{\mathbf{x}} = \frac{1}{2}\mathbf{y} + \frac{1}{2}\mathbf{x}$, and by the convexity of f , $f(\bar{\mathbf{x}}) \leq \frac{1}{2}f(\mathbf{y}) + \frac{1}{2}f(\mathbf{x})$ follows. Combining this inequality with (4.26), it follows that $f(\bar{\mathbf{x}}) - f(\mathbf{x}) \leq \varepsilon$, whence Definition 4.6(a) applies. We are done. ■

Note that convex functions need not be continuous everywhere; by the above theorem we know however that points of non-continuity must occur at the boundary of the effective domain of f . For example, check the continuity of the following convex function:

$$f(x) := \begin{cases} x^2, & \text{for } |x| < 1, \\ 2, & \text{for } |x| = 1. \end{cases}$$

4.6.2 The Separation Theorem

The previously established Weierstrass Theorem 4.7 will now be utilized together with the variational inequality characterization (4.10) of stationary points in order to finally establish the Separation Theorem 3.24. For simplicity, we rephrase the theorem.

Theorem 4.28 (Separation Theorem) *Suppose that the set $S \subseteq \mathbb{R}^n$ is closed and convex, and that the point \mathbf{y} does not lie in S . Then, there exist a vector $\boldsymbol{\pi} \neq \mathbf{0}^n$ and $\alpha \in \mathbb{R}$ such that $\boldsymbol{\pi}^T \mathbf{y} > \alpha$ and $\boldsymbol{\pi}^T \mathbf{x} \leq \alpha$ for all $\mathbf{x} \in S$.*

Proof. We may assume that S is nonempty, and define a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ through $f(\mathbf{x}) := \|\mathbf{x} - \mathbf{y}\|^2/2$, $\mathbf{x} \in \mathbb{R}^n$. By Weierstrass' Theorem 4.7 there exists a minimum $\tilde{\mathbf{x}}$ of f over S , which by the first order necessary condition given in Proposition 4.23(b) satisfies $(\mathbf{y} - \tilde{\mathbf{x}})^T(\mathbf{x} - \tilde{\mathbf{x}}) \leq 0$ for all $\mathbf{x} \in S$ (since $-\nabla f(\tilde{\mathbf{x}}) = \mathbf{y} - \tilde{\mathbf{x}}$). Setting $\boldsymbol{\pi} := \mathbf{y} - \tilde{\mathbf{x}} \neq \mathbf{0}^n$ and $\alpha := (\mathbf{y} - \tilde{\mathbf{x}})^T \tilde{\mathbf{x}}$ gives the result sought. ■

A slightly different separation theorem will be used in the Lagrangian duality theory in Chapter 6. We state it without proof.

Theorem 4.29 (separation of convex sets) *Each pair of disjoint ($A \cap B = \emptyset$) nonempty convex sets A and B in \mathbb{R}^n can be separated by a hyperplane in \mathbb{R}^n , that is, there exist a vector $\boldsymbol{\pi} \in \mathbb{R}^n$ and $\alpha \in \mathbb{R}$ such that $\boldsymbol{\pi}^T \mathbf{x} \leq \alpha$ for all $\mathbf{x} \in A$ and $\boldsymbol{\pi}^T \mathbf{y} \geq \alpha$ for all $\mathbf{y} \in B$.* ■

Remark 4.30 A main difference between the Separation Theorems 3.24 and 4.29 is that in Theorem 3.24 there exists a hyperplane that in fact *strictly* separates the point \mathbf{y} and the closed convex set C , that is, there exists a vector $\boldsymbol{\pi} \in \mathbb{R}^n$ and an $\alpha \in \mathbb{R}$ such that $\boldsymbol{\pi}^T \mathbf{y} > \alpha$ while $\boldsymbol{\pi}^T \mathbf{x} < \alpha$ holds for all $\mathbf{x} \in C$. In Theorem 4.29, however, this is not true. Consider, for example, the sets $A := \{\mathbf{x} \in \mathbb{R}^2 \mid x_2 \leq 0\}$ and $B := \{\mathbf{x} \in \mathbb{R}^2 \mid x_1 > 0; x_2 \geq 1/x_1\}$. Then, the line $\{\mathbf{x} \in \mathbb{R}^2 \mid x_2 = 0\}$ separates A and B , but the sets cannot be strictly separated. ■

4.6.3 Euclidean projection

We will finish our discussions on the projection operation, which was defined in Section 3.4, by establishing an interesting continuity property.

Definition 4.31 (non-expansive operator) *Suppose that $S \subseteq \mathbb{R}^n$ is closed and convex. Let $\mathbf{f} : S \rightarrow S$ denote a vector-valued operator from S to S . We say that \mathbf{f} is non-expansive if, as a result of applying the mapping \mathbf{f} , the distance between any two vectors \mathbf{x} and \mathbf{y} in S does not increase.*

In other words, the operator \mathbf{f} is non-expansive on S if

$$\|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})\| \leq \|\mathbf{x} - \mathbf{y}\|, \quad \mathbf{x}, \mathbf{y} \in S, \quad (4.27)$$

holds.

Theorem 4.32 (the projection operation is non-expansive) *Let S be a nonempty, closed and convex set in \mathbb{R}^n . For every $\mathbf{x} \in \mathbb{R}^n$, its projection $\text{Proj}_S(\mathbf{x})$ is uniquely defined. The operator $\text{Proj}_S : \mathbb{R}^n \rightarrow S$ is non-expansive on \mathbb{R}^n , and therefore in particular continuous.*

Proof. The uniqueness of the operation is the result of the fact that the function $\mathbf{x} \mapsto \|\mathbf{x} - \mathbf{z}\|^2$ is both coercive and strictly convex on S , so there exists a unique optimal solution to the projection problem for every $\mathbf{z} \in \mathbb{R}^n$. (Cf. Weierstrass' Theorem 4.7 and Proposition 4.11, respectively.)

Next, take $\mathbf{x}^1, \mathbf{x}^2 \in \mathbb{R}^n$. Then, by the characterization (4.15) of the Euclidean projection,

$$\begin{aligned} [\mathbf{x}^1 - \text{Proj}_S(\mathbf{x}^1)]^T (\text{Proj}_S(\mathbf{x}^2) - \text{Proj}_S(\mathbf{x}^1)) &\leq 0, \\ [\mathbf{x}^2 - \text{Proj}_S(\mathbf{x}^2)]^T (\text{Proj}_S(\mathbf{x}^1) - \text{Proj}_S(\mathbf{x}^2)) &\leq 0. \end{aligned}$$

Summing the two inequalities yields

$$\begin{aligned} \|\text{Proj}_S(\mathbf{x}^2) - \text{Proj}_S(\mathbf{x}^1)\|^2 &\leq [\text{Proj}_S(\mathbf{x}^2) - \text{Proj}_S(\mathbf{x}^1)]^T (\mathbf{x}^2 - \mathbf{x}^1) \\ &\leq \|\text{Proj}_S(\mathbf{x}^2) - \text{Proj}_S(\mathbf{x}^1)\| \cdot \|\mathbf{x}^2 - \mathbf{x}^1\|, \end{aligned}$$

where the last inequality is a consequence of Cauchy's inequality; we obtain that $\|\text{Proj}_S(\mathbf{x}^2) - \text{Proj}_S(\mathbf{x}^1)\| \leq \|\mathbf{x}^2 - \mathbf{x}^1\|$. Since this is true for every pair $\mathbf{x}^1, \mathbf{x}^2$ of vectors in \mathbb{R}^n , we have shown that the operator Proj_S is non-expansive on \mathbb{R}^n . In particular, non-expansive functions are continuous. (The proof of the latter is left as an exercise.) ■

The theorem is illustrated in Figure 4.5.

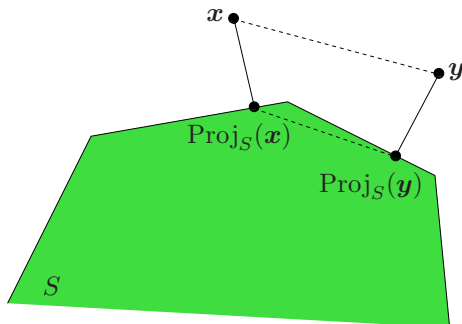


Figure 4.5: The projection operation is non-expansive.

4.6.4 Fixed point theorems

Fixed point theorems state properties of a problem of the following form: Suppose the mapping \mathbf{f} is defined on a closed, convex set S in \mathbb{R}^n and that $\mathbf{f}(\mathbf{x}) \subseteq S$ for every $\mathbf{x} \in S$. Is there an $\mathbf{x} \in S$ such that \mathbf{f} maps \mathbf{x} onto itself (that is, onto \mathbf{x}), or, in other words,

$$\exists \mathbf{x} \in S \text{ such that } \mathbf{x} \in \mathbf{f}(\mathbf{x})?$$

Such a point is called a *fixed point* of \mathbf{f} over the set S . If the mapping \mathbf{f} is single-valued rather than set-valued then the question boils down to:

$$\exists \mathbf{x} \in S \text{ such that } \mathbf{x} = \mathbf{f}(\mathbf{x})?$$

Many questions in optimization and analysis can be reduced to the analysis of a fixed point problem. For example, an optimization problem can in some circumstances be reduced to a fixed point problem, in which case the question of the existence of solutions to the optimization problem can be answered by studying the fixed point problem. Further, the optimality conditions analyzed in Section 4.4 can be written as the solution to a fixed point problem, cf. (4.11); we can therefore equate the search for a stationary point with that of finding a fixed point of a particular function \mathbf{f} . This type of analysis is quite useful also when analyzing the convergence of iterative algorithms for optimization problems.

4.6.4.1 Theory

We begin by studying some classic fixed point theorems, and then we provide examples of the connections between the results in Section 4.4 with fixed point theory.

Definition 4.33 (contractive operator) Let $S \subseteq \mathbb{R}^n$ be a nonempty set in \mathbb{R}^n . Let \mathbf{f} be a mapping from S to S . We say that \mathbf{f} is contractive on S if, as a result of applying the mapping \mathbf{f} , the distance between any two distinct vectors \mathbf{x} and \mathbf{y} in S decreases.

In other words, the operator \mathbf{f} is contractive on S if there exists $\alpha \in [0, 1)$ such that

$$\|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{y})\| \leq \alpha \|\mathbf{x} - \mathbf{y}\|, \quad \mathbf{x}, \mathbf{y} \in S, \quad (4.28)$$

holds. ■

Clearly, a contractive operator is non-expansive.

In the below result we utilize the notion of a *geometric convergence rate*; while its definition is in fact given in the result below, we also refer to Sections 6.4.1 and 11.10 for more detailed discussions on convergence rates.

Theorem 4.34 (fixed point theorems) Let S be a nonempty and closed set in \mathbb{R}^n .

(a) [Banach's Theorem] Let $\mathbf{f} : S \rightarrow S$ be a contraction mapping. Then, \mathbf{f} has a unique fixed point $\mathbf{x}^* \in S$. Further, for every initial vector $\mathbf{x}_0 \in S$, the iteration sequence $\{\mathbf{x}_k\}$ defined by the fixed-point iteration

$$\mathbf{x}_{k+1} := \mathbf{f}(\mathbf{x}_k), \quad k = 0, 1, \dots, \quad (4.29)$$

converges geometrically to the unique fixed point \mathbf{x}^* . In particular,

$$\|\mathbf{x}_k - \mathbf{x}^*\| \leq \alpha^k \|\mathbf{x}_0 - \mathbf{x}^*\|, \quad k = 0, 1, \dots$$

(b) [Brouwer's Theorem] Let S further be convex and bounded, and assume merely that \mathbf{f} is continuous. Then, \mathbf{f} has a fixed point.

Proof. (a) For any $\mathbf{x}_0 \in S$, consider the sequence $\{\mathbf{x}_k\}$ defined by (4.29). Then, for any $p \geq 1$,

$$\begin{aligned} \|\mathbf{x}_{k+p} - \mathbf{x}_k\| &\leq \sum_{i=1}^p \|\mathbf{x}_{k+i} - \mathbf{x}_{k+i-1}\| \\ &\leq (\alpha^{p-1} + \dots + 1) \|\mathbf{x}_{k+1} - \mathbf{x}_k\| \leq [\alpha^k / (1 - \alpha)] \|\mathbf{x}_1 - \mathbf{x}_0\|. \end{aligned}$$

Hence, $\{\mathbf{x}_k\}$ is a Cauchy sequence and thus converges as $k \rightarrow \infty$. By continuity and the contraction property, the limit point is the unique fixed point. In detail, this last part of the proof is as follows: Suppose that $\mathbf{x}_k \rightarrow \mathbf{x}^* \in S$. Then, for any iterate $\mathbf{x}_k \neq \mathbf{x}^*$ it holds that

$$\begin{aligned} \|\mathbf{f}(\mathbf{x}^*) - \mathbf{x}^*\| &\leq \|\mathbf{f}(\mathbf{x}^*) - \mathbf{f}(\mathbf{x}_k)\| + \|\mathbf{f}(\mathbf{x}_k) - \mathbf{x}^*\| \\ &\leq \alpha \|\mathbf{x}^* - \mathbf{x}_k\| + \|\mathbf{x}_{k+1} - \mathbf{x}^*\|, \end{aligned}$$

and according to our assumptions both of the latter terms converge to zero. Hence, \mathbf{x}^* is a fixed point. Suppose then that there is also another fixed point, \mathbf{x}^{**} , so that $\mathbf{x}^{**} \neq \mathbf{x}^*$. Then, $\|\mathbf{x}^* - \mathbf{x}^{**}\| = \|\mathbf{f}(\mathbf{x}^*) - \mathbf{f}(\mathbf{x}^{**})\| \leq \alpha \|\mathbf{x}^* - \mathbf{x}^{**}\|$, which yields a contradiction since $\alpha < 1$. Hence, $\mathbf{x}^{**} = \mathbf{x}^*$, and the fixed point \mathbf{x}^* is unique.

The convergence speed follows from the identification

$$\|\mathbf{x}_k - \mathbf{x}^*\| = \|\mathbf{f}(\mathbf{x}_{k-1}) - \mathbf{f}(\mathbf{x}^*)\| \leq \alpha \|\mathbf{x}_{k-1} - \mathbf{x}^*\|, \quad k = 1, 2, \dots$$

Applying this relation recursively yields the result.

(b) [Sketch] In short, the proof is to first establish that any C^1 function on the unit ball has a fixed point there. Extending the reasoning to merely continuous operators is possible, because of the Stone–Weierstrass Theorem (which states that for any continuous operator defined on the unit ball there is a sequence of C^1 functions defined on the unit ball that uniformly converges to it). Each of these functions can be established to have a fixed point, and because of the compactness of the unit ball, so does the merely continuous limit function. For our final argument, we can assume that the set S has a nonempty interior. Then there exists a homeomorphism⁶ $\mathbf{h} : S \rightarrow B$, where B is the unit ball. Since the composite mapping $\mathbf{h} \circ \mathbf{f} \circ \mathbf{h}^{-1}$ is a continuous operator from B to B it has a fixed point \mathbf{y} in B ; therefore, $\mathbf{h}^{-1}(\mathbf{y})$ is a fixed point of \mathbf{f} . ■

Notice that Banach’s Theorem holds without a convexity assumption on the set S ; the contraction property is indeed very strong. Brouwer’s Theorem constitutes a major improvement in that the mapping \mathbf{f} need only be continuous; on the other hand, the requirements on the set S increase dramatically. (This is a clear example of a case where a result (in this case: the existence of a fixed point) requires a sufficient “critical mass” of properties which however can be distributed in different ways.) An alternative proof of Banach’s Theorem is provided in Exercise 4.10.

A special case in one variable of the result in Brouwer’s Theorem is illustrated in Figure 4.6.

4.6.4.2 Applications

Particularly the result of Theorem 4.34(b) is quite remarkably strong. We provide some sample consequences of it below. In each case, we ask the reader to find the pair (S, \mathbf{f}) defining the corresponding fixed point problem.

⁶The given function \mathbf{h} is a *homeomorphism* if it is a continuous operator which is *onto*—that is, its range, $\mathbf{h}(S)$, is identical to the set B defining its image set—and has a continuous inverse.

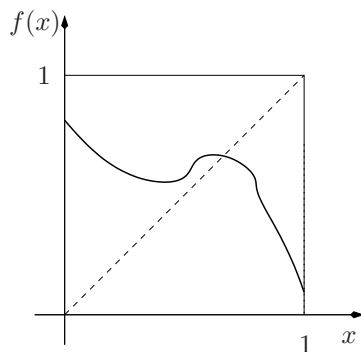


Figure 4.6: Consider the case $S = [0, 1]$, and a continuous function $f : S \rightarrow S$. Brouwer's Theorem states that there exists an $x^* \in S$ with $f(x^*) = x^*$. This is the same as saying that the continuous curve starting at $(0, f(0))$ and ending at $(1, f(1))$ must pass through the line $y = x$ inside the square.

- [Mountaineering] You climb a mountain, following a trail, in six hours (noon to 6 PM). You camp on top overnight. Then at noon the next day, you start descending. The descent is easier, and you make much better time. After an hour, you notice that your compass is missing, and you turn around and ascend a short distance, where you find your compass. You sit on a rock to admire the view. Then you descend the rest of the way. The entire descent takes four hours (noon to 4 PM). Along the trail there must then be a place where you were at the same place at the same time on both days.
- [Maps] Suppose you have two city maps over Gothenburg, which are not of the same scale. You crumple one of them up into a loose ball and place it on top of the other map entirely within the borders of the Gothenburg region on the flat map. Then, there is a point on the crumpled map (that represents the same place in Gothenburg on both maps) that is directly over its twin on the flat map. (A more simple problem is defined by a non-crumpled map and the city of Gothenburg itself; lay down the map anywhere in Gothenburg, and at least one point on the map will lie over that exact spot in real-life Gothenburg.)
- [Raking of gravel] Suppose you wish to rake the gravel in your garden; if the area is, say, circular, then any continuous raking will leave at least one tiny stone (which one is a function of time)

in the same place.

- [Stirring coffee] Stirring the contents of a (convex) coffee cup in a continuous way, no matter how long you stir, some particle (which one is a function of time) will stay in the same position as it did before you began stirring.
- [Meteorology] Even as the wind blows across the Earth there will be one location where the wind is perfectly vertical (or, perfectly calm). This fact actually implies the existence of cyclones; not to mention whorls, or crowns, in your hair no matter how you comb it. (The latter result also bears its own name: The Hairy Ball Theorem; cf. [BoL00, pp. 186–187].)

Applying fixed point theorems to our own development, we take a look at the variational inequality (4.10). Rephrasing it in a more general form, the *variational inequality problem* (VIP) is, for some set $S \subseteq \mathbb{R}^n$ and mapping $f : S \rightarrow \mathbb{R}^n$ to find $\mathbf{x}^* \in S$ such that

$$\mathbf{f}(\mathbf{x}^*)^\top(\mathbf{x} - \mathbf{x}^*) \geq 0, \quad \mathbf{x} \in S. \quad (4.30)$$

In order to turn it into a fixed point problem, we construct the following composite operator from \mathbb{R}^n to S :

$$\mathbf{F}(\mathbf{x}) := \text{Proj}_S(\mathbf{x} - \mathbf{f}(\mathbf{x})), \quad \mathbf{x} \in \mathbb{R}^n,$$

and consider finding a fixed point of \mathbf{F} on S . Why is this operator a correct one? Because it is equivalent to the statement that

$$\text{Proj}_S(\mathbf{x} - \mathbf{f}(\mathbf{x})) = \mathbf{x};$$

the special case for $\mathbf{f} = \nabla f$ is found in (4.11). Applying a fixed point theorem to the above problem then proves that the variational inequality problem (4.30) has solutions whenever \mathbf{f} is continuous on S and S is nonempty, convex and compact. (Moreover, we have immediately found an iterative algorithm for the variational inequality problem: if the operator $\mathbf{x} \mapsto \text{Proj}_S(\mathbf{x} - \alpha \mathbf{f}(\mathbf{x}))$ is contractive for some $\alpha > 0$, then it defines a convergent algorithm.)

At the same time, we saw that the fixed point problem was defined through the same type of stationarity condition that we derived in Section 4.4 for differentiable optimization problems over convex sets. We have thereby also illustrated that stationarity in an optimization problem is intimately associated with fixed points of a particular operator.⁷

⁷The book [Pat98] analyzes a large variety of optimization algorithms by utilizing this connection.

As an exercise, we consider the problem to find an $x \in \mathbb{R}$ such that $f(x) = 0$, where $f : \mathbb{R} \rightarrow \mathbb{R}$ is twice differentiable near a zero of f . The classic Newton–Raphson algorithm has an iteration formula of the form

$$x_0 \in \mathbb{R}; \quad x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}, \quad k = 0, 1, \dots$$

If we assume that x^* is a zero of f at which $f'(x^*) > 0$,⁸ then by starting close enough to x^* we can prove that the above iteration formula defines a contraction, and hence we can establish local convergence. (Why?) Further analyses of Newton methods will be performed in Chapter 11.

A similar technique can be used to establish that a system of linear equations with a symmetric matrix is solvable by the classic Jacobi algorithm in numerical analysis, if the matrix is diagonally dominant; this condition is equivalent to the Jacobi algorithm’s algorithm-defining operator being a contraction. (Similar, but stronger, results can also be obtained for the Gauss–Seidel algorithm; cf. [OrR70, Kre78, BeT89].)

An elegant application of fixed point theorems is the analysis of *matrix games*. The famous Minimax Theorem of von Neumann is associated with the existence of a saddle point of a function of the form $(\mathbf{v}, \mathbf{w}) \mapsto L(\mathbf{v}, \mathbf{w}) := \mathbf{v}^T \mathbf{A} \mathbf{w}$. Von Neumann’s minimax theorem states that if V and W both are nonempty, convex and compact, then

$$\underset{\mathbf{v} \in V}{\text{minimum}} \underset{\mathbf{w} \in W}{\text{maximum}} \mathbf{v}^T \mathbf{A} \mathbf{w} = \underset{\mathbf{w} \in W}{\text{maximum}} \underset{\mathbf{v} \in V}{\text{minimum}} \mathbf{v}^T \mathbf{A} \mathbf{w}. \quad (4.31)$$

In order to prove this theorem we can use the above existence theorem for variational inequalities. Let

$$\mathbf{x} = \begin{pmatrix} \mathbf{v} \\ \mathbf{w} \end{pmatrix}; \quad \mathbf{f}(\mathbf{x}) = \begin{pmatrix} -\mathbf{A}^T \mathbf{v} \\ \mathbf{A} \mathbf{w} \end{pmatrix}; \quad S = V \times W.$$

It is a reasonably simple exercise to prove that the variational inequality (4.30) with the above identifications is equivalent to the saddle point conditions, which can also be written as the existence of a pair $(\mathbf{v}^*, \mathbf{w}^*) \in V \times W$ such that

$$(\mathbf{v}^*)^T \mathbf{A} \mathbf{w} \leq (\mathbf{v}^*)^T \mathbf{A} \mathbf{w}^* \leq \mathbf{v}^T \mathbf{A} \mathbf{w}^*, \quad (\mathbf{v}, \mathbf{w}) \in V \times W;$$

and we are done immediately.

Saddle point results will be returned to in the study of (Lagrangian) duality in the coming chapters, especially for linear programming (which was also von Neumann’s special interest).

⁸The sign of $f'(x^*)$ is immaterial, as long as $f'(x^*) \neq 0$.

4.7 Notes and further reading

Most of the material of this chapter is elementary (as it relies mostly on the Taylor expansion of differentiable functions), and can be found in most basic books on nonlinear optimization, such as [Man69, Zan69, Avr76, BSS93, Ber99].

Complexity analysis, in the form of the analysis of NP-hard problems and the “P = NP?” question, is beyond the scope of this book; we refer to [GaJ79] for an excellent introduction.

Weierstrass’ Theorem 4.7 is the strongest existence result for optimal solutions that does not utilize convexity. The result is credited to Karl Weierstrass, but it was in fact known already by Bernard Bolzano in 1817 (although then only available in manuscript form); it has strong connections to the theorem of the existence of intermediate values as well as to that on the existence of limit points of every bounded sequence (now often referred to as the Bolzano–Weierstrass Theorem), and the notion of Cauchy sequences, often also credited to Weierstrass and Augustin-Louis Cauchy, respectively.

The Frank–Wolfe Theorem in Corollary 4.9 is found in [FrW56]. The stronger result in Theorem 4.8 is found in [Eav71, BIO72]. Proposition 4.26 on the invariance of the gradient on the solution set is found in [Man88, BuF91].

The result in Theorem 4.34(a) is due to Banach [Ban22]; Theorem 4.34(b) is due to Brouwer [Bro09, Bro12], and Hadamard [Had10]. Fixed point theorems are developed in greater detail in [GrD03]. Non-cooperative game theory was developed in work by John von Neumann, together with Oskar Morgenstern (see [vNe28, vNM43]), and by John Nash [Nas50, Nas51].

As far as iterative algorithms go, it is quite often the case that for the problem (4.9) involving a (convex) feasible set the sequences $\{\mathbf{x}_k\}$ of iterates do not necessarily stay inside the feasible set $S \subset \mathbb{R}^n$. The reason is that even if the constraints are convex inequalities it is difficult to check when one reaches the boundary of S . We mention however two cases where feasible algorithms (that is, those for which $\{\mathbf{x}_k\} \subset S$ holds) are viable:

- (I) When S is a polyhedral set, then it is only a matter of solving a series of simple linear systems to check for the maximum step length along a feasible direction. Among the algorithms that actually are feasible we count the *simplex method* for linear programming (LP) problems (see Chapters 9 and 10), the *Frank–Wolfe* and *Simplicial decomposition* methods (see Sections 12.2 and 12.3) which build on solving such LP problems, and the *projection method* (see Sec-

tion 12.4) which builds on the property (4.11); see also Exercise 4.5. More on these algorithms will be said in Chapter 12.

- (II) When the set S has an interior point, we may replace the constraints with an *interior penalty function* which has an asymptote whenever approaching the boundary, thus automatically ensuring that iterates stay (strictly) feasible. More on a class of methods based on this penalty function is said in Chapter 13.

4.8 Exercises

Exercise 4.1 (redundant constraints) Consider the problem to

$$\begin{aligned} &\text{minimize } f(\mathbf{x}), \\ &\text{subject to } g(\mathbf{x}) \leq b, \end{aligned} \tag{4.32}$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g : \mathbb{R}^n \rightarrow \mathbb{R}$ are continuous functions, and $b \in \mathbb{R}$. Suppose that this problem has a globally optimal solution, \mathbf{x}^* , and that $g(\mathbf{x}^*) < b$ holds. Consider also the problem to

$$\begin{aligned} &\text{minimize } f(\mathbf{x}), \\ &\text{subject to } \mathbf{x} \in \mathbb{R}^n, \end{aligned} \tag{4.33}$$

in which we have removed the constraint. The question is under which circumstances this is valid.

(a) Show by means of a counter-example that the vector \mathbf{x}^* may *not* solve (4.33); in other words, in general we cannot throw away a constraint that is not active without affecting the optimal solution, even if it is inactive. Hence, it would be wrong to call such constraints “redundant.”

(b) Suppose now that f is convex on \mathbb{R}^n . Show that \mathbf{x}^* solves (4.33).

Exercise 4.2 (unconstrained optimization, exam 020826) Consider the unconstrained optimization problem to minimize the function

$$f(\mathbf{x}) := \frac{3}{2}(x_1^2 + x_2^2) + (1+a)x_1x_2 - (x_1 + x_2) + b$$

over \mathbb{R}^2 , where a and b are real-valued parameters. Find all values of a and b such that the problem has a unique optimal solution.

Exercise 4.3 (spectral theory and unconstrained optimization) Let \mathbf{A} be a symmetric $n \times n$ matrix. For $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{x} \neq \mathbf{0}^n$, consider the function $\rho(\mathbf{x}) := \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}}$, and the related optimization problem to

$$\text{minimize}_{\mathbf{x} \neq \mathbf{0}^n} \rho(\mathbf{x}). \tag{P}$$

Determine the stationary points as well as the global minima in the problem (P). Interpret the result in terms of linear algebra.

Exercise 4.4 (non-convex QP over subspaces) The Frank–Wolfe Theorem 4.9 can be further improved for some special cases of linear constraints. Suppose that $f(\mathbf{x}) := \frac{1}{2}\mathbf{x}^T\mathbf{Q}\mathbf{x} - \mathbf{q}^T\mathbf{x}$, where $\mathbf{Q} \in \mathbb{R}^{n \times n}$ is a symmetric matrix and $\mathbf{q} \in \mathbb{R}^n$. Suppose further that the constraints are equalities, that is, that the ℓ constraints define the linear system $\mathbf{E}\mathbf{x} = \mathbf{d}$, where $\mathbf{E} \in \mathbb{R}^{\ell \times n}$ and $\mathbf{d} \in \mathbb{R}^\ell$. Note that the problem may not be convex, as we have not assumed that \mathbf{Q} is positive semidefinite.

For this set-up, establish the following:

- (a) Every locally optimal solution is a globally optimal solution.
- (b) A locally [hence globally, by (a)] optimal solution exists if and only if f is lower bounded on $S := \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{E}\mathbf{x} = \mathbf{d}\}$.

Exercise 4.5 (descent from projection) Consider the problem (4.9), where f is in C^1 on the closed and convex set S . Let $\mathbf{x} \in S$. Let $\alpha > 0$, and define

$$\mathbf{p} := \text{Proj}_S[\mathbf{x} - \alpha \nabla f(\mathbf{x})] - \mathbf{x}.$$

Notice that \mathbf{p} is a feasible direction at \mathbf{x} . Establish that

$$\nabla f(\mathbf{x})^T \mathbf{p} \leq -\frac{1}{\alpha} \|\mathbf{p}\|^2$$

holds. Hence, \mathbf{p} is zero if and only if \mathbf{x} is stationary [according to the characterization in (4.11)], and if \mathbf{p} is non-zero then it defines a feasible descent direction with respect to f at \mathbf{x} .

Exercise 4.6 (optimality conditions for a special problem) Suppose that $f \in C^1$ on the set $S := \{\mathbf{x} \in \mathbb{R}^n \mid x_j \geq 0, \quad j = 1, 2, \dots, n\}$, and consider the problem of finding a minimum of $f(\mathbf{x})$ over S . Develop the necessary optimality conditions for this problem in a compact form.

Exercise 4.7 (optimality conditions for a special problem) Consider the problem to

$$\begin{aligned} &\text{maximize} && f(\mathbf{x}) := x_1^{a_1} x_2^{a_2} \cdots x_n^{a_n}, \\ &\text{subject to} && \sum_{j=1}^n x_j = 1, \\ &&& x_j \geq 0, \quad j = 1, \dots, n, \end{aligned}$$

where the values of a_j ($j = 1, \dots, n$) are positive. Find a global maximum and show that it is unique.

Exercise 4.8 (extensions of convexity, exam 040602) We have stressed that convexity is a crucial property of functions when analyzing optimization models in general and studying optimality conditions in particular. There are, however, certain properties of convex functions that are shared also by classes of non-convex functions. The purpose of this exercise is to relate the convex functions to two such classes of non-convex functions by means of some example properties.

Suppose that $S \subseteq \mathbb{R}^n$ and that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuous on S .

(a) Suppose further that f is in C^1 on S . We say that the function f is *pseudo-convex* on S if, for every $\mathbf{x}, \mathbf{y} \in S$,

$$\nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) \geq 0 \quad \implies \quad f(\mathbf{y}) \geq f(\mathbf{x}).$$

Establish the following two statements: (1) if f is a convex function on S then f is pseudo-convex on S (that is, “convexity implies pseudo-convexity”); (2) the reverse statement (“pseudo-convexity implies convexity”) is not true.

[Note: The definition of a pseudo-convex function is due to Mangasarian [Man65].]

(b) A well-known property of a differentiable convex function is its role in necessary and sufficient conditions for globally optimal solutions. Suppose now that S is convex. Establish that the equivalence stated in Theorem 4.24 still holds if the convexity of f on S is replaced by the pseudo-convexity of f on S .

(c) Let S be convex. We say that the function f is *quasi-convex* on S if its level sets restricted to S are convex. In other words, f is quasi-convex on S if

$$\text{lev}_f^S(b) := \{\mathbf{x} \in S \mid f(\mathbf{x}) \leq b\}$$

is convex for every $b \in \mathbb{R}$.

Establish the following two statements for a function f which is in C^1 on S : (1) if f is a convex function on S then f is quasi-convex on S (that is, “convexity implies quasi-convexity”); (2) the reverse statement (“quasi-convexity implies convexity”) is not true.

[Note: The definition of a quasi-convex function is due to De Finetti [DeF49].]

Exercise 4.9 (illustrations of fixed point results) (a) Let $S := \{x \in \mathbb{R} \mid x \geq 1\}$ and $f(x) := x/2 + 1/x$. Show that f is a contraction and find the smallest value of α .

(b) In analysis, a usual condition for the convergence of an iteration $x_k = g(x_{k-1})$ is that g is continuously differentiable and

$$|g'(x)| \leq \alpha < 1, \quad x \in S. \quad (4.34)$$

Establish that (4.34) implies convergence, by using Banach’s Theorem 4.34(a).

(c) Show that a fixed-point iteration for calculating the square root of a given positive number c is

$$x_0 > 0; \quad x_{k+1} = g(x_k) := \frac{1}{2} \left(x_k + \frac{c}{x_k} \right), \quad k = 0, 1, \dots$$

What condition do we get from (b)? Starting at $x_0 = 1$, calculate the approximations x_1, x_2, x_3, x_4 of $\sqrt{2}$.

Exercise 4.10 (Ekeland's Variational Principle and Banach's Theorem) Ekeland's variational principle states that if a lower semicontinuous function f attains a value close to its infimum at some point then a nearby point minimizes a slightly perturbed function exactly. We state and prove this result and then utilize it to prove Banach's Theorem 4.34(a).

Theorem 4.35 (Ekeland's variational principle) Suppose the function $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is lower semicontinuous and that the point $\mathbf{x} \in \mathbb{R}^n$ satisfies $f(\mathbf{x}) \leq f^* + \varepsilon$ for some $\varepsilon > 0$. Then for any real $\lambda > 0$ there is a point $\mathbf{v} \in \mathbb{R}^n$ such that

- (a) $\|\mathbf{x} - \mathbf{v}\| \leq \lambda$;
- (b) $f(\mathbf{v}) \leq f(\mathbf{x})$; and
- (c) \mathbf{v} uniquely minimizes the function $\mathbf{x} \mapsto f(\mathbf{x}) + \frac{\varepsilon}{\lambda} \|\mathbf{x} - \mathbf{v}\|$ over \mathbb{R}^n .

Proof. We can assume that f is proper, and by assumption it is bounded below. Since the function $f(\cdot) + \frac{\varepsilon}{\lambda} \|\cdot - \mathbf{x}\|$ therefore has compact level sets, its set of minimizers $M \subset \mathbb{R}^n$ is nonempty and compact.

Choose a minimizer \mathbf{v} of f over M . Then for points $\mathbf{z} \neq \mathbf{v}$ in M we know that

$$f(\mathbf{v}) \leq f(\mathbf{z}) < f(\mathbf{z}) + \frac{\varepsilon}{\lambda} \|\mathbf{z} - \mathbf{v}\|,$$

while for \mathbf{z} not in M we have that

$$f(\mathbf{v}) + \frac{\varepsilon}{\lambda} \|\mathbf{v} - \mathbf{x}\| < f(\mathbf{z}) + \frac{\varepsilon}{\lambda} \|\mathbf{z} - \mathbf{x}\|.$$

Part (c) follows by the triangle inequality. Since \mathbf{v} lies in M we have that

$$f(\mathbf{z}) + \frac{\varepsilon}{\lambda} \|\mathbf{z} - \mathbf{x}\| \geq f(\mathbf{v}) + \frac{\varepsilon}{\lambda} \|\mathbf{v} - \mathbf{x}\|, \quad \mathbf{z} \in \mathbb{R}^n.$$

Setting $\mathbf{z} = \mathbf{x}$ shows the inequalities

$$f(\mathbf{v}) + \varepsilon \geq f^* + \varepsilon \geq f(\mathbf{x}) \geq f(\mathbf{v}) + \frac{\varepsilon}{\lambda} \|\mathbf{v} - \mathbf{x}\|.$$

Properties (a) and (b) follow. ■

Using this result, prove Banach's Theorem 4.34(a).

[Hint: Apply the Ekeland variational principle to the function

$$\mathbb{R}^n \ni \mathbf{z} \mapsto \begin{cases} \|\mathbf{z} - \mathbf{f}(\mathbf{z})\|, & \text{if } \mathbf{z} \in S, \\ +\infty, & \text{otherwise} \end{cases}$$

at an arbitrary point \mathbf{x} in S , with the choice of reals

$$\varepsilon := \|\mathbf{x} - \mathbf{f}(\mathbf{x})\| \quad \text{and} \quad \lambda := \frac{\varepsilon}{1 - \alpha},$$

where α is the contraction parameter for \mathbf{f} .]

[Note: Ekeland's variational principle is found in [Eke74], while its use in the proof of Banach's Theorem can be found, for example, in [BoL00, Theorem 8.1.2].]

Optimality conditions

V

5.1 Relations between optimality conditions and CQs at a glance

Optimality conditions are introduced as an attempt to construct an easily verifiable criterion that allows us to examine points in a feasible set, one after another, and classify them into optimal and non-optimal ones. Unfortunately, this is impossible in practice, and not only due to the fact that there are far too many feasible points, but also because it is impossible to construct such a universal criterion. It is usually possible to construct either practical (that is, computationally verifiable) conditions that admit some mistakes in the characterization, or perfect ones which are impossible to use in the computations. It is of course the first group that is of practical value for us, and it may further be classified into two distinct subgroups based on the type of mistakes allowed in the decision-making process. Namely, optimality conditions encountered in practice are divided into two classes, known as *necessary* and *sufficient* conditions.

Necessary conditions must be satisfied at every locally optimal point; on the other hand, we cannot guarantee that every point satisfying the necessary optimality conditions is indeed locally optimal. On the contrary, sufficient optimality conditions provide such guarantees; however, there may be some locally optimal points that violate the optimality conditions. Arguably, it is much more important to be able to find a few candidates for local minima that can be further investigated by other means, than to eliminate some local (or even global) minima from the beginning. Therefore, this chapter is dedicated to the development of *necessary optimality conditions*. However, for *convex* optimization problems these conditions turn out to be *sufficient*.

Now, we can concentrate on what should be meant by easily verifiable conditions. A human being can immediately state whether a given point belongs to a simple set or not, by just glancing at a picture of it; for a numerical algorithm, a clear algebraic description of a set in terms of equalities and inequalities is vital. Therefore, we start our development with geometric optimality conditions (Section 5.3), to gain an understanding about the relationships between the gradient of the objective function and the feasible set that must hold at every local minimum point. Given a specific description of a feasible set in terms of inequalities, the geometric conditions immediately imply some relationships between the gradients of the objective function and the constraints that are active at the point under consideration (see Section 5.4); these conditions are known as the *Fritz John optimality conditions*, and are rather weak (i.e., they can be satisfied by many points that have nothing in common with locally optimal points). However, if we assume an additional regularity of the system of inequalities and equalities that define our feasible set, then the geometric optimality conditions imply stronger conditions, known as the *Karush–Kuhn–Tucker optimality conditions* (see Section 5.5). The additional regularity assumptions are known under the name *constraint qualifications* (CQs), and they vary from very abstract and difficult to check, but enjoyed by many feasible sets (such as, e.g., Abadie’s CQ, see Definition 5.23) to more specific, easily verifiable but also somewhat restrictive in many situations (such as the linear independence CQ, see Definition 5.41, or the Slater CQ, see Definition 5.38). In Section 5.8 we show that for convex problems the KKT conditions are sufficient for local, hence global, optimality.

The contents of this chapter are in principle summarized in the flow-chart in Figure 5.1. Various optimality conditions and constraint qualifications that are discussed in this chapter constitute the nodes of the flow-chart. Logical relationships between them are denoted with edges, and the direction of the arrow shows the direction of the logical implication; each implication is further labeled with the result that establishes it. We note that the KKT conditions “follow” from *both* geometric conditions and constraint qualifications satisfied at a given point; also, global optimality holds if *both* the KKT conditions are verified and the optimization problem is convex.

5.2 A note of caution

In this chapter we will discuss various *necessary* optimality conditions for a given point to be a local minimum to a nonlinear programming model. If the NLP is a convex program, any point satisfying these necessary

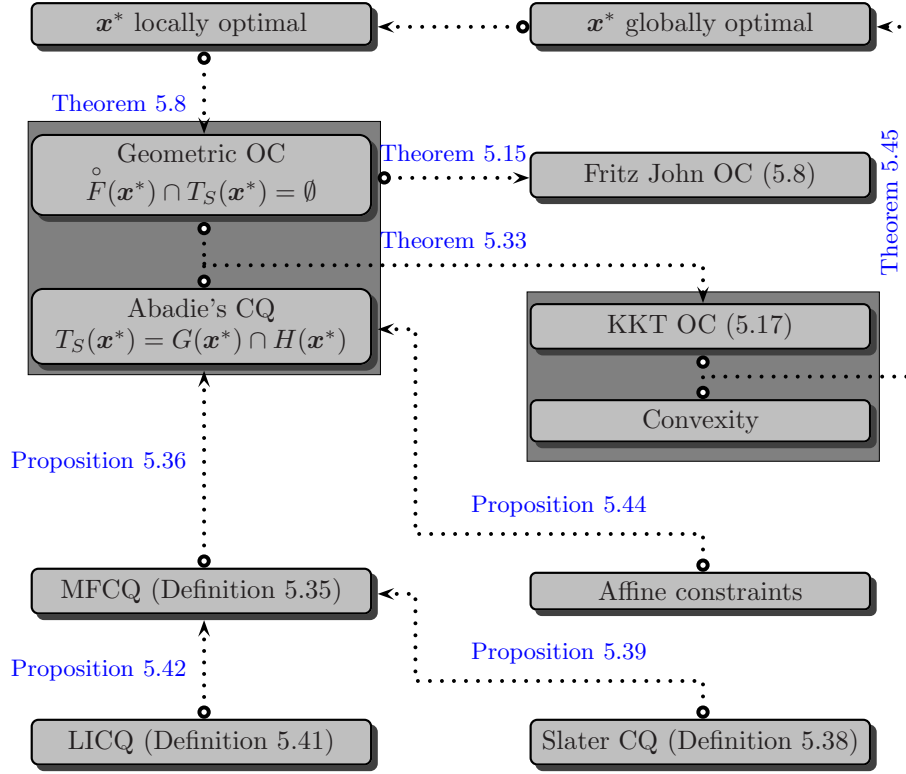


Figure 5.1: Relations between optimality conditions and CQs at a glance.

optimality conditions is not only a local minimum, but actually a global minimum (see Section 5.8). Arguably, most NLP models that arise in real world applications tend to be nonconvex, and for such a problem, a point satisfying the necessary optimality conditions may not even be a local minimum. Algorithms for NLP are usually designed to converge to a point satisfying the necessary optimality conditions, and as mentioned earlier, one should not blindly accept such a point as an optimum solution to the problem without checking (e.g., using the second order necessary optimality conditions, see [BSS93, Section 4.4], or by means of some local search in the vicinity of the point) that it is at least better than all the other nearby points. Also, the system of necessary optimality conditions may have many solutions. Finding alternate solutions of this system, and selecting the best among them, usually leads to a good

point to investigate further.

We will illustrate the importance of this with the story of the US Air Force's controversial B-2 Stealth bomber program in the Reagan era of the 1980s. There were many design variables, such as the various dimensions, the distribution of volume between the wing and the fuselage, flying speed, thrust, fuel consumption, drag, lift, air density, etc., that could be manipulated for obtaining the best range (i.e., the distance it can fly starting with full tank, without refueling). The problem of maximizing the range subject to all the constraints was modeled as an NLP in a secret Air Force study going back to the 1940s. A solution to the necessary optimality conditions of this problem was found; it specified values for the design variables that put almost all of the total volume in the wing, leading to the *flying wing design* for the B-2 bomber. After spending billions of dollars, building test planes, etc., it was found that the design solution implemented works, but that its range was too low in comparison with other bomber designs being experimented subsequently in the US and abroad.

A careful review of the model was then carried out. The review indicated that all the formulas used, and the model itself, are perfectly valid. However, the model was a nonconvex NLP, and the review revealed a second solution to the system of necessary optimality conditions for it, besides the one found and implemented as a result of earlier studies. The second solution makes the wing volume much less than the total volume, and seems to maximize the range; while the first solution that is implemented for the B-2 bomber seems to actually minimize the range. (The second solution also looked like an aircraft should, while the flying wing design was counter-intuitive.) In other words, the design implemented was the aerodynamically *worst* possible choice of configuration, leading to a very costly error. The aircraft does fly, but apparently, then, has the only advantage that it is a "stealth" plane.

For an account, see "Skeleton Alleged in the Stealth Bomber's Closet," *Science*, vol. 244, 12 May 1989 issue, pages 650–651.

5.3 Geometric optimality conditions

In this section we will discuss the optimality conditions for the following optimization problem [cf. (4.1)]:

$$\begin{aligned} &\text{minimize } f(\mathbf{x}), \\ &\text{subject to } \mathbf{x} \in S, \end{aligned} \tag{5.1}$$

where $S \subseteq \mathbb{R}^n$ is a nonempty closed set and $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a given differentiable function. Since we do not have any particular description

of the feasible set S in terms of equality or inequality constraints, the optimality conditions will be based on purely geometrical ideas. Being very general, the optimality conditions we develop in this section are almost useless for computations, because they are also not very easy (in fact they are even impossible) to verify for an optimization algorithm. Therefore, in the sections that follow we will use an algebraic description of the set S and geometric optimality conditions to further develop the classic Fritz John and Karush–Kuhn–Tucker optimality conditions in the form of easily verifiable systems of equations and inequalities.

The basic idea behind the optimality conditions is that if the point $\mathbf{x}^* \in S$ is a local minimum of f over S , it should not be possible to draw a curve starting at the point \mathbf{x}^* inside S , such that f decreases along it arbitrarily close to \mathbf{x}^* . Linearizing the objective function and the constraints along such curves, we eventually establish relationships between their gradients that are necessary to hold at points of local minima.

We first define the meaning of “possible to draw a curve starting at \mathbf{x}^* inside S .” Arguably, the simplest curves are the straight lines; the following definition gives exactly the set of lines that locally around \mathbf{x}^* belong to S .

Definition 5.1 (cone of feasible directions) *Let $S \subseteq \mathbb{R}^n$ be a nonempty closed set. The cone of feasible directions for S at $\mathbf{x} \in \mathbb{R}^n$, known also as the radial cone, is defined as:*

$$R_S(\mathbf{x}) := \{ \mathbf{p} \in \mathbb{R}^n \mid \exists \tilde{\delta} > 0 \text{ such that } \mathbf{x} + \delta \mathbf{p} \in S, 0 \leq \delta \leq \tilde{\delta} \}. \quad (5.2)$$

Thus, this is nothing else but the cone containing all feasible directions at \mathbf{x} in the sense of Definition 4.20. ■

This cone is used in some optimization algorithms, but unfortunately it is too small to develop optimality conditions that are general enough. Therefore, we consider less intuitive, but bigger and more well-behaving sets (cf. Proposition 5.3 and the examples that follow).

Definition 5.2 (tangent cone) *Let $S \subseteq \mathbb{R}^n$ be a nonempty closed set. The tangent cone for S at $\mathbf{x} \in \mathbb{R}^n$ is defined as*

$$T_S(\mathbf{x}) := \{ \mathbf{p} \in \mathbb{R}^n \mid \exists \{ \mathbf{x}_k \} \subset S, \{ \lambda_k \} \subset (0, \infty) : \lim_{k \rightarrow \infty} \mathbf{x}_k = \mathbf{x}; \lim_{k \rightarrow \infty} \lambda_k (\mathbf{x}_k - \mathbf{x}) = \mathbf{p} \}. \quad (5.3)$$

■

Thus, to construct a tangent cone we consider all the sequences $\{\mathbf{x}_k\}$ in S that converge to the given $\mathbf{x} \in \mathbb{R}^n$, and then calculate all the directions $\mathbf{p} \in \mathbb{R}^n$ that are tangential to the sequences at \mathbf{x} ; such tangential vectors are described as the limits of $\{\lambda_k(\mathbf{x}_k - \mathbf{x})\}$ for arbitrary positive sequences $\{\lambda_k\}$. Note that to generate a nonzero vector $\mathbf{p} \in T_S(\mathbf{x})$ the sequence $\{\lambda_k\}$ must converge to $+\infty$.

While it is possible that $\text{cl } R_S(\mathbf{x}) = T_S(\mathbf{x})$, or even that $R_S(\mathbf{x}) = T_S(\mathbf{x})$, in general we have only the following proposition, and examples that follow show that the two cones might be very different.

Proposition 5.3 (relationship between the radial and the tangent cones) *The tangent cone is a closed set, and the inclusion $\text{cl } R_S(\mathbf{x}) \subseteq T_S(\mathbf{x})$ holds for every $\mathbf{x} \in \mathbb{R}^n$.*

Proof. Consider a sequence $\{\mathbf{p}_k\} \subset T_S(\mathbf{x})$, and assume that $\{\mathbf{p}_k\} \rightarrow \mathbf{p}$. Since every $\mathbf{p}_k \in T_S(\mathbf{x})$, there exist $\mathbf{x}_k \in S$ and $\lambda_k > 0$, such that $\|\mathbf{x}_k - \mathbf{x}\| < k^{-1}$ and $\|\lambda_k(\mathbf{x}_k - \mathbf{x}) - \mathbf{p}_k\| < k^{-1}$. Then, clearly, $\{\mathbf{x}_k\} \rightarrow \mathbf{x}$, and, by the triangle inequality, $\|\lambda_k(\mathbf{x}_k - \mathbf{x}) - \mathbf{p}\| \leq \|\lambda_k(\mathbf{x}_k - \mathbf{x}) - \mathbf{p}_k\| + \|\mathbf{p}_k - \mathbf{p}\|$, and the two terms in the right-hand side converge to 0, which implies that $\mathbf{p} \in T_S(\mathbf{x})$ and thus the latter set is closed.

In view of the closedness of the tangent cone, it is enough to show the inclusion $R_S(\mathbf{x}) \subseteq T_S(\mathbf{x})$. Let $\mathbf{p} \in R_S(\mathbf{x})$. Then, for all large integers k it holds that $\mathbf{x} + k^{-1}\mathbf{p} \in S$, and, therefore, setting $\mathbf{x}_k = \mathbf{x} + k^{-1}\mathbf{p}$ and $\lambda_k = k$ we see that $\mathbf{p} \in T_S(\mathbf{x})$ as defined by Definition 5.2. ■

Example 5.4 Let $S := \{\mathbf{x} \in \mathbb{R}^2 \mid -x_1 \leq 0; (x_1 - 1)^2 + x_2^2 \leq 1\}$. Then, $R_S(\mathbf{0}^2) = \{\mathbf{p} \in \mathbb{R}^2 \mid p_1 > 0\}$, and $T_S(\mathbf{0}^2) = \{\mathbf{p} \in \mathbb{R}^2 \mid p_1 \geq 0\}$, i.e., $T_S(\mathbf{0}^2) = \text{cl } R_S(\mathbf{0}^2)$ (see Figure 5.2). ■

Example 5.5 (complementarity constraint) Let $S := \{\mathbf{x} \in \mathbb{R}^2 \mid -x_1 \leq 0; -x_2 \leq 0; x_1 x_2 \leq 0\}$. In this case, S is a (non-convex) cone, and $R_S(\mathbf{0}^2) = T_S(\mathbf{0}^2) = S$ (see Figure 5.3). ■

Example 5.6 Let $S := \{\mathbf{x} \in \mathbb{R}^2 \mid -x_1^3 + x_2 \leq 0; x_1^5 - x_2 \leq 0; -x_2 \leq 0\}$. Then, $R_S(\mathbf{0}^2) = \emptyset$, and $T_S(\mathbf{0}^2) = \{\mathbf{p} \in \mathbb{R}^2 \mid p_1 \geq 0; p_2 = 0\}$ (see Figure 5.4). ■

Example 5.7 Let $S := \{\mathbf{x} \in \mathbb{R}^2 \mid -x_2 \leq 0; (x_1 - 1)^2 + x_2^2 = 1\}$. Then, $R_S(\mathbf{0}^2) = \emptyset$, $T_S(\mathbf{0}^2) = \{\mathbf{p} \in \mathbb{R}^2 \mid p_1 = 0; p_2 \geq 0\}$ (see Figure 5.5). ■

We already know that f decreases along any descent direction (cf. Definition 4.15), and that for a vector $\mathbf{p} \in \mathbb{R}^n$ it is sufficient to verify the inequality $\nabla f(\mathbf{x}^*)^T \mathbf{p} < 0$ to be a descent direction for f at $\mathbf{x}^* \in \mathbb{R}^n$

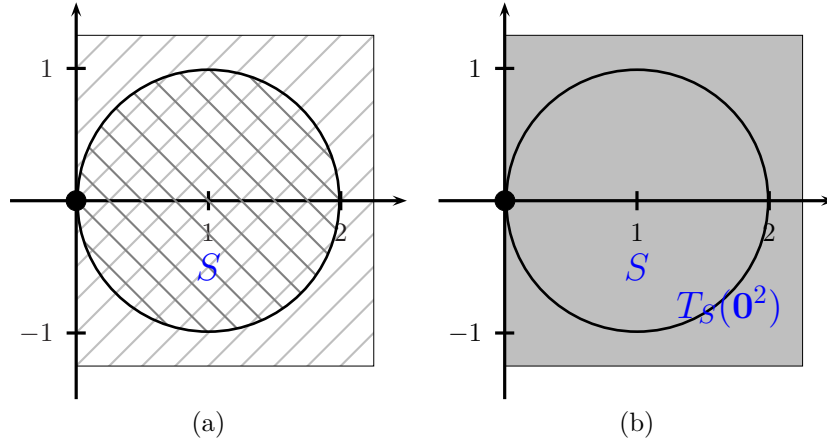


Figure 5.2: (a) The set S obtained as the intersection of the solution set of two constraints; (b) the tangent cone $T_S(\mathbf{0}^2)$ (see Example 5.4).

(see Proposition 4.16). Even though this condition is not necessary, it is very easy to check in practice and therefore we will use it to develop optimality conditions. Therefore, it would be convenient to define a cone of such directions (which is empty if $\nabla f(\mathbf{x}^*)$ happens to be $\mathbf{0}^n$):

$$\overset{\circ}{F}(\mathbf{x}^*) := \{\mathbf{p} \in \mathbb{R}^n \mid \nabla f(\mathbf{x}^*)^\top \mathbf{p} < 0\}. \quad (5.4)$$

Now we have the necessary notation in order to state and prove the main theorem of this section.

Theorem 5.8 (geometric necessary optimality conditions) *Consider the optimization problem (5.1). Then,*

$$\mathbf{x}^* \text{ is a local minimum of } f \text{ over } S \implies \overset{\circ}{F}(\mathbf{x}^*) \cap T_S(\mathbf{x}^*) = \emptyset,$$

where $\overset{\circ}{F}(\mathbf{x}^*)$ is defined by (5.4), and $T_S(\mathbf{x}^*)$ by Definition 5.2.

Proof. Assume that $\mathbf{p} \in T_S(\mathbf{x}^*)$, i.e., $\exists \{\mathbf{x}_k\} \subset S$, and $\{\lambda_k\} \subset (0, \infty)$ such that $\lim_{k \rightarrow \infty} \mathbf{x}_k = \mathbf{x}^*$ and $\lim_{k \rightarrow \infty} \lambda_k(\mathbf{x}_k - \mathbf{x}^*) = \mathbf{p}$. Using the first order Taylor expansion (2.1) we get:

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) = \nabla f(\mathbf{x}^*)^\top (\mathbf{x}_k - \mathbf{x}^*) + o(\|\mathbf{x}_k - \mathbf{x}^*\|) \geq 0,$$

where the last inequality holds for all enough large k by the local opti-

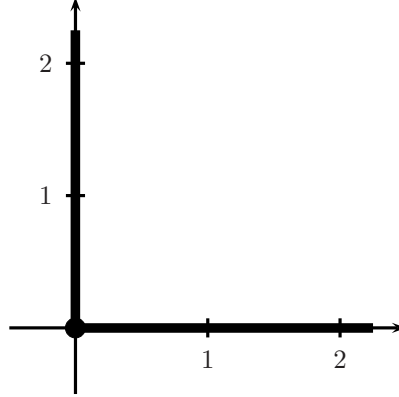


Figure 5.3: $S = R_S(\mathbf{0}^2) = T_S(\mathbf{0}^2)$ (see Example 5.5).

ality of \mathbf{x}^* . Multiplying by $\lambda_k > 0$ and taking limit we get

$$\begin{aligned} 0 &\leq \lim_{k \rightarrow \infty} \left[\lambda_k \nabla f(\mathbf{x}^*)^T (\mathbf{x}_k - \mathbf{x}^*) + \|\lambda_k (\mathbf{x}_k - \mathbf{x}^*)\| \frac{o(\|\mathbf{x}_k - \mathbf{x}^*\|)}{\|\mathbf{x}_k - \mathbf{x}^*\|} \right] \\ &= \nabla f(\mathbf{x}^*)^T \mathbf{p} + \|\mathbf{p}\| \cdot 0, \end{aligned}$$

and thus $\mathbf{p} \notin \overset{\circ}{F}(\mathbf{x}^*)$. ■

Combining Proposition 5.3 and Theorem 5.8 we get that

$$\mathbf{x}^* \text{ is a local minimum of } f \text{ over } S \implies \overset{\circ}{F}(\mathbf{x}^*) \cap R_S(\mathbf{x}^*) = \emptyset;$$

but this statement is weaker than Theorem 5.8.

Example 5.9 Consider the differentiable (linear) function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ defined by $f(\mathbf{x}) = x_1$. Then, $\nabla f = (1, 0)^T$, and $\overset{\circ}{F}(\mathbf{0}^2) = \{\mathbf{x} \in \mathbb{R}^2 \mid x_1 < 0\}$. It is easy to see from geometric considerations that $\mathbf{x}^* = \mathbf{0}^2$ is a local (in fact, even global) minimum in either problem (5.1) with S given by Examples 5.4–5.7, and equally easy it is to check that the geometric necessary optimality condition $\overset{\circ}{F}(\mathbf{0}^2) \cap T_S(\mathbf{0}^2) = \emptyset$ is satisfied in all these examples (which is no surprise, in view of Theorem 5.8). ■

5.4 The Fritz John conditions

Theorem 5.8 gives a very elegant criterion for checking whether a given point $\mathbf{x}^* \in S$ is a candidate for a local minimum of the problem (5.1),

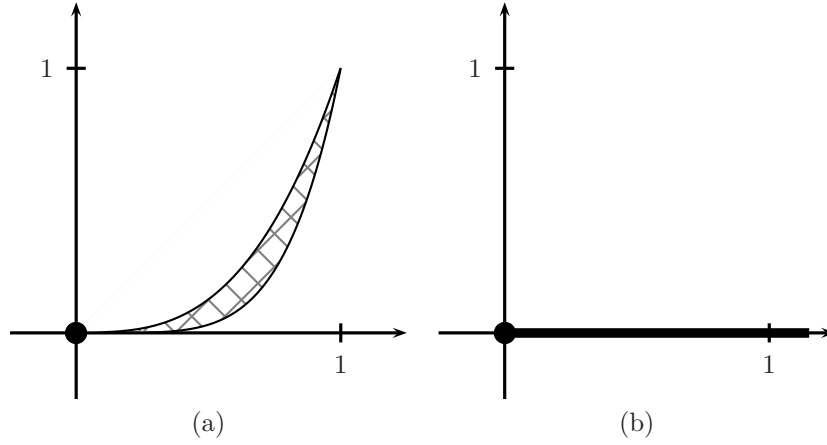


Figure 5.4: (a) The set S ; (b) the tangent cone $T_S(\mathbf{0}^2)$ (see Example 5.6).

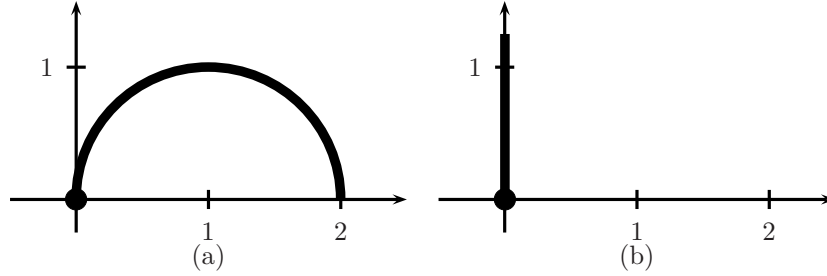


Figure 5.5: (a) The set S ; (b) the tangent cone $T_S(\mathbf{0}^2)$ (see Example 5.7).

but there is a catch: the set $T_S(\mathbf{x}^*)$ is close to impossible to compute for general sets S ! Therefore, in this section we will use an algebraic characterization of the set S to compute other cones that we hope could approximate $T_S(\mathbf{x}^*)$ in many practical situations.

Namely, we assume that the set S is defined as the solution set of a system of differentiable inequality constraints defined by the functions $g_i \in C^1(\mathbb{R}^n)$, $i = 1, \dots, m$, such that

$$S := \{ \mathbf{x} \in \mathbb{R}^n \mid g_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m \}. \quad (5.5)$$

We can always assume this structure, because any equality constraint $h(\mathbf{x}) = 0$ may be written in the form $h(\mathbf{x}) \leq 0$, $-h(\mathbf{x}) \leq 0$. Treating the equality constraints in this way we obtain the Fritz John conditions, that however are somewhat too weak to be practical; on the positive side, it significantly simplifies the notation and does not affect the development

Optimality conditions

of the KKT conditions. Therefore, we keep this assumption for some time, and state the KKT system that specifically distinguishes between the inequality and equality constraints in Section 5.6. We will use the symbol $\mathcal{I}(\mathbf{x})$ to denote the index set of active inequality constraints at $\mathbf{x} \in \mathbb{R}^n$ (see Definition 4.21), and $|\mathcal{I}(\mathbf{x})|$ to denote the cardinality of this set, i.e., the number of active inequality constraints at $\mathbf{x} \in \mathbb{R}^n$.

In order to compute approximations to the tangent cone $T_S(\mathbf{x})$, similarly to Example 4.22 we consider cones associated with the active constraints at a given point:

$$\overset{\circ}{G}(\mathbf{x}) := \{ \mathbf{p} \in \mathbb{R}^n \mid \nabla g_i(\mathbf{x})^\top \mathbf{p} < 0, \ i \in \mathcal{I}(\mathbf{x}) \}, \quad (5.6)$$

and

$$G(\mathbf{x}) := \{ \mathbf{p} \in \mathbb{R}^n \mid \nabla g_i(\mathbf{x})^\top \mathbf{p} \leq 0, \ i \in \mathcal{I}(\mathbf{x}) \}. \quad (5.7)$$

The following proposition verifies that $\overset{\circ}{G}(\mathbf{x})$ is an inner approximation for $R_S(\mathbf{x})$ (and, therefore, for $T_S(\mathbf{x})$ as well, see Proposition 5.3), and $G(\mathbf{x})$ is an outer approximation for $T_S(\mathbf{x})$.

Lemma 5.10 *For every $\mathbf{x} \in \mathbb{R}^n$ it holds that $\overset{\circ}{G}(\mathbf{x}) \subseteq R_S(\mathbf{x})$, and $T_S(\mathbf{x}) \subseteq G(\mathbf{x})$.*

Proof. Let $\mathbf{p} \in \overset{\circ}{G}(\mathbf{x})$. For every $i \notin \mathcal{I}(\mathbf{x})$ the function g_i is continuous and $g_i(\mathbf{x}) < 0$; therefore $g_i(\mathbf{x} + \delta \mathbf{p}) < 0$ for all small $\delta > 0$. Moreover, by Proposition 4.16, \mathbf{p} is a direction of descent for every g_i at \mathbf{x} , $i \in \mathcal{I}(\mathbf{x})$, which means that $g_i(\mathbf{x} + \delta \mathbf{p}) < g_i(\mathbf{x}) = 0$ for all such i and all small $\delta > 0$. Thus, $\mathbf{p} \in R_S(\mathbf{x})$, and, hence, $\overset{\circ}{G}(\mathbf{x}) \subseteq R_S(\mathbf{x})$.

Now, let $\mathbf{p} \in T_S(\mathbf{x})$, i.e., $\exists \{\mathbf{x}_k\} \subset S$, and $\{\lambda_k\} \subset (0, \infty)$ such that $\lim_{k \rightarrow \infty} \mathbf{x}_k = \mathbf{x}$ and $\lim_{k \rightarrow \infty} \lambda_k(\mathbf{x}_k - \mathbf{x}) = \mathbf{p}$. Exactly as in the proof of Theorem 5.8, we use the first order Taylor expansion (2.1) of the functions g_i , $i \in \mathcal{I}(\mathbf{x})$, to get:

$$0 \geq g_i(\mathbf{x}_k) = g_i(\mathbf{x}_k) - g_i(\mathbf{x}) = \nabla g_i(\mathbf{x})^\top (\mathbf{x}_k - \mathbf{x}) + o(\|\mathbf{x}_k - \mathbf{x}\|),$$

where the first inequality is by the feasibility of \mathbf{x}_k . Multiplying by $\lambda_k > 0$ and taking limit we get, for $i \in \mathcal{I}(\mathbf{x})$,

$$\begin{aligned} 0 &\geq \lim_{k \rightarrow \infty} \left[\lambda_k \nabla g_i(\mathbf{x})^\top (\mathbf{x}_k - \mathbf{x}) + \|\lambda_k(\mathbf{x}_k - \mathbf{x})\| \frac{o(\|\mathbf{x}_k - \mathbf{x}\|)}{\|\mathbf{x}_k - \mathbf{x}\|} \right] \\ &= \nabla g_i(\mathbf{x})^\top \mathbf{p} + \|\mathbf{p}\| \cdot 0, \end{aligned}$$

and thus $\mathbf{p} \in G(\mathbf{x})$. ■

Example 5.11 (Example 5.4 continued) The set S is defined by the two inequality constraints $g_1(\mathbf{x}) := -x_1 \leq 0$ and $g_2(\mathbf{x}) := (x_1 - 1)^2 + x_2^2 - 1 \leq 0$. Let us calculate $\overset{\circ}{G}(\mathbf{0}^2)$ and $G(\mathbf{0}^2)$. Both constraints are satisfied with equality at the given point, so that $\mathcal{I}(\mathbf{0}^2) = \{1, 2\}$. Then, $\nabla g_1(\mathbf{0}^2) = (-1, 0)^T$, $\nabla g_2(\mathbf{0}^2) = (-2, 0)^T$, and thus $\overset{\circ}{G}(\mathbf{0}^2) = \{\mathbf{x} \in \mathbb{R}^2 \mid x_1 > 0\} = R_S(\mathbf{0}^2)$, $G(\mathbf{0}^2) = \{\mathbf{x} \in \mathbb{R}^2 \mid x_1 \geq 0\} = T_S(\mathbf{0}^2)$ in this case. ■

Example 5.12 (Example 5.5 continued) The set S is defined by the three inequality constraints $g_1(\mathbf{x}) := -x_1 \leq 0$, $g_2(\mathbf{x}) := -x_2 \leq 0$, $g_3(\mathbf{x}) := x_1 x_2 \leq 0$, which are all active at $\mathbf{x} = \mathbf{0}^2$; $\nabla g_1(\mathbf{0}^2) = (-1, 0)^T$, $\nabla g_2(\mathbf{0}^2) = (0, -1)^T$, and $\nabla g_3(\mathbf{0}^2) = (0, 0)^T$. Therefore, $\overset{\circ}{G}(\mathbf{0}^2) = \emptyset \subsetneq R_S(\mathbf{0}^2)$, and $G(\mathbf{0}^2) = \{\mathbf{x} \in \mathbb{R}^2 \mid x_1 \geq 0, x_2 \geq 0\} \supsetneq T_S(\mathbf{0}^2)$. ■

Example 5.13 (Example 5.6 continued) The set S is defined by the three inequality constraints $g_1(\mathbf{x}) := -x_1^3 + x_2 \leq 0$, $g_2(\mathbf{x}) := x_1^5 - x_2 \leq 0$, $g_3(\mathbf{x}) := -x_2 \leq 0$, which are all active at $\mathbf{x} = \mathbf{0}^2$; $\nabla g_1(\mathbf{0}^2) = (0, 1)^T$, $\nabla g_2(\mathbf{0}^2) = (0, -1)^T$, and $\nabla g_3(\mathbf{0}^2) = (0, -1)^T$. Therefore, $\overset{\circ}{G}(\mathbf{0}^2) = \emptyset = R_S(\mathbf{0}^2)$, and $G(\mathbf{0}^2) = \{\mathbf{x} \in \mathbb{R}^2 \mid x_2 = 0\} \supsetneq T_S(\mathbf{0}^2)$. ■

Example 5.14 (Example 5.7 continued) The set S is defined by the inequality constraint $g_1(\mathbf{x}) := -x_2 \leq 0$, and the equality constraint $h_1(\mathbf{x}) := (x_1 - 1)^2 + x_2^2 - 1 = 0$; we split the latter into two inequality constraints $g_2(\mathbf{x}) := h_1(\mathbf{x}) \leq 0$, and $g_3(\mathbf{x}) := -h_1(\mathbf{x}) \leq 0$. Thus, we end up with three active inequality constraints at $\mathbf{x} = \mathbf{0}^2$; $\nabla g_1(\mathbf{0}^2) = (0, -1)^T$, $\nabla g_2(\mathbf{0}^2) = (-2, 0)^T$, and $\nabla g_3(\mathbf{0}^2) = (2, 0)^T$. Therefore, $\overset{\circ}{G}(\mathbf{0}^2) = \emptyset = R_S(\mathbf{0}^2)$, and $G(\mathbf{0}^2) = \{\mathbf{x} \in \mathbb{R}^2 \mid x_1 = 0, x_2 \geq 0\} = T_S(\mathbf{0}^2)$. ■

Now we are ready to establish the Fritz John optimality conditions.

Theorem 5.15 (Fritz John necessary optimality conditions) *Let the set S be defined by (5.5). If $\mathbf{x}^* \in S$ is a local minimum of f over S then there exist multipliers $\mu_0 \in \mathbb{R}$, $\boldsymbol{\mu} \in \mathbb{R}^m$, such that*

$$\mu_0 \nabla f(\mathbf{x}^*) + \sum_{i=1}^m \mu_i \nabla g_i(\mathbf{x}^*) = \mathbf{0}^n, \quad (5.8a)$$

$$\mu_i g_i(\mathbf{x}^*) = 0, \quad i = 1, \dots, m, \quad (5.8b)$$

$$\mu_0, \mu_i \geq 0, \quad i = 1, \dots, m, \quad (5.8c)$$

$$(\mu_0, \boldsymbol{\mu}^T)^T \neq \mathbf{0}^{m+1}. \quad (5.8d)$$

Optimality conditions

In other words,

$$\mathbf{x}^* \text{ local minimum of } f \text{ over } S \implies \exists(\mu_0, \boldsymbol{\mu}) \in \mathbb{R} \times \mathbb{R}^m : (5.8) \text{ holds.}$$

Proof. Combining the results of Lemma 5.10 with the geometric optimality conditions provided by Theorem 5.8, we conclude that there is no direction $\mathbf{p} \in \mathbb{R}^n$ such that $\nabla f(\mathbf{x}^*)^T \mathbf{p} < 0$ and $\nabla g_i(\mathbf{x}^*)^T \mathbf{p} < 0$, $i \in \mathcal{I}(\mathbf{x}^*)$. Define the matrix \mathbf{A} with columns $\nabla f(\mathbf{x}^*)$, $\nabla g_i(\mathbf{x}^*)$, $i \in \mathcal{I}(\mathbf{x}^*)$; then the system $\mathbf{A}^T \mathbf{p} < \mathbf{0}^{1+|\mathcal{I}(\mathbf{x}^*)|}$ is unsolvable. By Farkas' Lemma (cf. Theorem 3.30) there exists a nonzero vector $\boldsymbol{\lambda} \in \mathbb{R}^{1+|\mathcal{I}(\mathbf{x}^*)|}$ such that $\boldsymbol{\lambda} \geq \mathbf{0}^{1+|\mathcal{I}(\mathbf{x}^*)|}$ and $\mathbf{A}\boldsymbol{\lambda} = \mathbf{0}^n$. (Why?) Now, let $(\mu_0, \boldsymbol{\mu}_{\mathcal{I}(\mathbf{x}^*)}^T)^T := \boldsymbol{\lambda}$, and set $\mu_i = 0$ for $i \notin \mathcal{I}(\mathbf{x}^*)$. It is an easy exercise to verify that so defined μ_0 and $\boldsymbol{\mu}$ satisfy the conditions (5.8). ■

Remark 5.16 (terminology) The solutions $(\mu_0, \boldsymbol{\mu})$ to the system (5.8) are known as *Lagrange multipliers* (or just *multipliers*) associated with a given candidate $\mathbf{x}^* \in \mathbb{R}^n$ for a local minimum. Note that every multiplier (except μ_0) corresponds to some constraint in the algebraic representation of S . The conditions (5.8a) and (5.8c) are known as the *dual feasibility* conditions, and (5.8b) as the *complementarity* conditions, respectively; this terminology will become more clear in Chapter 6. Owing to the complementarity constraints, the multipliers μ_i corresponding to *inactive* inequality constraints $i \notin \mathcal{I}(\mathbf{x}^*)$ must be zero. In general, the Lagrange multiplier μ_i bears the important information about how *sensitive* a particular local minimum is with respect to small changes in the constraint g_i . ■

In the following examples, as before, we assume that $f(\mathbf{x}) := x_1$, so that $\nabla f = (1, 0)^T$ and $\mathbf{x}^* = \mathbf{0}^2$ is the point of local minimum.

Example 5.17 (Example 5.4 continued) The Fritz John system (5.8) at the point $\mathbf{x}^* = \mathbf{0}^2$ reduces to:

$$\begin{aligned} \mu_0 \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \begin{pmatrix} -1 & -2 \\ 0 & 0 \end{pmatrix} \boldsymbol{\mu} &= \mathbf{0}^2, \\ (\mu_0, \boldsymbol{\mu}^T)^T &\succeq \mathbf{0}^3, \end{aligned}$$

where $\boldsymbol{\mu} \in \mathbb{R}^2$ is a vector of Lagrange multipliers for the inequality constraints. We do not write the complementarity constraints (5.8b), because in our case the two constraints are active, and therefore the equation (5.8b) is automatically satisfied for all $\boldsymbol{\mu}$. The solutions to this

system are the pairs $(\mu_0, \boldsymbol{\mu})$, with $\boldsymbol{\mu} = (\mu_1, 2^{-1}(\mu_0 - \mu_1))^T$, for every $\mu_0 > 0$, $0 \leq \mu_1 \leq \mu_0$. There are infinitely many Lagrange multipliers, that form an unbounded set, but μ_0 must always be positive. ■

Example 5.18 (Example 5.5 continued) Similarly to the previous example, the Fritz John system (5.8) at the point $\mathbf{x}^* = \mathbf{0}^2$ reduces to:

$$\begin{aligned} \mu_0 \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \begin{pmatrix} -1 & 0 & 0 \\ 0 & -1 & 0 \end{pmatrix} \boldsymbol{\mu} &= \mathbf{0}^2, \\ (\mu_0, \boldsymbol{\mu}^T)^T &\succeq \mathbf{0}^4, \end{aligned}$$

where $\boldsymbol{\mu} \in \mathbb{R}^3$ is a vector of Lagrange multipliers for the inequality constraints. The solution to the Fritz John system is every pair $(\mu_0, \boldsymbol{\mu})$ with $\boldsymbol{\mu} = (\mu_0, 0, \mu_3)^T$ for every $\mu_0 \geq 0$, $\mu_3 \geq 0$ such that either of them is strictly bigger than zero. That is, there are infinitely many Lagrange multipliers, that form an unbounded set, and it is possible for μ_0 to assume the value zero. ■

Example 5.19 (Example 5.6 continued) The Fritz John system (5.8) at the point $\mathbf{x}^* = \mathbf{0}^2$ reduces to:

$$\begin{aligned} \mu_0 \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \begin{pmatrix} 0 & 0 & 0 \\ 1 & -1 & -1 \end{pmatrix} \boldsymbol{\mu} &= \mathbf{0}^2, \\ (\mu_0, \boldsymbol{\mu}^T)^T &\succeq \mathbf{0}^4, \end{aligned}$$

where $\boldsymbol{\mu} \in \mathbb{R}^3$ is a vector of Lagrange multipliers for the inequality constraints. Thus, $\mu_0 = 0$, $\boldsymbol{\mu} = (\mu_1, \mu_2, \mu_1 - \mu_2)^T$ for every $\mu_1 > 0$, $0 \leq \mu_2 \leq \mu_1$. That is, there are infinitely many Lagrange multipliers, that form an unbounded set, and μ_0 must assume the value zero. ■

Example 5.20 (Example 5.7 continued) The Fritz John system (5.8) at the point $\mathbf{x}^* = \mathbf{0}^2$ reduces to:

$$\begin{aligned} \mu_0 \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \begin{pmatrix} 0 & -2 & 2 \\ -1 & 0 & 0 \end{pmatrix} \boldsymbol{\mu} &= \mathbf{0}^2, \\ (\mu_0, \boldsymbol{\mu}^T)^T &\succeq \mathbf{0}^4, \end{aligned}$$

where $\boldsymbol{\mu} \in \mathbb{R}^3$ is a vector of Lagrange multipliers for the inequality constraints. The solution to the Fritz John system is every pair $(\mu_0, \boldsymbol{\mu})$ with $\boldsymbol{\mu} = (0, \mu_2, \mu_2 - 2^{-1}\mu_0)^T$ for every $\mu_2 > 0$, $0 \leq \mu_0 \leq 2\mu_2$. That is, there are infinitely many Lagrange multipliers, that form an unbounded set, and it is possible for μ_0 to assume the value zero. ■

The fact that μ_0 may be zero in the system (5.8) essentially means that the objective function f plays no role in the optimality conditions. This is of course a rather unexpected and unwanted situation, and the rest of the chapter is dedicated to describing how one can avoid it.

Since the cone of feasible directions $R_S(\mathbf{x})$ may be a bad approximation of the tangent cone $T_S(\mathbf{x})$, so may $\overset{\circ}{G}(\mathbf{x})$ owing to Lemma 5.10. Therefore, in the most general case we cannot improve on the conditions (5.8); however, it is possible to improve upon (5.8) if we assume that the set S is “regular” in some sense, i.e., that either $\overset{\circ}{G}(\mathbf{x})$ or $G(\mathbf{x})$ is a tight enough approximation of $T_S(\mathbf{x})$. Requirements of this type are called *constraint qualifications*, and they will be discussed in more detail in Section 5.7. However, to get a feeling of what can be achieved with a regular constraint set S , we show that the multiplier μ_0 in the system (5.8) cannot vanish (i.e., the KKT conditions hold, see Section 5.5) if the constraint qualification $\overset{\circ}{G}(\mathbf{x}^*) \neq \emptyset$ holds (which is quite a restrictive one, in view of Example 5.22; however, see the much weaker assumption denoted MFCQ in Definition 5.35).

Proposition 5.21 (KKT optimality conditions, preview) *Assume the conditions of Theorem 5.8, and assume that $\overset{\circ}{G}(\mathbf{x}^*) \neq \emptyset$. Then, the multiplier μ_0 in (5.8) cannot be zero; dividing all equations by μ_0 we may assume that it equals one.*

Proof. Assume that $\mu_0 = 0$ in (5.8), and define the matrix \mathbf{A} with columns $\nabla g_i(\mathbf{x}^*)$, $i \in \mathcal{I}(\mathbf{x}^*)$. Since $\mathbf{A}\boldsymbol{\mu} = \mathbf{0}^n$, $\boldsymbol{\mu} \geq \mathbf{0}^{|\mathcal{I}(\mathbf{x}^*)|}$, and $\boldsymbol{\mu} \neq \mathbf{0}^{|\mathcal{I}(\mathbf{x}^*)|}$, the system $\mathbf{A}^T \mathbf{p} < \mathbf{0}^{|\mathcal{I}(\mathbf{x}^*)|}$ is unsolvable (see Farkas’ Lemma, Theorem 3.30), i.e., $\overset{\circ}{G}(\mathbf{x}^*) = \emptyset$. ■

Example 5.22 Out of the four Examples 5.4–5.7, only the first verifies the condition $\overset{\circ}{G}(\mathbf{x}^*) \neq \emptyset$ assumed in Proposition 5.21, while as we will see later (and as Examples 5.17–5.20 may suggest), three out of the four problems admit solutions to the corresponding KKT systems. ■

5.5 The Karush–Kuhn–Tucker conditions

In this section we develop the famous and classic Karush–Kuhn–Tucker optimality conditions for constrained optimization problems with inequality constraints, which are essentially the Fritz John conditions (5.8) with the additional requirement $\mu_0 \neq 0$ (in fact, $\mu_0 = 1$). We establish these conditions as before, for inequality constrained problems (5.5)

(which we do without any loss of generality or sharpness of the theory), and then discuss the possible modifications of the conditions if one wants to specifically distinguish between equality and inequality constraints in Section 5.6. Abadie’s *constraint qualification* (see Definition 5.23) which we impose is very abstract and extremely general (this is *almost* the weakest condition one can require); of course it is impossible to check it when it comes to practical problems. Therefore, in Section 5.7 we list some computationally verifiable assumptions that all imply Abadie’s constraint qualification.

We start with a formal definition.

Definition 5.23 (Abadie’s constraint qualification) *We say that at the point $\mathbf{x} \in S$ Abadie’s constraint qualification holds if $T_S(\mathbf{x}) = G(\mathbf{x})$, where $T_S(\mathbf{x})$ is defined by Definition 5.2 and $G(\mathbf{x})$ by (5.7). ■*

Example 5.24 Out of the four Examples 5.4–5.7, the first and the last satisfy Abadie’s constraint qualification (see Examples 5.11–5.14). ■

Then, we are ready to prove the main theorem in this chapter.

Theorem 5.25 (Karush–Kuhn–Tucker optimality conditions) *Assume that at a given point $\mathbf{x}^* \in S$ Abadie’s constraint qualification holds. If $\mathbf{x}^* \in S$ is a local minimum of f over S then there exists a vector $\boldsymbol{\mu} \in \mathbb{R}^m$ such that*

$$\nabla f(\mathbf{x}^*) + \sum_{i=1}^m \mu_i \nabla g_i(\mathbf{x}^*) = \mathbf{0}^n, \quad (5.9a)$$

$$\mu_i g_i(\mathbf{x}^*) = 0, \quad i = 1, \dots, m, \quad (5.9b)$$

$$\boldsymbol{\mu} \geq \mathbf{0}^m. \quad (5.9c)$$

In other words,

$$\left. \begin{array}{l} \mathbf{x}^* \text{ local minimum of } f \text{ over } S \\ \text{Abadie's CQ holds at } \mathbf{x}^* \end{array} \right\} \implies \exists \boldsymbol{\mu} \in \mathbb{R}^m : (5.9) \text{ holds.}$$

The system will be referred to as the *Karush–Kuhn–Tucker optimality conditions*.

Proof. By Theorem 5.8 we have that $\overset{\circ}{F}(\mathbf{x}^*) \cap T_S(\mathbf{x}^*) = \emptyset$, which due to our assumptions implies that $\overset{\circ}{F}(\mathbf{x}^*) \cap G(\mathbf{x}^*) = \emptyset$.

As in the proof of Theorem 5.15, construct a matrix \mathbf{A} with columns $\nabla g_i(\mathbf{x}^*)$, $i \in \mathcal{I}(\mathbf{x}^*)$. Then, the system $\mathbf{A}^T \mathbf{p} \leq \mathbf{0}^{|\mathcal{I}(\mathbf{x}^*)|}$ and $-\nabla f(\mathbf{x}^*)^T \mathbf{p} >$

Optimality conditions

0 has no solutions. By Farkas' Lemma (cf. Theorem 3.30), the system $A\xi = -\nabla f(\mathbf{x}^*)$, $\xi \geq \mathbf{0}^{|\mathcal{I}(\mathbf{x}^*)|}$ has a solution. Define the vector $\mu_{\mathcal{I}(\mathbf{x}^*)} = \xi$, and $\mu_i = 0$, for $i \notin \mathcal{I}(\mathbf{x}^*)$. Then, the so defined μ verifies the KKT conditions (5.9). ■

Remark 5.26 (terminology) Similarly to the case of the Fritz John necessary optimality conditions, the solutions μ to the system (5.9) are known as *Lagrange multipliers* (or just *multipliers*) associated with a given candidate $\mathbf{x}^* \in \mathbb{R}^n$ for a local minimum. The conditions (5.9a) and (5.9c) are known as the *dual feasibility* conditions, and (5.9b) as the *complementarity* conditions, respectively; this terminology will become more clear in Chapter 6. Owing to the complementarity constraints, the multipliers μ_i corresponding to *inactive* inequality constraints $i \notin \mathcal{I}(\mathbf{x}^*)$ must be zero. In general, the Lagrange multiplier μ_i bears the important information about how *sensitive* a particular local minimum is with respect to small changes in the constraint g_i . ■

Remark 5.27 (geometric interpretation) The system of equations and inequalities defining (5.9) can (and should) be interpreted geometrically as $-\nabla f(\mathbf{x}^*) \in N_S(\mathbf{x}^*)$ (see Figure 5.6), the latter cone being the *normal cone* to S at $\mathbf{x}^* \in S$ (see Definition 4.25); according to the figure, the normal cone to S at \mathbf{x}^* is furthermore spanned by the gradients of the active constraints at \mathbf{x}^* .¹

Notice the specific roles played by the different parts of the system (5.9) in this respect: the complementarity conditions (5.9b) force μ_i to be equal to 0 for the inactive constraints, whence the summation in the left-hand side of the linear system (5.9a) involves the active constraints only. Further, the sign conditions in (5.9c) ensures that each vector $\mu_i \nabla g_i(\mathbf{x}^*)$, $i \in \mathcal{I}(\mathbf{x}^*)$, is an *outward normal* to S at \mathbf{x}^* . ■

Remark 5.28 Note that in the unconstrained case the KKT system (5.9) reduces to the single requirement $\nabla f(\mathbf{x}^*) = \mathbf{0}^n$, which we have already encountered in Theorem 4.14.

It is possible to further develop the KKT theory (with some technical complications) for twice differentiable functions as it has been done for the unconstrained case in Theorem 4.17. We refer the interested reader to [BSS93, Section 4.4]. ■

¹Compare with the normal cone characterization (4.17) and Figure 4.4 in the case of convex feasible sets: we could, roughly, say that the role of a constraint qualification in the more general context of this chapter is to ensure that the normal cone to the feasible set at the vector \mathbf{x}^* is a *finitely generated convex cone*, which moreover is generated by the gradients of the active constraints' describing functions g_i at \mathbf{x}^* , thus extending the normal cone inclusion in (4.17) to more general sets.

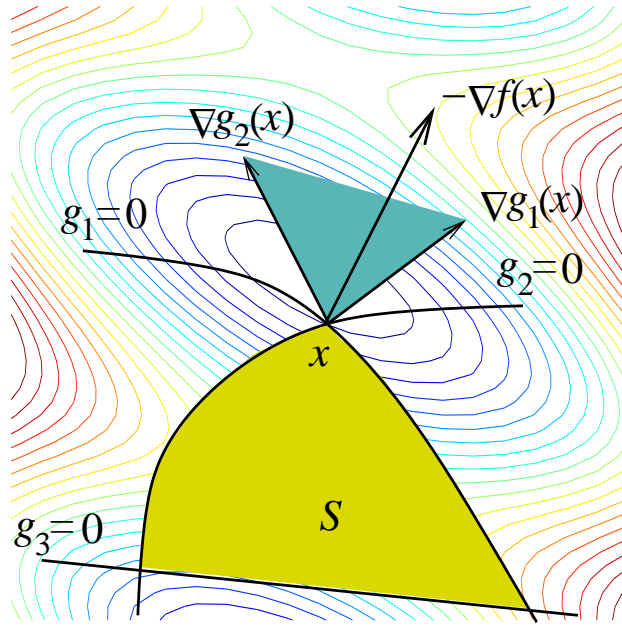


Figure 5.6: Geometrical interpretation of the KKT system.

Example 5.29 (Example 5.4 continued) In this example Abadie's constraint qualification is fulfilled; the KKT system must be solvable. Indeed, the system

$$\begin{pmatrix} 1 \\ 0 \end{pmatrix} + \begin{pmatrix} -1 & -2 \\ 0 & 0 \end{pmatrix} \mu = \mathbf{0}^2, \\ \mu \geq \mathbf{0}^2,$$

possesses solutions $\mu = (\mu_1, 2^{-1}(1 - \mu_1))^T$ for every $0 \leq \mu_1 \leq 1$. Therefore, there are infinitely many multipliers, that all belong to a bounded set. ■

Example 5.30 (Example 5.5 continued) This is one of the rare cases when Abadie's constraint qualification is violated, and nevertheless the KKT system happens to be solvable:

$$\begin{pmatrix} 1 \\ 0 \end{pmatrix} + \begin{pmatrix} -1 & 0 & 0 \\ 0 & -1 & 0 \end{pmatrix} \mu = \mathbf{0}^2, \\ \mu \geq \mathbf{0}^3,$$

Optimality conditions

admits solutions $\boldsymbol{\mu} = (1, 0, \mu_3)^\top$ for every $\mu_3 \geq 0$. That is, the set of Lagrange multipliers is unbounded in this case. ■

Example 5.31 (Example 5.6 continued) Since, for this example, in the Fritz John system the multiplier μ_0 is necessarily zero, the KKT system admits no solutions:

$$\begin{pmatrix} 1 \\ 0 \end{pmatrix} + \begin{pmatrix} 0 & 0 & 0 \\ 1 & -1 & -1 \end{pmatrix} \boldsymbol{\mu} = \mathbf{0}^2, \\ \boldsymbol{\mu} \geq \mathbf{0}^3,$$

is clearly inconsistent. In this example Abadie's constraint qualification is violated. ■

Example 5.32 (Example 5.7 continued) This example satisfies Abadie's constraint qualification, and therefore, since a global optimum exists, the KKT system is solvable:

$$\begin{pmatrix} 1 \\ 0 \end{pmatrix} + \begin{pmatrix} 0 & -2 & 2 \\ -1 & 0 & 0 \end{pmatrix} \boldsymbol{\mu} = \mathbf{0}^2, \\ \boldsymbol{\mu} \geq \mathbf{0}^3,$$

admits the solutions $\boldsymbol{\mu} = (0, \mu_2, \mu_2 - 2^{-1})^\top$, for all $\mu_2 \geq 2^{-1}$. The set of Lagrange multipliers is unbounded in this case, but this is because we have split the original equality constraint into two inequalities. In Section 5.6 we formulate the KKT system that keeps the original equality-representation of the set, and thus reduce the number of multipliers for the equality constraint to just one! ■

5.6 Proper treatment of equality constraints

Now we consider both inequality and equality constraints, that is, we assume that the feasible set S is given by

$$S := \{ \mathbf{x} \in \mathbb{R}^n \mid g_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m; \\ h_j(\mathbf{x}) = 0, \quad j = 1, \dots, \ell \}, \quad (5.10)$$

instead of (5.5), where $g_i \in C^1(\mathbb{R}^n)$, $i = 1, \dots, m$, and $h_j \in C^1(\mathbb{R}^n)$, $j = 1, \dots, \ell$. As it was done in Section 5.4, we write S using only inequality constraints, by defining the functions $\tilde{g}_i \in C^1(\mathbb{R}^n)$, $i = 1, \dots, m + 2\ell$, via:

$$\tilde{g}_i := \begin{cases} g_i, & i = 1, \dots, m, \\ h_{i-m}, & i = m + 1, \dots, m + \ell, \\ -h_{i-m-\ell}, & i = m + \ell + 1, \dots, m + 2\ell, \end{cases} \quad (5.11)$$

so that

$$S = \{ \mathbf{x} \in \mathbb{R}^n \mid \tilde{g}_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m + 2\ell \}. \quad (5.12)$$

Now, let $\tilde{G}(\mathbf{x})$ be defined by (5.7) for the inequality representation (5.12) of S . We will use the old notation $G(\mathbf{x})$ for the cone defined only by the gradients of the functions defining the *inequality* constraints active at \mathbf{x} in the representation (5.10), and in addition define the null space of the matrix defined by the gradients of the functions defining the *equality* constraints:

$$H(\mathbf{x}) := \{ \mathbf{p} \in \mathbb{R}^n \mid \nabla h_i(\mathbf{x})^T \mathbf{p} = 0, \quad i = 1, \dots, \ell \}. \quad (5.13)$$

Since all inequality constraint functions \tilde{g}_i , $i = m + 1, \dots, m + 2\ell$, are necessarily active at every $\mathbf{x} \in S$, it holds that

$$\tilde{G}(\mathbf{x}) = G(\mathbf{x}) \cap H(\mathbf{x}), \quad (5.14)$$

and thus Abadie's constraint qualification (see Definition 5.23) for the set (5.10) may be equivalently written as

$$T_S(\mathbf{x}) = G(\mathbf{x}) \cap H(\mathbf{x}). \quad (5.15)$$

Assuming that the latter constraint qualification holds we can write the KKT system (5.9) for $\mathbf{x}^* \in S$, corresponding to the inequality representation (5.12) (see Theorem 5.25):

$$\begin{aligned} \sum_{i=1}^m \mu_i \nabla g_i(\mathbf{x}^*) + \sum_{i=m+1}^{m+\ell} \mu_i \nabla h_{i-m}(\mathbf{x}^*) - \sum_{i=m+\ell+1}^{m+2\ell} \mu_i \nabla h_{i-m-\ell}(\mathbf{x}^*) \\ + \nabla f(\mathbf{x}^*) = \mathbf{0}^n, \end{aligned} \quad (5.16a)$$

$$\mu_i g_i(\mathbf{x}^*) = 0, \quad i = 1, \dots, m, \quad (5.16b)$$

$$\mu_i h_{i-m}(\mathbf{x}^*) = 0, \quad i = m + 1, \dots, m + \ell, \quad (5.16c)$$

$$-\mu_i h_{i-m-\ell}(\mathbf{x}^*) = 0, \quad i = m + \ell + 1, \dots, m + 2\ell, \quad (5.16d)$$

$$\boldsymbol{\mu} \geq \mathbf{0}^{m+2\ell}. \quad (5.16e)$$

Define the pair of vectors $(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\lambda}}) \in \mathbb{R}^m \times \mathbb{R}^\ell$ as $\tilde{\mu}_i = \mu_i$, $i = 1, \dots, m$; $\tilde{\lambda}_j = \mu_{m+j} - \mu_{m+\ell+j}$, $j = 1, \dots, \ell$. We also note that the equations (5.16c) and (5.16d) are superfluous, because $\mathbf{x}^* \in S$ implies that $h_j(\mathbf{x}^*) = 0$, $j = 1, \dots, \ell$. Therefore, we get the following system for $(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\lambda}})$, known as the KKT necessary optimality conditions for the sets represented by

Optimality conditions

differentiable equality and inequality constraints:

$$\nabla f(\mathbf{x}^*) + \sum_{i=1}^m \tilde{\mu}_i \nabla g_i(\mathbf{x}^*) + \sum_{j=1}^{\ell} \tilde{\lambda}_j \nabla h_j(\mathbf{x}^*) = \mathbf{0}^n, \quad (5.17a)$$

$$\tilde{\mu}_i g_i(\mathbf{x}^*) = 0, \quad i = 1, \dots, m, \quad (5.17b)$$

$$\tilde{\boldsymbol{\mu}} \geq \mathbf{0}^m. \quad (5.17c)$$

Thus, we have established the following theorem.

Theorem 5.33 (KKT optimality conditions for inequality and equality constraints) *Assume that at a given point $\mathbf{x}^* \in S$ Abadie's constraint qualification (5.15) holds, where S is given by (5.10). If \mathbf{x}^* is a local minimum of a differentiable function f over S then there exists a pair of vectors $(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\lambda}}) \in \mathbb{R}^m \times \mathbb{R}^{\ell}$ such that the system (5.17) is satisfied.*

In other words,

$$\left. \begin{array}{l} \mathbf{x}^* \text{ local minimum of } f \text{ over } S \\ \text{Abadie's CQ holds at } \mathbf{x}^* \end{array} \right\} \implies \exists (\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\lambda}}) \in \mathbb{R}^m \times \mathbb{R}^{\ell} : (5.17) \text{ holds.}$$

■

Example 5.34 (Example 5.32 revisited) Let us write the system of KKT conditions for the original representation of the set with one inequality and one equality constraint (see Example 5.14). As has already been mentioned, Abadie's constraint qualification is satisfied, and therefore, since an optimum exists, the KKT system is necessarily solvable:

$$\begin{pmatrix} 1 \\ 0 \end{pmatrix} + \mu_1 \begin{pmatrix} 0 \\ -1 \end{pmatrix} + \lambda_1 \begin{pmatrix} -2 \\ 0 \end{pmatrix} = \mathbf{0}^2, \\ \mu_1 \geq 0,$$

which admits the unique solution $\mu_1 = 0$, $\lambda_1 = 1/2$.

■

5.7 Constraint qualifications

In this section we discuss conditions on the functions involved in the representation (5.10) of a given feasible set S , that all imply Abadie's constraint qualification (5.15).

5.7.1 Mangasarian–Fromovitz CQ (MFCQ)

Definition 5.35 (Mangasarian–Fromovitz CQ) *We say that at the point $\mathbf{x} \in S$, where S is given by (5.10), the Mangasarian–Fromovitz CQ holds if the gradients $\nabla h_j(\mathbf{x})$ of the functions h_j , $j = 1, \dots, \ell$, defining the equality constraints, are linearly independent, and the intersection $\overset{\circ}{G}(\mathbf{x}) \cap H(\mathbf{x})$ is nonempty.* ■

We state the following result without a “real” proof, but we outline the ideas.

Proposition 5.36 *The MFCQ implies Abadie’s CQ.*

Proof. [Sketch] Since the gradients $\nabla h_j(\mathbf{x})$, $j = 1, \dots, \ell$, are linearly independent, it can be shown that $\text{cl}(\overset{\circ}{G}(\mathbf{x}) \cap H(\mathbf{x})) \subseteq T_S(\mathbf{x})$ (in the absence of equality constraints, it follows directly from Lemma 5.10).

Furthermore, from Lemma 5.10 applied to the inequality representation of S , i.e., to $\tilde{G}(\mathbf{x})$ defined by (5.14), we know that $T_S(\mathbf{x}) \subseteq (G(\mathbf{x}) \cap H(\mathbf{x}))$.

Finally, since $\overset{\circ}{G}(\mathbf{x}) \cap H(\mathbf{x}) \neq \emptyset$, it can be shown that $\text{cl}(\overset{\circ}{G}(\mathbf{x}) \cap H(\mathbf{x})) = G(\mathbf{x}) \cap H(\mathbf{x})$. ■

Example 5.37 Since MFCQ implies Abadie’s constraint qualification, Example 5.5 and 5.6 must necessarily violate it. On the other hand, both Examples 5.4 and 5.7 verify it (since they also satisfy stronger constraint qualifications, see Example 5.40 and 5.43). ■

5.7.2 Slater CQ

Definition 5.38 (Slater CQ) *We say that the system of constraints describing the feasible set S via (5.10) satisfies the Slater CQ, if the functions g_i , $i = 1, \dots, m$, defining the inequality constraints are convex, the functions h_j , $j = 1, \dots, \ell$, defining the equality constraints are affine with linearly independent gradients $\nabla h_j(\mathbf{x})$, $j = 1, \dots, \ell$, and, finally, that there exists $\bar{\mathbf{x}} \in S$ such that $g_i(\bar{\mathbf{x}}) < 0$, for all $i \in \{1, \dots, m\}$.* ■

Proposition 5.39 *The Slater CQ implies the MFCQ.*

Proof. Suppose the Slater CQ holds at $\mathbf{x} \in S$. By the convexity of the inequality constraints we get:

$$0 > g_i(\bar{\mathbf{x}}) = g_i(\bar{\mathbf{x}}) - g_i(\mathbf{x}) \geq \nabla g_i(\mathbf{x})^T(\bar{\mathbf{x}} - \mathbf{x}),$$

Optimality conditions

for all $i \in \mathcal{I}(\mathbf{x})$. Furthermore, since the equality constraints are affine, we have that

$$0 = h_j(\bar{\mathbf{x}}) - h_j(\mathbf{x}) = \nabla h_j(\mathbf{x})^T(\bar{\mathbf{x}} - \mathbf{x}),$$

$j = 1, \dots, \ell$. Then, $\bar{\mathbf{x}} - \mathbf{x} \in \overset{\circ}{G}(\mathbf{x}) \cap H(\mathbf{x})$. ■

Example 5.40 Only Example 5.4 verifies the Slater CQ (which in particular explains why it satisfies MFCQ as well, see Example 5.37). ■

5.7.3 Linear independence CQ (LICQ)

Definition 5.41 (LICQ) We say that at the point $\mathbf{x} \in S$, where S is given by (5.10), the linear independence CQ holds if the gradients $\nabla g_i(\mathbf{x})$ of the functions g_i , $i \in \mathcal{I}(\mathbf{x})$, defining the active inequality constraints, as well as the gradients $\nabla h_j(\mathbf{x})$ of the functions h_j , $j = 1, \dots, \ell$, defining the equality constraints, are linearly independent. ■

Proposition 5.42 The LICQ implies the MFCQ.

Proof. [Sketch] Assume that $\overset{\circ}{G}(\mathbf{x}^*) \cap H(\mathbf{x}^*) = \emptyset$, i.e., the system $\mathbf{G}^T \mathbf{p} < \mathbf{0}^{|\mathcal{I}(\mathbf{x}^*)|}$ and $\mathbf{H}^T \mathbf{p} = \mathbf{0}^\ell$ is unsolvable, where \mathbf{G} and \mathbf{H} are the matrices having the gradients of the active inequality and equality constraints, respectively, as their columns. Using a separation result similar to Farkas' Lemma (cf. Theorem 3.30) one can show that the system $\mathbf{G}\boldsymbol{\mu} + \mathbf{H}\boldsymbol{\lambda} = \mathbf{0}^n$, $\boldsymbol{\mu} \geq \mathbf{0}^{|\mathcal{I}(\mathbf{x}^*)|}$ has a nonzero solution $(\boldsymbol{\mu}^T, \boldsymbol{\lambda}^T)^T \in \mathbb{R}^{|\mathcal{I}(\mathbf{x}^*)| + \ell}$, which contradicts the linear independence assumption. ■

In fact, the solution $(\boldsymbol{\mu}, \boldsymbol{\lambda})$ to the KKT system (5.17), if one exists, is necessarily unique in this case, and therefore LICQ is a rather strong assumption in many practical situations.

Example 5.43 Only Example 5.7 in the original description using both inequality and equality constraints verifies the LICQ (which in particular explains why it satisfies the MFCQ, see Example 5.37, and why the Lagrange multipliers are unique in this case, see Example 5.34). ■

5.7.4 Affine constraints

Assume that both the functions g_i , $i = 1, \dots, m$, defining the inequality constraints and the functions h_j , $j = 1, \dots, \ell$, defining the equality constraints in the representation (5.10) are affine, that is, the feasible

set S is a polyhedron. Then, the radial cone $R_S(\mathbf{x})$ (see Definition 5.1) is equal to $G(\mathbf{x}) \cap H(\mathbf{x})$ (see Example 4.22). Owing to the inclusions $R_S(\mathbf{x}) \subseteq T_S(\mathbf{x})$ (Proposition 5.3) and $T_S(\mathbf{x}) \subseteq \tilde{G}(\mathbf{x}) = G(\mathbf{x}) \cap H(\mathbf{x})$ (Lemma 5.10), where $\tilde{G}(\mathbf{x})$ was defined in Section 5.6 (cf. (5.12) and the discussion thereafter), Abadie's CQ (5.15) holds in this case.

Thus, the following claim is established.

Proposition 5.44 *If all (inequality and equality) constraints are affine, then Abadie's CQ is satisfied.* ■

5.8 Sufficiency of the KKT conditions under convexity

In general, the KKT necessary conditions do not imply local optimality, as has been mentioned before (see, e.g., the example right after the proof of Theorem 4.14). However, if the optimization problem (5.1) is convex, then the KKT conditions are *sufficient* for global optimality.

Theorem 5.45 (sufficiency of the KKT conditions for convex problems) *Assume that the problem (5.1) with the feasible set S given by (5.10) is convex, i.e., the objective function f as well as the functions g_i , $i = 1, \dots, m$, are convex, and the functions h_j , $j = 1, \dots, \ell$, are affine. Assume further that for $\mathbf{x}^* \in S$ the KKT conditions (5.17) are satisfied. Then, \mathbf{x}^* is a globally optimal solution of the problem (5.1).*

In other words,

$$\left. \begin{array}{l} \text{the problem (5.1) is convex} \\ \text{KKT conditions (5.17) hold at } \mathbf{x}^* \end{array} \right\} \implies \mathbf{x}^* \text{ global minimum in (5.1).}$$

Proof. Choose an arbitrary $\mathbf{x} \in S$. Then, by the convexity of the functions g_i , $i = 1, \dots, m$, it holds that

$$-\nabla g_i(\mathbf{x}^*)^T(\mathbf{x} - \mathbf{x}^*) \geq g_i(\mathbf{x}^*) - g_i(\mathbf{x}) = -g_i(\mathbf{x}) \geq 0, \quad (5.18)$$

for all $i \in \mathcal{I}(\mathbf{x}^*)$, and using the affinity of the functions h_j , $j = 1, \dots, \ell$, we get that

$$-\nabla h_j(\mathbf{x}^*)^T(\mathbf{x} - \mathbf{x}^*) = h_j(\mathbf{x}^*) - h_j(\mathbf{x}) = 0, \quad (5.19)$$

for all $j = 1, \dots, \ell$. Using the convexity of the objective function, equations (5.17a) and (5.17b), non-negativity of the Lagrange multipliers μ_i ,

Optimality conditions

$i \in \mathcal{I}(\mathbf{x}^*)$, and equations (5.18) and (5.19) we obtain the inequality

$$\begin{aligned} f(\mathbf{x}) - f(\mathbf{x}^*) &\geq \nabla f(\mathbf{x}^*)^\top (\mathbf{x} - \mathbf{x}^*) \\ &= - \sum_{i \in \mathcal{I}(\mathbf{x}^*)} \mu_i \nabla g_i(\mathbf{x}^*)^\top (\mathbf{x} - \mathbf{x}^*) - \sum_{j=1}^{\ell} \lambda_j \nabla h_j(\mathbf{x}^*)^\top (\mathbf{x} - \mathbf{x}^*) \\ &\geq 0. \end{aligned}$$

The point $\mathbf{x} \in S$ was arbitrary, whence \mathbf{x}^* solves (5.1). ■

Remark 5.46 (alternative proof of Theorem 5.45) An alternative proof of Theorem 5.45 is available from the sufficient global optimality condition in unconstrained optimization. Suppose that the conditions of Theorem 5.45 are fulfilled. By Theorem 4.19 this is equivalent to the Lagrangian function $\mathbf{x} \mapsto L(\mathbf{x}, \boldsymbol{\mu}^*, \boldsymbol{\lambda}^*) = f(\mathbf{x}) + (\boldsymbol{\mu}^*)^\top \mathbf{g}(\mathbf{x}) + (\boldsymbol{\lambda}^*)^\top \mathbf{h}(\mathbf{x})$ having a global minimum over \mathbb{R}^n at \mathbf{x}^* . In particular, then,

$$L(\mathbf{x}^*, \boldsymbol{\mu}^*, \boldsymbol{\lambda}^*) \leq L(\mathbf{x}, \boldsymbol{\mu}^*, \boldsymbol{\lambda}^*)$$

holds for every *feasible* \mathbf{x} . The rest of the proof is a simple matter of writing out this inequality explicitly and utilizing the remaining parts of the KKT conditions (5.17):

$$\begin{aligned} f(\mathbf{x}^*) + (\boldsymbol{\mu}^*)^\top \mathbf{g}(\mathbf{x}^*) + (\boldsymbol{\lambda}^*)^\top \mathbf{h}(\mathbf{x}^*) &\leq f(\mathbf{x}) + (\boldsymbol{\mu}^*)^\top \mathbf{g}(\mathbf{x}) + (\boldsymbol{\lambda}^*)^\top \mathbf{h}(\mathbf{x}) \\ &\iff \\ f(\mathbf{x}^*) &\leq f(\mathbf{x}) + \underbrace{(\boldsymbol{\mu}^*)^\top \mathbf{g}(\mathbf{x})}_{\leq 0 \text{ [(5.17c)+feas.]}} - \underbrace{(\boldsymbol{\mu}^*)^\top \mathbf{g}(\mathbf{x}^*)}_{=0 \text{ [(5.17b)]}} + \underbrace{(\boldsymbol{\lambda}^*)^\top [\mathbf{h}(\mathbf{x}) - \mathbf{h}(\mathbf{x}^*)]}_{=0 \text{ [feas.]}}. \end{aligned}$$

(Here, “feas.” stands for “primal feasibility.”) We are done. ■

Theorem 5.45 combined with the necessity of the KKT conditions under an appropriate CQ leads to the following statement.

Corollary 5.47 *Assume that the problem (5.1) is convex and verifies the Slater CQ (Definition 5.38). Then, for $\mathbf{x}^* \in S$ to be a globally optimal solution of (5.1) it is both necessary and sufficient to verify the KKT system (5.17).* ■

Not surprisingly, without the Slater constraint qualification the KKT conditions remain only sufficient (i.e., they are unnecessarily strong), as the following example demonstrates.

Example 5.48 Consider the optimization problem to

$$\begin{aligned} & \text{minimize } x_1, \\ & \text{subject to } x_1^2 + x_2 \leq 0, \\ & \quad -x_2 \leq 0, \end{aligned}$$

which is convex but has only one feasible point $\mathbf{0}^2 \in \mathbb{R}^2$. At this unique point both the inequality constraints are active, and thus the Slater CQ is violated, which however does not contradict the global optimality of $\mathbf{0}^2$. It is easy to check that the KKT system

$$\begin{aligned} \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \begin{pmatrix} 0 & 0 \\ 1 & -1 \end{pmatrix} \boldsymbol{\mu} &= \mathbf{0}^2, \\ \boldsymbol{\mu} &\geq \mathbf{0}^2, \end{aligned}$$

is unsolvable, and therefore the KKT conditions are not necessary without a CQ even for convex problems. ■

5.9 Applications and examples

Example 5.49 Consider a symmetric square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, and the optimization problem to

$$\begin{aligned} & \text{minimize } -\mathbf{x}^T \mathbf{A} \mathbf{x}, \\ & \text{subject to } \mathbf{x}^T \mathbf{x} \leq 1. \end{aligned}$$

The only constraint of this problem is convex; furthermore, $(\mathbf{0}^n)^T \mathbf{0}^n = 0 < 1$, and thus Slater's CQ (Definition 5.38) is verified. Therefore, the KKT conditions are necessary for the local optimality in this problem. We will find all the possible KKT points, and then choose a globally optimal point among them.

$\nabla(-\mathbf{x}^T \mathbf{A} \mathbf{x}) = -2\mathbf{A} \mathbf{x}$ (\mathbf{A} is symmetric), and $\nabla(\mathbf{x}^T \mathbf{x}) = 2\mathbf{x}$. Thus, the KKT system is as follows: $\mathbf{x}^T \mathbf{x} \leq 1$ and

$$\begin{aligned} -2\mathbf{A} \mathbf{x} + 2\mu \mathbf{x} &= \mathbf{0}^n, \\ \mu &\geq 0, \\ \mu(\mathbf{x}^T \mathbf{x} - 1) &= 0. \end{aligned}$$

From the first two equations we immediately see that either $\mathbf{x} = \mathbf{0}^n$, or the pair (μ, \mathbf{x}) is, respectively, a nonnegative eigenvalue and a corresponding eigenvector of \mathbf{A} (recall that $\mathbf{A} \mathbf{x} = \mu \mathbf{x}$ holds). In the former case, from the complementarity condition we deduce that $\mu = 0$.

Thus, we can characterize the KKT points of the problem into the following groups:

Optimality conditions

1. Let μ_1, \dots, μ_k be all the *positive* eigenvalues of \mathbf{A} (if any), and define $X_i := \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{x}^T \mathbf{x} = 1; \mathbf{A}\mathbf{x} = \mu_i \mathbf{x}\}$ to be the set of corresponding eigenvectors of length 1, $i = 1, \dots, k$. Then, (\mathbf{x}, μ_i) is a KKT point with the corresponding multiplier for every $\mathbf{x} \in X_i$, $i = 1, \dots, k$. Moreover, $-\mathbf{x}^T \mathbf{A}\mathbf{x} = -\mu_i \mathbf{x}^T \mathbf{x} = -\mu_i < 0$, for every $\mathbf{x} \in X_i$, $i = 1, \dots, k$.
2. Define also $X_0 := \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{x}^T \mathbf{x} \leq 1; \mathbf{A}\mathbf{x} = \mathbf{0}^n\}$. Then, the pair $(\mathbf{x}, 0)$ is a KKT point with the corresponding multiplier for every $\mathbf{x} \in X_0$. We note that if the matrix \mathbf{A} is nonsingular, then $X_0 = \{\mathbf{0}^n\}$. In any case, $-\mathbf{x}^T \mathbf{A}\mathbf{x} = 0$ for every $\mathbf{x} \in X_0$.

Therefore, if the matrix \mathbf{A} has any positive eigenvalue, then the global minima of the problem we consider are the eigenvectors of length one, corresponding to the largest positive eigenvalue; otherwise, every vector $\mathbf{x} \in X_0$ is globally optimal. ■

Example 5.50 Similarly to the previous example, consider the following equality-constrained minimization problem associated with a symmetric matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$:

$$\begin{aligned} &\text{minimize } -\mathbf{x}^T \mathbf{A}\mathbf{x}, \\ &\text{subject to } \mathbf{x}^T \mathbf{x} = 1. \end{aligned}$$

The gradient of the only equality constraint equals $2\mathbf{x}$, and since $\mathbf{0}^n$ is infeasible, LICQ is satisfied at every feasible point (see Definition 5.41), and the KKT conditions are necessary for local optimality. In this case, the KKT system is extremely simple: $\mathbf{x}^T \mathbf{x} = 1$ and

$$-2\mathbf{A}\mathbf{x} + 2\lambda\mathbf{x} = \mathbf{0}^n.$$

Let $\lambda_1 < \lambda_2 < \dots < \lambda_k$ denote all distinct eigenvalues of \mathbf{A} , and define as before $X_i := \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{x}^T \mathbf{x} = 1; \mathbf{A}\mathbf{x} = \lambda_i \mathbf{x}\}$ to be the set of corresponding eigenvectors of length 1, $i = 1, \dots, k$. Then, (\mathbf{x}, λ_i) is a KKT point with the corresponding multiplier for every $\mathbf{x} \in X_i$, $i = 1, \dots, k$. Furthermore, since $-\mathbf{x}^T \mathbf{A}\mathbf{x} = -\lambda_i$ for every $\mathbf{x} \in X_i$, $i = 1, \dots, k$, it holds that every $\mathbf{x} \in X_k$, that is, every eigenvector corresponding to the *largest* eigenvalue, is globally optimal.

Considering the problem for $\mathbf{A}^T \mathbf{A}$ and using the spectral theorem, we deduce the well known fact that $\|\mathbf{A}\| = \max_{1 \leq i \leq k} \{|\lambda_i|\}$. ■

Example 5.51 Consider the problem of finding the projection of a given point \mathbf{y} onto the polyhedron $\{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{x} = \mathbf{b}\}$, where $\mathbf{A} \in \mathbb{R}^{k \times n}$, $\mathbf{b} \in \mathbb{R}^k$. Thus, we consider the following minimization problem with

affine constraints (so that the KKT conditions are necessary for the local optimality, see Section 5.7.4):

$$\begin{aligned} & \text{minimize } \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|^2, \\ & \text{subject to } \mathbf{Ax} = \mathbf{b}. \end{aligned}$$

The KKT system in this case is written as follows:

$$\begin{aligned} \mathbf{Ax} &= \mathbf{b}, \\ (\mathbf{x} - \mathbf{y}) + \mathbf{A}^T \boldsymbol{\lambda} &= \mathbf{0}^n, \end{aligned}$$

for some $\boldsymbol{\lambda} \in \mathbb{R}^k$. Pre-multiplying the last equation with \mathbf{A} , and using the fact that $\mathbf{Ax} = \mathbf{b}$ we get:

$$\mathbf{AA}^T \boldsymbol{\lambda} = \mathbf{Ay} - \mathbf{b}.$$

Substituting an arbitrary solution of this equation into the KKT system, we calculate \mathbf{x} via $\mathbf{x} := \mathbf{y} - \mathbf{A}^T \boldsymbol{\lambda}$. It can be shown that the vector $\mathbf{A}^T \boldsymbol{\lambda}$ is the same constant for every Lagrange multiplier $\boldsymbol{\lambda}$, so using this formula we obtain the globally optimal solution to our minimization problem.

Now assume that the columns of \mathbf{A}^T are linearly independent, i.e., LICQ holds. Then, the matrix \mathbf{AA}^T is nonsingular, and the multiplier $\boldsymbol{\lambda}$ is therefore unique:

$$\boldsymbol{\lambda} = (\mathbf{AA}^T)^{-1}(\mathbf{Ay} - \mathbf{b}).$$

Substituting this into the KKT system, we finally obtain

$$\mathbf{x} = \mathbf{y} - \mathbf{A}^T (\mathbf{AA}^T)^{-1} (\mathbf{Ay} - \mathbf{b}),$$

the well known formula for calculating the projection. ■

5.10 Notes and further reading

One cannot overemphasize the importance of the Karush–Kuhn–Tucker optimality conditions for any development in optimization. We essentially follow the ideas presented in [BSS93, Chapters 4 and 5]; see also [Ber99, Chapter 3]. The original papers by Fritz John [Joh48], and Kuhn and Tucker [KuT51] might also be interesting. The work of Karush is a 1939 M.Sc. thesis from the University of Chicago.

Various forms of constraint qualifications play an especially important role in sensitivity analyses and studies of parametric optimization problems (e.g., [Fia83, BoS00]). Original presentations of constraint

qualifications, some of which we considered in this chapter, may be found in the works of Arrow, Hurwitz, and Uzawa [AHU61], Abadie [Aba67], Mangasarian and Fromowitz [MaF67], Guignard [Gui69], Zangwill [Zan69], and Evans [Eva70].

5.11 Exercises

Exercise 5.1 Consider the following problem:

$$\begin{aligned} &\text{minimize} && f(\mathbf{x}) := 2x_1^2 + 2x_1x_2 + x_2^2 - 10x_1 - 10x_2, \\ &\text{subject to} && x_1^2 + x_2^2 \leq 5, \\ &&& 3x_1 + x_2 \leq 6. \end{aligned}$$

Check if the point $\mathbf{x}^0 = (1, 2)^T$ is a KKT point for this problem. Is this an optimal solution? Which CQs are satisfied at the point \mathbf{x}^0 ?

Exercise 5.2 (optimality conditions, exam 020529) (a) Consider the following optimization problem:

$$\begin{aligned} &\text{minimize} && x^2, \\ &\text{subject to} && \sin(x) \leq -1. \end{aligned} \tag{5.20}$$

Find every locally and every globally optimal solution. Write down the KKT conditions. Are they necessary/sufficient for this problem?

(b) Do the locally/globally optimal solutions to the problem (5.20) satisfy the FJ optimality conditions?

(c) Question the usefulness of the FJ optimality conditions by finding a point (x, y) , which satisfies the FJ conditions for the problem:

$$\begin{aligned} &\text{minimize} && y, \\ &\text{subject to} && x^2 + y^2 \leq 1, \\ &&& x^3 \geq y^4, \end{aligned}$$

but, nevertheless, is neither a local nor a global minimum.

Exercise 5.3 Consider the following *linear* programming problem:

$$\begin{aligned} &\text{minimize} && \mathbf{c}^T \mathbf{x}, \\ &\text{subject to} && \mathbf{A}\mathbf{x} \geq \mathbf{b}. \end{aligned}$$

State the KKT conditions for this problem. Verify that every KKT point \mathbf{x} satisfies $\mathbf{c}^T \mathbf{x} = \mathbf{b}^T \boldsymbol{\lambda}$, where $\boldsymbol{\lambda}$ is a vector of KKT multipliers.

Exercise 5.4 (optimality conditions, exam 020826) (a) Consider the nonlinear programming problem with equality constraints:

$$\begin{aligned} &\text{minimize} && f(\mathbf{x}), \\ &\text{subject to} && h_i(\mathbf{x}) = 0, \quad i = 1, \dots, m, \end{aligned} \tag{5.21}$$

where f, h_1, \dots, h_m are continuously differentiable functions.

Show that the problem (5.21) is equivalent to the following problem with one inequality constraint:

$$\begin{aligned} & \text{minimize } f(\mathbf{x}), \\ & \text{subject to } \sum_{i=1}^m (h_i(\mathbf{x}))^2 \leq 0. \end{aligned} \quad (5.22)$$

Show (by a formal argument or an illustrative example) that the KKT conditions for the latter problem are not necessary for local optimality.

Can Slater's CQ or LICQ be satisfied for the problem (5.22)?

(b) Consider the unconstrained minimization problem to

$$\text{minimize } \max\{f_1(\mathbf{x}), f_2(\mathbf{x})\},$$

where $f_1 : \mathbb{R}^n \rightarrow \mathbb{R}, f_2 : \mathbb{R}^n \rightarrow \mathbb{R}$ are in C^1 . Show that if \mathbf{x}^* is a local minimum for this problem, then there exist $\mu_1, \mu_2 \in \mathbb{R}$ such that

$$\mu_1 \geq 0, \mu_2 \geq 0, \quad \mu_1 \nabla f_1(\mathbf{x}^*) + \mu_2 \nabla f_2(\mathbf{x}^*) = \mathbf{0}^n, \quad \mu_1 + \mu_2 = 1,$$

and $\mu_i = 0$ if $f_i(\mathbf{x}^*) < \max\{f_1(\mathbf{x}^*), f_2(\mathbf{x}^*)\}$, $i = 1, 2$.

Exercise 5.5 Consider the following optimization problem:

$$\begin{aligned} & \text{minimize } \frac{1}{2} \mathbf{x}^T \mathbf{x}, \\ & \text{subject to } \mathbf{A} \mathbf{x} = \mathbf{b}. \end{aligned}$$

Assume that the matrix \mathbf{A} has full row rank. Find the globally optimal solution to this problem.

Exercise 5.6 Consider the following optimization problem:

$$\begin{aligned} & \text{minimize } \sum_{j=1}^n c_j x_j, \\ & \text{subject to } \sum_{j=1}^n x_j^2 \leq 1, \\ & \quad -x_j \leq 0, \quad j = 1, \dots, n. \end{aligned} \quad (5.23)$$

Assume that $\min\{c_1, \dots, c_n\} < 0$, and let us introduce KKT multipliers $\lambda \geq 0$ and $\mu_j \geq 0, j = 1, \dots, n$ for the inequality constraints.

(a) Show that the equalities

$$\begin{aligned} x_j^* &= \min\{0, c_j\} / (2\lambda^*), \quad j = 1, \dots, n, \\ \lambda^* &= \frac{1}{2} \left(\sum_{j=1}^n [\min\{0, c_j\}]^2 \right)^{1/2}, \\ \mu_j^* &= \max\{0, c_j\}, \quad j = 1, \dots, n, \end{aligned}$$

define a KKT point for (5.23).

(b) Show that there is only one optimal solution to (5.23).

Optimality conditions

Exercise 5.7 (optimality conditions, exam 040308) Consider the following optimization problem:

$$\begin{aligned} & \underset{(x,y) \in \mathbb{R} \times \mathbb{R}}{\text{minimize}} && f(x,y) := \frac{1}{2}(x-2)^2 + \frac{1}{2}(y-1)^2, \\ & \text{subject to} && x - y \geq 0, \\ & && y \geq 0, \\ & && y(x-y) = 0. \end{aligned} \tag{5.24}$$

(a) Find *all* points of global and local minima (you may do this graphically), as well as *all* KKT points. Is this a convex problem? Are the KKT optimality conditions necessary and/or sufficient for local optimality *in this problem*?

(b) Demonstrate that LICQ is violated at *every feasible point* of the problem (5.24). Show that instead of solving the problem (5.24) we can solve *two* convex optimization problems that furthermore verify some constraint qualification, and then choose the best point out of the two.

(c) Generalize the procedure from the previous part to the more general optimization problem to

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && g(\mathbf{x}), \\ & \text{subject to} && \mathbf{a}_i^T \mathbf{x} \geq b_i, \quad i = 1, \dots, n, \\ & && x_i \geq 0, \quad i = 1, \dots, n, \\ & && x_i(\mathbf{a}_i^T \mathbf{x} - b_i) = 0, \quad i = 1, \dots, n, \end{aligned}$$

where $\mathbf{x} = (x_1, \dots, x_n)^T \in \mathbb{R}^n$, $\mathbf{a}_i \in \mathbb{R}^n$, $b_i \in \mathbb{R}$, $i = 1, \dots, n$, and $g : \mathbb{R}^n \rightarrow \mathbb{R}$ is a convex differentiable function.

Exercise 5.8 Determine the values of the parameter c for which the point $(x, y) = (4, 3)$ is an optimal solution to the following problem:

$$\begin{aligned} & \underset{(x,y) \in \mathbb{R} \times \mathbb{R}}{\text{minimize}} && cx + y, \\ & \text{subject to} && x^2 + y^2 \leq 25, \\ & && x - y \leq 1. \end{aligned}$$

Exercise 5.9 Consider the following optimization problem:

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && f(\mathbf{x}) := \sum_{j=1}^n \frac{x_j^2}{c_j}, \\ & \text{subject to} && \sum_{j=1}^n x_j = D, \\ & && x_j \geq 0, \quad j = 1, \dots, n, \end{aligned}$$

where $c_j > 0$, $j = 1, \dots, n$, and $D > 0$. Find the unique globally optimal solution to this problem.

Lagrangian duality

VI

This chapter collects some basic results on Lagrangian duality, in particular as it applies to convex programs with a zero duality gap.

6.1 The relaxation theorem

Given the problem to find

$$f^* := \infimum_x f(\mathbf{x}), \quad (6.1a)$$

$$\text{subject to } \mathbf{x} \in S, \quad (6.1b)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a given function and $S \subseteq \mathbb{R}^n$, we define a *relaxation* to (6.1) to be a problem of the following form: find

$$f_R^* := \infimum_x f_R(\mathbf{x}), \quad (6.2a)$$

$$\text{subject to } \mathbf{x} \in S_R, \quad (6.2b)$$

where $f_R : \mathbb{R}^n \rightarrow \mathbb{R}$ is a function with the property that $f_R \leq f$ on S , and where $S_R \supseteq S$. For this pair of problems, we have the following basic result.

Theorem 6.1 (Relaxation Theorem) (a) [relaxation] $f_R^* \leq f^*$.

(b) [infeasibility] If (6.2) is infeasible, then so is (6.1).

(c) [optimal relaxation] If the problem (6.2) has an optimal solution, \mathbf{x}_R^* , for which it holds that

$$\mathbf{x}_R^* \in S \quad \text{and} \quad f_R(\mathbf{x}_R^*) = f(\mathbf{x}_R^*), \quad (6.3)$$

then \mathbf{x}_R^* is an optimal solution to (6.1) as well.

Lagrangian duality

Proof. The result in (a) is obvious, as every solution feasible in (6.1) is both feasible in (6.2) and has a lower objective value in the latter problem. The result in (b) follows for similar reasons. For the result in (c), we note that

$$f(\mathbf{x}_R^*) = f_R(\mathbf{x}_R^*) \leq f_R(\mathbf{x}) \leq f(\mathbf{x}), \quad \mathbf{x} \in S,$$

from which the result follows. ■

This basic result will be utilized both in this chapter and later on to motivate why Lagrangian relaxation, objective function linearization and penalization constitute relaxations, and to derive optimality conditions and algorithms based on them.

6.2 Lagrangian duality

In this section we formulate the Lagrangian dual problem and establish its convexity. The Weak Duality Theorem is also established, and we introduce the terms “Lagrangian relaxation,” “Lagrange multiplier,” and “duality gap.”

6.2.1 Lagrangian relaxation and the dual problem

Consider the optimization problem to find

$$\begin{aligned} f^* &:= \inf_{\mathbf{x}} f(\mathbf{x}), \\ \text{subject to} \quad &\mathbf{x} \in X, \\ &g_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m, \end{aligned} \tag{6.4}$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$ ($i = 1, 2, \dots, m$) are given functions, and $X \subseteq \mathbb{R}^n$.

For this problem, we assume that

$$-\infty < f^* < \infty, \tag{6.5}$$

that is, that f is bounded from below on the feasible set and the problem has at least one feasible solution.

Definition 6.2 (Lagrange function, relaxation, multiplier) (a) For an arbitrary vector $\boldsymbol{\mu} \in \mathbb{R}^m$, the Lagrange function is

$$L(\mathbf{x}, \boldsymbol{\mu}) := f(\mathbf{x}) + \sum_{i=1}^m \mu_i g_i(\mathbf{x}) = f(\mathbf{x}) + \boldsymbol{\mu}^T \mathbf{g}(\mathbf{x}). \tag{6.6}$$

(b) Consider the problem to

$$\begin{aligned} & \text{minimize } L(\mathbf{x}, \boldsymbol{\mu}), \\ & \text{subject to } \mathbf{x} \in X. \end{aligned} \quad (6.7)$$

Whenever $\boldsymbol{\mu}$ is non-negative, the problem (6.7) is referred to as a Lagrangian relaxation.

(c) We call the vector $\boldsymbol{\mu}^* \in \mathbb{R}^m$ a Lagrange multiplier vector if it is non-negative and if $f^* = \inf_{\mathbf{x} \in X} L(\mathbf{x}, \boldsymbol{\mu}^*)$ holds. ■

Note that the Lagrangian relaxation (6.7) is a relaxation, in terms of Section 6.1.

Theorem 6.3 (Lagrange multipliers and global optima) *Let $\boldsymbol{\mu}^*$ be a Lagrange multiplier vector. Then, \mathbf{x}^* is an optimal solution to (6.4) if and only if \mathbf{x}^* is feasible in (6.4) and*

$$\mathbf{x}^* \in \arg \min_{\mathbf{x} \in X} L(\mathbf{x}, \boldsymbol{\mu}^*), \quad \text{and} \quad \mu_i^* g_i(\mathbf{x}^*) = 0, \quad i = 1, \dots, m. \quad (6.8)$$

Proof. If \mathbf{x}^* is an optimal solution to (6.4), then it is in particular feasible, and

$$f^* = f(\mathbf{x}^*) \geq L(\mathbf{x}^*, \boldsymbol{\mu}^*) \geq \inf_{\mathbf{x} \in X} L(\mathbf{x}, \boldsymbol{\mu}^*),$$

where the first inequality stems from the feasibility of \mathbf{x}^* and the definition of a Lagrange multiplier vector. The second part of that definition implies that $f^* = \inf_{\mathbf{x} \in X} L(\mathbf{x}, \boldsymbol{\mu}^*)$, so that equality holds throughout in the above line of inequalities. Hence, (6.8) follows.

Conversely, if \mathbf{x}^* is feasible and (6.8) holds, then by the use of the definition of a Lagrange multiplier vector,

$$f(\mathbf{x}^*) = L(\mathbf{x}^*, \boldsymbol{\mu}^*) = \min_{\mathbf{x} \in X} L(\mathbf{x}, \boldsymbol{\mu}^*) = f^*,$$

so \mathbf{x}^* is a global optimum. ■

Let

$$q(\boldsymbol{\mu}) := \inf_{\mathbf{x} \in X} L(\mathbf{x}, \boldsymbol{\mu}) \quad (6.9)$$

be the *Lagrangian dual function*, defined by the infimum value of the Lagrange function over X ; the *Lagrangian dual problem* is to

$$\begin{aligned} & \text{maximize } q(\boldsymbol{\mu}), \\ & \text{subject to } \boldsymbol{\mu} \geq \mathbf{0}^m. \end{aligned} \quad (6.10)$$

Lagrangian duality

For some $\boldsymbol{\mu}$, $q(\boldsymbol{\mu}) = -\infty$ is possible; if it is true for all $\boldsymbol{\mu} \geq \mathbf{0}^m$, then

$$q^* := \sup_{\boldsymbol{\mu} \geq \mathbf{0}^m} q(\boldsymbol{\mu})$$

equals $-\infty$. (We can then say that the dual problem is infeasible.)

The *effective domain* of q is

$$D_q := \{ \boldsymbol{\mu} \in \mathbb{R}^m \mid q(\boldsymbol{\mu}) > -\infty \}.$$

Theorem 6.4 (convex dual problem) *The effective domain D_q of q is convex, and q is concave on D_q .*

Proof. Let $\mathbf{x} \in \mathbb{R}^n$, $\boldsymbol{\mu}, \bar{\boldsymbol{\mu}} \in \mathbb{R}^m$, and $\alpha \in [0, 1]$. We have that

$$L(\mathbf{x}, \alpha\boldsymbol{\mu} + (1 - \alpha)\bar{\boldsymbol{\mu}}) = \alpha L(\mathbf{x}, \boldsymbol{\mu}) + (1 - \alpha)L(\mathbf{x}, \bar{\boldsymbol{\mu}}).$$

Take the infimum over $\mathbf{x} \in X$ on both sides; then,

$$\begin{aligned} \inf_{\mathbf{x} \in X} L(\mathbf{x}, \alpha\boldsymbol{\mu} + (1 - \alpha)\bar{\boldsymbol{\mu}}) &= \inf_{\mathbf{x} \in X} \{ \alpha L(\mathbf{x}, \boldsymbol{\mu}) + (1 - \alpha)L(\mathbf{x}, \bar{\boldsymbol{\mu}}) \} \\ &\geq \inf_{\mathbf{x} \in X} \alpha L(\mathbf{x}, \boldsymbol{\mu}) + \inf_{\mathbf{x} \in X} (1 - \alpha)L(\mathbf{x}, \bar{\boldsymbol{\mu}}) \\ &= \alpha \inf_{\mathbf{x} \in X} L(\mathbf{x}, \boldsymbol{\mu}) + (1 - \alpha) \inf_{\mathbf{x} \in X} L(\mathbf{x}, \bar{\boldsymbol{\mu}}), \end{aligned}$$

since $\alpha \in [0, 1]$, and the sum of infimum values may be smaller than the infimum of the sum, since in the former case we have the possibility to choose different optimal solutions in the two problems. Hence,

$$q(\alpha\boldsymbol{\mu} + (1 - \alpha)\bar{\boldsymbol{\mu}}) \geq \alpha q(\boldsymbol{\mu}) + (1 - \alpha)q(\bar{\boldsymbol{\mu}})$$

holds. This inequality has two implications: if $\boldsymbol{\mu}$ and $\bar{\boldsymbol{\mu}}$ lie in D_q , then so does $\alpha\boldsymbol{\mu} + (1 - \alpha)\bar{\boldsymbol{\mu}}$, so D_q is convex; also, q is concave on D_q . ■

That the Lagrangian dual problem always is convex (we indeed maximize a concave function) is good news, because it means that it can be solved efficiently. What remains is to show how a Lagrangian dual optimal solution can be used to generate a primal optimal solution.

Next, we establish that every feasible point in the Lagrangian *dual* problem always underestimates the objective function value of every feasible point in the *primal* problem; hence, also their optimal values have this relationship.

Theorem 6.5 (Weak Duality Theorem) (a) *Let \mathbf{x} and $\boldsymbol{\mu}$ be feasible in the problems (6.4) and (6.10), respectively. Then,*

$$q(\boldsymbol{\mu}) \leq f(\mathbf{x}).$$

In particular,

$$q^* \leq f^*.$$

(b) If $q(\boldsymbol{\mu}) = f(\mathbf{x})$, then the pair $(\mathbf{x}, \boldsymbol{\mu})$ is optimal in its respective problem.

Proof. For all $\boldsymbol{\mu} \geq \mathbf{0}^m$ and $\mathbf{x} \in X$ with $\mathbf{g}(\mathbf{x}) \leq \mathbf{0}^m$,

$$q(\boldsymbol{\mu}) = \inf_{\mathbf{z} \in X} L(\mathbf{z}, \boldsymbol{\mu}) \leq f(\mathbf{x}) + \boldsymbol{\mu}^T \mathbf{g}(\mathbf{x}) \leq f(\mathbf{x}),$$

so

$$q^* = \sup_{\boldsymbol{\mu} \geq \mathbf{0}^m} q(\boldsymbol{\mu}) \leq \inf_{\mathbf{x} \in X: \mathbf{g}(\mathbf{x}) \leq \mathbf{0}^m} f(\mathbf{x}) = f^*.$$

The result follows. ■

Weak duality is also a consequence of the Relaxation Theorem: For any $\boldsymbol{\mu} \geq \mathbf{0}^m$, let

$$S := X \cap \{ \mathbf{x} \in \mathbb{R}^n \mid \mathbf{g}(\mathbf{x}) \leq \mathbf{0}^m \}, \quad (6.11a)$$

$$S_R := X, \quad (6.11b)$$

$$f_R := L(\boldsymbol{\mu}, \cdot). \quad (6.11c)$$

Then, the weak duality statement is the result in Theorem 6.1(a).

If our initial feasibility assumption (6.5) is false, then what does weak duality imply? Suppose that $f^* = -\infty$. Then, weak duality implies that $q(\boldsymbol{\mu}) = -\infty$ for all $\boldsymbol{\mu} \geq \mathbf{0}^m$, that is, the dual problem is infeasible. Suppose then that $X \neq \emptyset$ but that $X \cap \{ \mathbf{x} \in \mathbb{R}^n \mid \mathbf{g}(\mathbf{x}) \leq \mathbf{0}^m \}$ is empty. Then, $f^* = \infty$, by convention. The dual function satisfies $q(\boldsymbol{\mu}) < \infty$ for all $\boldsymbol{\mu} \geq \mathbf{0}^m$, but it is possible that $q^* = -\infty$, $-\infty < q^* < \infty$, or $q^* = \infty$ (see [Ber99, Figure 5.1.8]). For linear programs, $-\infty < q^* < \infty$ implies $-\infty < f^* < \infty$; see below.

If $q^* = f^*$, then we say that the *duality gap* (as given by $\Gamma := f^* - q^*$) is zero, or that *there is no duality gap*. If there exists a Lagrange multiplier vector, then by the weak duality theorem, this implies that there is no duality gap. The converse is not true in general: there may be cases where no Lagrange multipliers exist even when there is no duality gap;¹ in that case though, the Lagrangian dual problem cannot have an optimal solution, as implied by the following result.

Proposition 6.6 (duality gap and the existence of Lagrange multipliers)

(a) If there is no duality gap, then the set of Lagrange multiplier vectors equals the set of optimal dual solutions (which however may be empty).

(b) If there is a duality gap, then there are no Lagrange multipliers.

¹Take the example of minimizing $f(x) := x$ subject to $g(x) := x^2 \leq 0$; $x \in X := \mathbb{R}$.

Lagrangian duality

Proof. By definition, a vector $\boldsymbol{\mu}^* \geq \mathbf{0}^m$ is a Lagrange multiplier vector if and only if $f^* = q(\boldsymbol{\mu}^*) \leq q^*$, the equality following from the definition of $q(\boldsymbol{\mu}^*)$ and the inequality from the definition of q^* as the supremum of $q(\boldsymbol{\mu})$ over \mathbb{R}_+^m . By weak duality, this relation holds if and only if there is no duality gap and $\boldsymbol{\mu}^*$ is an optimal dual solution. ■

Above we have developed properties of the min-max problem for finding

$$q^* := \sup_{\boldsymbol{\mu} \geq \mathbf{0}^m} \inf_{\mathbf{x} \in X} L(\mathbf{x}, \boldsymbol{\mu}).$$

What then is the max-min problem to find

$$p^* := \inf_{\mathbf{x} \in X} \sup_{\boldsymbol{\mu} \geq \mathbf{0}^m} L(\mathbf{x}, \boldsymbol{\mu})?$$

Fix $\mathbf{x} \in X$. Then,

$$p(\mathbf{x}) := \sup_{\boldsymbol{\mu} \geq \mathbf{0}^m} L(\mathbf{x}, \boldsymbol{\mu}) = \begin{cases} f(\mathbf{x}), & \text{if } g(\mathbf{x}) \leq \mathbf{0}^m, \\ +\infty, & \text{otherwise.} \end{cases}$$

(We call the function $p : \mathbb{R}_+^m \rightarrow \mathbb{R} \cup \{+\infty\}$ the *primal function*, in contrast to the *dual function* q .) Hence, the max-min problem is essentially equivalent to minimizing f over the set $X \cap \{\mathbf{x} \in \mathbb{R}^n \mid g(\mathbf{x}) \leq \mathbf{0}^m\}$, that is, the original problem (6.4), and $p^* = f^*$ holds. Moreover, we have seen above that in general $q^* \leq f^*$ holds, that is, the min-max problem has an optimal value which is always at least as large as that of the max-min problem. This is a general statement, and equality holds precisely when there exists a saddle point of the function L . The above development extends that of Von Neumann's matrix game; cf. (4.31).

Before moving on, we remark on the *statement* of the problem (6.4). There are several ways in which the original set of constraints of the problem can be placed either within the definition of the *ground set* X (which is kept intact), or within the explicit constraints defined by the functions g_i (which are Lagrangian relaxed). How to distinguish between the two, that is, how to decide whether a constraint should be kept or be Lagrangian relaxed, depends on several factors. For example, keeping more constraints within X may result in a smaller duality gap, and fewer multipliers also result in a simpler Lagrangian dual problem. On the other hand, the Lagrangian subproblem defining the dual function simultaneously becomes more complex and difficult to solve. There are no immediate rules to follow, but experimentation and experience.

6.2.2 Global optimality conditions

The following result characterizes every optimal primal and dual solution. It is however applicable only in the presence of Lagrange multipliers; in other words, the below system (6.12) is consistent if and only if there exists a Lagrange multiplier vector and there is no duality gap.

Theorem 6.7 (global optimality conditions in the absence of a duality gap) *The vector $(\mathbf{x}^*, \boldsymbol{\mu}^*)$ is a pair of primal optimal solution and Lagrange multiplier vector if and only if*

$$\boldsymbol{\mu}^* \geq \mathbf{0}^m, \quad (\text{Dual feasibility}) \quad (6.12a)$$

$$\mathbf{x}^* \in \arg \min_{\mathbf{x} \in X} L(\mathbf{x}, \boldsymbol{\mu}^*), \quad (\text{Lagrangian optimality}) \quad (6.12b)$$

$$\mathbf{x}^* \in X, \quad \mathbf{g}(\mathbf{x}^*) \leq \mathbf{0}^m, \quad (\text{Primal feasibility}) \quad (6.12c)$$

$$\mu_i^* g_i(\mathbf{x}^*) = 0, \quad i = 1, \dots, m. \quad (\text{Complementary slackness}) \quad (6.12d)$$

Proof. Suppose that the pair $(\mathbf{x}^*, \boldsymbol{\mu}^*)$ satisfies (6.12). Then, from (6.12a) we have that the Lagrangian problem to minimize $L(\mathbf{x}, \boldsymbol{\mu}^*)$ over $\mathbf{x} \in X$ is a (Lagrangian) relaxation of (6.4). Moreover, according to (6.12b) \mathbf{x}^* solves this problem, (6.12c) shows that \mathbf{x}^* is feasible in (6.4), and (6.12d) implies that $L(\mathbf{x}^*, \boldsymbol{\mu}^*) = f(\mathbf{x}^*)$. The Relaxation Theorem 6.1 then yields that \mathbf{x}^* is optimal in (6.4), which in turn implies that $\boldsymbol{\mu}^*$ is a Lagrange multiplier vector.

Conversely, if $(\mathbf{x}^*, \boldsymbol{\mu}^*)$ is a pair of optimal primal solution and Lagrange multiplier vector, then they are primal and dual feasible, respectively. The relations (6.12b) and (6.12d) follow from Theorem 6.3. ■

Theorem 6.8 (global optimality and saddle points) *The vector $(\mathbf{x}^*, \boldsymbol{\mu}^*)$ is a pair of optimal primal solution and Lagrange multiplier vector if and only if $\mathbf{x}^* \in X$, $\boldsymbol{\mu}^* \geq \mathbf{0}^m$, and $(\mathbf{x}^*, \boldsymbol{\mu}^*)$ is a saddle point of the Lagrangian function on $X \times \mathbb{R}_+^m$, that is,*

$$L(\mathbf{x}^*, \boldsymbol{\mu}) \leq L(\mathbf{x}^*, \boldsymbol{\mu}^*) \leq L(\mathbf{x}, \boldsymbol{\mu}^*), \quad (\mathbf{x}, \boldsymbol{\mu}) \in X \times \mathbb{R}_+^m \quad (6.13)$$

holds.

Proof. We establish that (6.12) and (6.13) are equivalent; Theorem 6.7 then gives the result. The first inequality in (6.13) is equivalent to

$$-\mathbf{g}(\mathbf{x}^*)^\top (\boldsymbol{\mu} - \boldsymbol{\mu}^*) \geq 0, \quad \boldsymbol{\mu} \in \mathbb{R}_+^m, \quad (6.14)$$

Lagrangian duality

for the given pair $(\mathbf{x}^*, \boldsymbol{\mu}^*) \in X \times \mathbb{R}_+^m$. This variational inequality is equivalent to stating that²

$$\mathbf{0}^m \geq \mathbf{g}(\mathbf{x}^*) \perp \boldsymbol{\mu}^* \geq \mathbf{0}^m, \quad (6.15)$$

where \perp denotes orthogonality: that is, for any vectors $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$, $\mathbf{a} \perp \mathbf{b}$ means that $\mathbf{a}^T \mathbf{b} = 0$. Because of the sign restrictions posed on $\boldsymbol{\mu}$ and \mathbf{g} , that is, the vectors \mathbf{a} and \mathbf{b} , the relation $\mathbf{a} \perp \mathbf{b}$ actually means that not only does it hold that $\mathbf{a}^T \mathbf{b} = 0$ but in fact $a_i b_i = 0$ must hold for all $i = 1, \dots, n$. This complementarity system is, for the given $\boldsymbol{\mu}^* \in \mathbb{R}_+^m$, the same as (6.12a), (6.12c) and (6.12d). The second inequality in (6.13) is equivalent to (6.12b). ■

The above two theorems also imply that the set of primal–dual optimal solutions $(\mathbf{x}^*, \boldsymbol{\mu}^*)$ is a Cartesian product set, $X^* \times U^*$. For example, given any optimal dual solution $\boldsymbol{\mu}^* \in U^*$, every optimal primal solution $\mathbf{x}^* \in X^*$ satisfies (6.12). Hence, we can write, for an *arbitrary* dual vector $\boldsymbol{\mu}^* \in U^*$,

$$\begin{aligned} X^* &= \{ \mathbf{x}^* \in \mathbb{R}^n \mid \mathbf{x}^* \text{ satisfies (6.12) for } \boldsymbol{\mu} = \boldsymbol{\mu}^* \} \\ &= \left\{ \mathbf{x}^* \in \arg \min_{\mathbf{x} \in X} L(\mathbf{x}, \boldsymbol{\mu}^*) \mid \mathbf{g}(\mathbf{x}^*) \leq \mathbf{0}^m; (\boldsymbol{\mu}^*)^T \mathbf{g}(\mathbf{x}^*) = 0 \right\}. \end{aligned}$$

We note that structurally similar results to the above two theorems which are valid for the general problem (6.4) with any size of the duality gap can be found in [LaP05].³

We finally note a *practical* connection between the KKT system (5.9) and the above system (6.12). The practical use of the KKT system is normally to investigate whether a primal vector \mathbf{x} —obtained perhaps from a solver for our problem—is a candidate for a locally optimal solution; in other words, we have access to \mathbf{x} and generate a vector $\boldsymbol{\mu}$ of

²We establish the equivalence between (6.14) and (6.15) as follows. (The proof extends that for line search problems in unconstrained optimization in a footnote in Section 11.3.1.)

First, suppose that (6.15) is fulfilled. Then, $-\mathbf{g}(\mathbf{x}^*)^T(\boldsymbol{\mu} - \boldsymbol{\mu}^*) = -\mathbf{g}(\mathbf{x}^*)^T \boldsymbol{\mu} \geq 0$, for all $\boldsymbol{\mu} \geq \mathbf{0}^m$, that is, (6.14) is fulfilled. Conversely, suppose that (6.14) is fulfilled. Setting $\boldsymbol{\mu} = \mathbf{0}^m$ yields that $\mathbf{g}(\mathbf{x}^*)^T \boldsymbol{\mu}^* \geq 0$. On the other hand, the choice $\boldsymbol{\mu} = 2\boldsymbol{\mu}^*$ yields that $-\mathbf{g}(\mathbf{x}^*)^T \boldsymbol{\mu}^* \geq 0$. Hence, $\mathbf{g}(\mathbf{x}^*)^T \boldsymbol{\mu}^* = 0$ holds. Last, let $\boldsymbol{\mu} = \boldsymbol{\mu}^* + \mathbf{e}_i$, where \mathbf{e}_i is the i^{th} unit vector in \mathbb{R}^m . Then, $-\mathbf{g}(\mathbf{x}^*)^T(\boldsymbol{\mu} - \boldsymbol{\mu}^*) = -g_i(\mathbf{x}^*) \geq 0$. Since this is true for all $i \in \{1, 2, \dots, m\}$ we have obtained that $-\mathbf{g}(\mathbf{x}^*) \geq \mathbf{0}^m$, that is, $\mathbf{g}(\mathbf{x}^*) \leq \mathbf{0}^m$. We are done.

³The system (6.12) is there appended with two relaxation parameters which measure, respectively, the near-optimality of \mathbf{x}^* in the Lagrangian subproblem [that is, the ε -optimality in (6.12b)], and the violation of the complementarity conditions (6.12d). The saddle point condition (6.13) is similarly perturbed, and at an optimal solution, the sum of these two parameter values equals the duality gap.

Lagrange multipliers in the investigation of the KKT system (5.9). In contrast, the system (6.12) is normally investigated in the reverse order; we formulate and solve the Lagrangian dual problem, thereby obtaining an optimal dual vector $\boldsymbol{\mu}$. Starting from that vector, we investigate the global optimality conditions stated in (6.12) to obtain, if possible, an optimal primal vector \boldsymbol{x} . In the section to follow, we show when this is possible, and provide strong connections between the systems (5.9) and (6.12) in the convex and differentiable case.

6.2.3 Strong duality for convex programs

So far the results have been rather non-technical to achieve: the convexity of the Lagrangian dual problem comes with very few assumptions on the original, primal problem, and the characterization of the primal-dual set of optimal solutions is simple and also quite easily established. In order to establish *strong duality*, that is, to establish sufficient conditions under which there is no duality gap, however, takes much more. In particular, as is the case with the KKT conditions we need regularity conditions (that is, constraint qualifications), and we also need to utilize separation theorems such as Theorem 4.29. Most importantly, however, is that strong duality is deeply associated with the convexity of the original problem, and it is in particular under convexity that the primal and dual optimal solutions are linked through the global optimality conditions provided in the previous section. We begin by concentrating on the inequality constrained case, proving this result in detail. We will also specialize the result to quadratic and linear optimization problems.

Consider the inequality constrained *convex* program (6.4), where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and g_i ($i = 1, \dots, m$) are convex functions and $X \subseteq \mathbb{R}^n$ is a convex set. For this problem, we introduce the following regularity condition, due to Slater (cf. Definition 5.38):

$$\exists \boldsymbol{x} \in X \text{ with } \boldsymbol{g}(\boldsymbol{x}) < \mathbf{0}^m. \quad (6.16)$$

Theorem 6.9 (Strong Duality, inequality constrained convex programs) *Suppose that the feasibility condition (6.5) and Slater's constraint qualification (6.16) hold for the convex problem (6.4).*

(a) *There is no duality gap and there exists at least one Lagrange multiplier vector $\boldsymbol{\mu}^*$. Moreover, the set of Lagrange multipliers is bounded and convex.*

(b) *If the infimum in (6.4) is attained at some \boldsymbol{x}^* , then the pair $(\boldsymbol{x}^*, \boldsymbol{\mu}^*)$ satisfies the global optimality conditions (6.12).*

(c) *If further f and \boldsymbol{g} are differentiable at \boldsymbol{x}^* , then the condition*

Lagrangian duality

(6.12b) can equivalently be written as the variational inequality

$$\nabla_x L(\mathbf{x}^*, \boldsymbol{\mu}^*)^\top (\mathbf{x} - \mathbf{x}^*) \geq 0, \quad \mathbf{x} \in X. \quad (6.17)$$

If, in addition, X is open (such as is the case when $X = \mathbb{R}^n$), then this reduces to the condition that

$$\nabla_x L(\mathbf{x}^*, \boldsymbol{\mu}^*) = \nabla f(\mathbf{x}^*) + \sum_{i=1}^m \mu_i^* \nabla g_i(\mathbf{x}^*) = \mathbf{0}^n, \quad (6.18)$$

and the global optimality conditions (6.12) reduce to the Karush–Kuhn–Tucker conditions stated in Theorem 5.25.

Proof. (a) We begin by establishing the existence of a Lagrange multiplier vector (and the presence of a zero duality gap).

First, we consider the following subset of \mathbb{R}^{m+1} :

$$A := \{(z_1, \dots, z_m, w)^\top \mid \exists \mathbf{x} \in X \text{ with } g_i(\mathbf{x}) \leq z_i, i = 1, \dots, m; f(\mathbf{x}) \leq w\}.$$

It is elementary to show that A is convex.

Next, we observe that $((\mathbf{0}^m)^\top, f^*)^\top$ is not an interior point of A ; otherwise, for some $\varepsilon > 0$ the point $((\mathbf{0}^m)^\top, f^* - \varepsilon)^\top \in A$ holds, which would contradict the definition of f^* . Therefore, by the (possibly non-proper) separation result in Theorem 4.29, we can find a hyperplane passing through $((\mathbf{0}^m)^\top, f^*)^\top$ such that A lies in one of the two corresponding half-spaces. In particular, there then exists a vector $(\boldsymbol{\mu}^\top, \beta)^\top \neq ((\mathbf{0}^m)^\top, 0)^\top$ such that

$$\beta f^* \leq \beta w + \boldsymbol{\mu}^\top \mathbf{z}, \quad (\mathbf{z}^\top, w)^\top \in A. \quad (6.19)$$

This implies that

$$\beta \geq 0; \quad \boldsymbol{\mu} \geq \mathbf{0}^m, \quad (6.20)$$

since for each $(\mathbf{z}^\top, w)^\top \in A$, $(\mathbf{z}^\top, w + \gamma)^\top \in A$ and $(z_1, \dots, z_{i-1}, z_i + \gamma, z_{i+1}, \dots, z_m, w)^\top \in A$ for all $\gamma > 0$ and $i = 1, \dots, m$.

We claim that $\beta > 0$ in fact holds. Indeed, if it was not the case, then $\beta = 0$ and (6.19) then implies that $\boldsymbol{\mu}^\top \mathbf{z} \geq 0$ for every pair $(\mathbf{z}^\top, w)^\top \in A$. But since $(\mathbf{g}(\bar{\mathbf{x}})^\top, f(\bar{\mathbf{x}}))^\top \in A$ [where $\bar{\mathbf{x}}$ is such that it satisfies the Slater condition (6.16)], we would obtain that $0 \leq \sum_{i=1}^m \mu_i g_i(\bar{\mathbf{x}})$ which in view of $\boldsymbol{\mu} \geq \mathbf{0}^m$ [cf. (6.20)] and the assumption that $\bar{\mathbf{x}}$ satisfies the Slater condition (6.16) implies that $\boldsymbol{\mu} = \mathbf{0}^m$. This means, however, that $(\boldsymbol{\mu}^\top, \beta)^\top = ((\mathbf{0}^m)^\top, 0)^\top$ —a contradiction. We may therefore claim that $\beta > 0$. We further, with no loss of generality, assume that $\beta = 1$.

Thus, since $(\mathbf{g}(\mathbf{x})^\top, f(\mathbf{x}))^\top \in A$ for every $\mathbf{x} \in X$, (6.19) yields that

$$f^* \leq f(\mathbf{x}) + \boldsymbol{\mu}^\top \mathbf{g}(\mathbf{x}), \quad \mathbf{x} \in X.$$

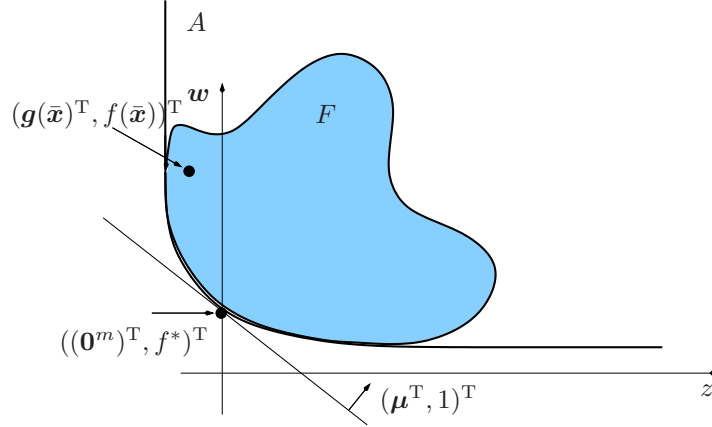


Figure 6.1: Illustration of the set $F := \{ (g(x)^T, f(x))^T \mid x \in X \}$ and the set A used in the proof of Theorem 6.9. The idea of the proof is to show that A is convex and that $((0^m)^T, f^*)^T$ is not an interior point of A . A hyperplane passing through $((0^m)^T, f^*)^T$ and supporting A is used to construct a Lagrange multiplier.

Taking the infimum over $x \in X$ and using that $\mu \geq 0^m$ we obtain

$$f^* \leq \inf_{x \in X} \{ f(x) + \mu^T g(x) \} = q(\mu) \leq \sup_{\mu \geq 0^m} q(\mu) = q^*.$$

From the Weak Duality Theorem 6.5 follows that μ is a Lagrange multiplier vector, and that there is no duality gap.

Take any vector $\bar{x} \in X$ satisfying (6.16) and a Lagrange multiplier vector μ^* . By the definition of a Lagrange multiplier vector, $f^* \leq L(\bar{x}, \mu^*)$ holds, which implies that

$$\sum_{i=1}^m \mu_i^* \leq \frac{[f(\bar{x}) - f^*]}{\min_{i=1, \dots, m} \{-g_i(\bar{x})\}}.$$

Since $\mu^* \geq 0^m$, boundedness follows. As by Proposition 6.6(a) the set of Lagrange multipliers is the set of optimal solutions to the dual problem (6.10), convexity follows from the identification of the dual solution set with the set of vectors $\mu \in \mathbb{R}_+^m$ for which

$$q(\mu) \geq q^*$$

holds. This is the upper level set for q at the level q^* ; this set is convex, by the concavity of q (cf. Theorem 6.4 and Proposition 3.44).

Lagrangian duality

(b) The result follows from Theorem 6.7.

(c) The first part follows from Theorem 4.24, as the Lagrangian function $L(\cdot, \boldsymbol{\mu}^*)$ is convex. The second part follows by identification. ■

Consider next the extension of the inequality constrained convex program (6.4) in which we seek to find

$$f^* := \infimum_x f(\mathbf{x}), \quad (6.21)$$

$$\begin{aligned} \text{subject to} \quad & \mathbf{x} \in X, \\ & g_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m, \\ & \boldsymbol{\varepsilon}_j^T \mathbf{x} - d_j = 0, \quad j = 1, \dots, \ell, \end{aligned}$$

under the same conditions as stated following (6.4), and where $\boldsymbol{\varepsilon}_j \in \mathbb{R}^n$, $j = 1, \dots, \ell$. For this problem, we replace the Slater condition (6.16) with the following (cf. [BSS93, Theorem 6.2.4]):

$$\exists \mathbf{x} \in X \text{ with } \mathbf{g}(\mathbf{x}) < \mathbf{0}^m \text{ and } \mathbf{0}^m \in \text{int} \{ \mathbf{E}\mathbf{x} - \mathbf{d} \mid \mathbf{x} \in X \}, \quad (6.22)$$

where $\mathbf{E} \in \mathbb{R}^{\ell \times n}$ has rows $\boldsymbol{\varepsilon}_j^T$, and $\mathbf{d} = (d_j)_{j \in \{1, \dots, \ell\}} \in \mathbb{R}^\ell$.

Note that in the statement (6.22), the “int” can be stricken whenever X is polyhedral, so that the latter part simply states that $\mathbf{E}\mathbf{x} = \mathbf{d}$.

For this problem, the Lagrangian dual problem is to find

$$q^* := \supremum_{(\boldsymbol{\mu}, \boldsymbol{\lambda})} q(\boldsymbol{\mu}, \boldsymbol{\lambda}), \quad (6.23)$$

$$\text{subject to } \boldsymbol{\mu} \geq \mathbf{0}^m,$$

where

$$\begin{aligned} q(\boldsymbol{\mu}, \boldsymbol{\lambda}) &:= \infimum_x L(\mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\lambda}) := f(\mathbf{x}) + \boldsymbol{\mu}^T \mathbf{g}(\mathbf{x}) + \boldsymbol{\lambda}^T (\mathbf{E}\mathbf{x} - \mathbf{d}), \\ &\text{subject to } \mathbf{x} \in X. \end{aligned}$$

Theorem 6.10 (Strong Duality, general convex programs) *Suppose that in addition to the feasibility condition (6.5), Slater’s constraint qualification (6.22) holds for the problem (6.21).*

(a) *The duality gap is zero and there exists at least one Lagrange multiplier vector pair $(\boldsymbol{\mu}^*, \boldsymbol{\lambda}^*)$.*

(b) *If the infimum in (6.21) is attained at some \mathbf{x}^* , then the triple $(\mathbf{x}^*, \boldsymbol{\mu}^*, \boldsymbol{\lambda}^*)$ satisfies the global optimality conditions*

$$\boldsymbol{\mu}^* \geq \mathbf{0}^m, \quad (\text{Dual feasibility}) \quad (6.24a)$$

$$\mathbf{x}^* \in \arg \min_{\mathbf{x} \in X} L(\mathbf{x}, \boldsymbol{\mu}^*, \boldsymbol{\lambda}^*), \quad (\text{Lagrangian optimality}) \quad (6.24b)$$

$$\mathbf{x}^* \in X, \mathbf{g}(\mathbf{x}^*) \leq \mathbf{0}^m, \mathbf{E}\mathbf{x}^* = \mathbf{d}, \quad (\text{Primal feasibility}) \quad (6.24c)$$

$$\mu_i^* g_i(\mathbf{x}^*) = 0, \quad i = 1, \dots, m. \quad (\text{Complementary slackness}) \quad (6.24d)$$

(c) If further f and \mathbf{g} are differentiable at \mathbf{x}^* , then the condition (6.24b) can equivalently be written as

$$\nabla_{\mathbf{x}} L(\mathbf{x}^*, \boldsymbol{\mu}^*, \boldsymbol{\lambda}^*)^\top (\mathbf{x} - \mathbf{x}^*) \geq 0, \quad \mathbf{x} \in X. \quad (6.25)$$

If, in addition, X is open (such as is the case when $X = \mathbb{R}^n$), then this reduces to the condition that

$$\nabla_{\mathbf{x}} L(\mathbf{x}^*, \boldsymbol{\mu}^*, \boldsymbol{\lambda}^*) = \nabla f(\mathbf{x}^*) + \sum_{i=1}^m \mu_i^* \nabla g_i(\mathbf{x}^*) + \sum_{j=1}^{\ell} \lambda_j^* \boldsymbol{\varepsilon}_j = \mathbf{0}^n, \quad (6.26)$$

and the global optimality conditions (6.24) reduce to the Karush–Kuhn–Tucker conditions stated in Theorem 5.33.

Proof. The proof is similar to that of Theorem 6.9. ■

We finally consider a special case where automatically a regularity condition holds.

Consider the affinely constrained convex program to find

$$\begin{aligned} f^* &:= \inf_{\mathbf{x}} f(\mathbf{x}), \\ \text{subject to} \quad &\mathbf{x} \in X, \\ &\mathbf{a}_i^\top \mathbf{x} - b_i \leq 0, \quad i = 1, \dots, m, \\ &\boldsymbol{\varepsilon}_j^\top \mathbf{x} - d_j = 0, \quad j = 1, \dots, \ell, \end{aligned} \quad (6.27)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex and $X \subseteq \mathbb{R}^n$ is polyhedral.

Theorem 6.11 (Strong Duality, affine constraints) *If the feasibility condition (6.5) holds for the problem (6.27), then there is no duality gap and there exists at least one Lagrange multiplier vector.*

Proof. Again, the proof is similar to that of Theorem 6.9, except that no additional regularity conditions are needed.⁴ ■

The existence of a multiplier vector [which by Proposition 6.6 and the absence of a duality gap implies the existence of an optimal solution to the dual problem (6.10)] does not imply the existence of an optimal solution to the primal problem (6.27) without any additional assumptions. However, when f is either weakly coercive, quadratic or linear, the existence results are stronger; see the primal existence results in Theorems 4.7, 4.8, and 6.12 below, for example.

⁴For a detailed proof, see [Ber99, Proposition 5.2.1]. (The special case where f is moreover differentiable is covered in [Ber99, Proposition 3.4.2].)

For convex programs where a Slater CQ holds, the Lagrange multipliers defined in this section, and those that appear in the Karush–Kuhn–Tucker conditions, clearly are identical. Next, we specialize the above to linear and quadratic programs.

6.2.4 Strong duality for linear and quadratic programs

The following result will be established and analyzed in detail in Chapter 10 on linear programming duality (cf. Theorem 10.6), but can in fact also be established similarly to above. (See [BSS93, Theorem 2.7.3] or [Ber99, Proposition 5.2.2], for example.) Its proof will however be relegated to that of Theorem 10.6.

Theorem 6.12 (Strong Duality, linear programs) *Assume, in addition to the conditions of Theorem 6.11, that f is linear, so that (6.27) is a linear program. Then, the primal and dual problems have optimal solutions and there is no duality gap.* ■

The above result states a strong duality result for a general linear program. We next develop an explicit Lagrangian dual problem for a linear program.

Let $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{c} \in \mathbb{R}^n$, and $\mathbf{b} \in \mathbb{R}^m$; consider the linear program

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} \quad \mathbf{c}^T \mathbf{x}, \\ & \text{subject to} \quad \mathbf{A}\mathbf{x} = \mathbf{b}, \\ & \quad \quad \quad \mathbf{x} \geq \mathbf{0}^n. \end{aligned} \tag{6.28}$$

If we let $X := \mathbb{R}_+^n$, then the Lagrangian dual problem is to

$$\begin{aligned} & \underset{\boldsymbol{\lambda} \in \mathbb{R}^m}{\text{maximize}} \quad \mathbf{b}^T \boldsymbol{\lambda}, \\ & \text{subject to} \quad \mathbf{A}^T \boldsymbol{\lambda} \leq \mathbf{c}. \end{aligned} \tag{6.29}$$

The reason why we can write it in this form is that

$$q(\boldsymbol{\lambda}) := \inf_{\mathbf{x} \geq \mathbf{0}^n} \left\{ \mathbf{c}^T \mathbf{x} + \boldsymbol{\lambda}^T (\mathbf{b} - \mathbf{A}\mathbf{x}) \right\} = \mathbf{b}^T \boldsymbol{\lambda} + \inf_{\mathbf{x} \geq \mathbf{0}^n} (\mathbf{c} - \mathbf{A}^T \boldsymbol{\lambda})^T \mathbf{x},$$

so that

$$q(\boldsymbol{\lambda}) = \begin{cases} \mathbf{b}^T \boldsymbol{\lambda}, & \text{if } \mathbf{A}^T \boldsymbol{\lambda} \leq \mathbf{c}, \\ -\infty, & \text{otherwise.} \end{cases}$$

(The infimum is attained at zero if and only if these inequalities are satisfied; otherwise, the inner problem is unbounded below.)

Further, why is it that λ here is not restricted in sign? Suppose we were to split the system $Ax = b$ into an inequality system of the form

$$\begin{aligned} Ax &\leq b, \\ -Ax &\leq -b. \end{aligned}$$

Let $((\mu^+)^T, (\mu^-)^T)^T$ be the corresponding vector of multipliers, and take the Lagrangian dual for this formulation. Then, we would have a Lagrange function of the form

$$(x, \mu^+, \mu^-) \mapsto L(x, \mu^+, \mu^-) := c^T x + (\mu^+ - \mu^-)^T (b - Ax),$$

and since $\mu^+ - \mu^-$ can take on any value in \mathbb{R}^m we can simply replace it with the unrestricted vector $\lambda \in \mathbb{R}^m$. This motivates why the multiplier for an equality constraint never is sign restricted; the same was the case, as we saw in Section 5.6, for the multipliers in the KKT conditions.

As applied to this problem, Theorem 6.12 states that if both the primal or dual problems have feasible solutions, then they both have optimal solutions, satisfying strong duality ($c^T x^* = b^T \lambda^*$). On the other hand, if any of the two problems has an unbounded solution, then the other problem is infeasible.

Consider next the quadratic programming problem to

$$\begin{aligned} &\underset{x}{\text{minimize}} \quad \left\{ \frac{1}{2} x^T Q x + c^T x \right\}, \\ &\text{subject to} \quad Ax \leq b, \end{aligned} \tag{6.30}$$

where $Q \in \mathbb{R}^{n \times n}$, $c \in \mathbb{R}^n$, $A \in \mathbb{R}^{m \times n}$, and $b \in \mathbb{R}^m$. We develop an explicit dual problem under the assumption that Q is positive definite.

By Lagrangian relaxing the inequality constraints, we obtain that the inner problem in x is solved by letting

$$x = -Q^{-1}(c + A^T \mu). \tag{6.31}$$

Substituting this expression into the Lagrangian function yields the Lagrangian dual problem to

$$\begin{aligned} &\underset{\mu}{\text{maximize}} \quad \left\{ -\frac{1}{2} \mu^T A Q^{-1} A^T \mu - (b + A Q^{-1} c)^T \mu - \frac{1}{2} c^T Q^{-1} c \right\}, \\ &\text{subject to} \quad \mu \geq 0^m, \end{aligned} \tag{6.32}$$

Strong duality follows for this convex primal–dual pair of quadratic programs, in much the same way as for linear programs.

Theorem 6.13 (Strong Duality, quadratic programs) *For the primal–dual pair of convex quadratic programs (6.30), (6.32), the following holds:*

(a) *If both problems have feasible solutions, then both problems also have optimal solutions, and the primal problem (6.30) also has a unique optimal solution, given by (6.31) for any optimal Lagrange multiplier vector, and in the two problems the optimal values are equal.*

(b) *If either of the two problems has an unbounded solution, then the other one is infeasible.*

(c) *Suppose that \mathbf{Q} is positive semidefinite, and that the feasibility condition (6.5) holds. Then, both the problem (6.30) and its Lagrangian dual have nonempty, closed and convex sets of optimal solutions, and their optimal values are equal.* ■

In the result (a) it is important to note that the Lagrangian dual problem (6.32) is not necessarily strictly convex; the matrix $\mathbf{A}\mathbf{Q}^{-1}\mathbf{A}^T$ need not be positive definite, especially so when \mathbf{A} does not have full rank. The result (c) extends the strong duality result from linear programming, since \mathbf{Q} in (c) can be the zero matrix. In the case of (c) we of course cannot write the Lagrangian dual problem in the form of (6.32) because \mathbf{Q} is not necessarily invertible.

6.2.5 Two illustrative examples

Example 6.14 (an explicit, differentiable dual problem) Consider the problem to

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && f(\mathbf{x}) := x_1^2 + x_2^2, \\ & \text{subject to} && x_1 + x_2 \geq 4, \\ & && x_j \geq 0, \quad j = 1, 2. \end{aligned}$$

We consider the first constraint to be the complicated one, and hence define $g(\mathbf{x}) := -x_1 - x_2 + 4$ and let $X := \{(x_1, x_2)^T \mid x_j \geq 0, j = 1, 2\}$.

Then, the Lagrangian dual function is

$$\begin{aligned} q(\mu) &= \underset{\mathbf{x} \in X}{\text{minimum}} L(\mathbf{x}, \mu) := f(\mathbf{x}) - \mu(x_1 + x_2 - 4) \\ &= 4\mu + \underset{\mathbf{x} \in X}{\text{minimum}} \{x_1^2 + x_2^2 - \mu x_1 - \mu x_2\} \\ &= 4\mu + \underset{x_1 \geq 0}{\text{minimum}} \{x_1^2 - \mu x_1\} + \underset{x_2 \geq 0}{\text{minimum}} \{x_2^2 - \mu x_2\}, \quad \mu \geq 0. \end{aligned}$$

For a fixed $\mu \geq 0$, the minimum is attained at $x_1(\mu) = \frac{\mu}{2}, x_2(\mu) = \frac{\mu}{2}$.

Substituting this expression yields $q(\mu) = f(\mathbf{x}(\mu)) - \mu(x_1(\mu) + x_2(\mu) - 4) = 4\mu - \frac{\mu^2}{2}$.

Note that q is strictly concave, and it is differentiable everywhere (due to the fact that f, g are differentiable and $\mathbf{x}(\mu)$ is unique), by Danskin's Theorem 6.17(d).

We have that $q'(\mu) = 4 - \mu = 0 \iff \mu = 4$. As $\mu = 4 \geq 0$, it is the optimum in the dual problem: $\mu^* = 4$; $\mathbf{x}^* = (x_1(\mu^*), x_2(\mu^*))^T = (2, 2)^T$. Also, $f(\mathbf{x}^*) = q(\mu^*) = 8$.

This is an example where the dual function is differentiable, and therefore we can utilize Proposition 6.29(c). In this case, the optimum \mathbf{x}^* is also unique, so it is automatically given as $\mathbf{x}^* = \mathbf{x}(\mu)$. ■

Example 6.15 (an implicit, non-differentiable dual problem) Consider the linear programming problem to

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && f(\mathbf{x}) := -x_1 - x_2, \\ & \text{subject to} && 2x_1 + 4x_2 \leq 3, \\ & && 0 \leq x_1 \leq 2, \\ & && 0 \leq x_2 \leq 1. \end{aligned}$$

The optimal solution is $\mathbf{x}^* = (3/2, 0)^T$, $f(\mathbf{x}^*) = -3/2$.

Consider Lagrangian relaxing the first constraint, obtaining

$$\begin{aligned} L(\mathbf{x}, \mu) &= -x_1 - x_2 + \mu(2x_1 + 4x_2 - 3); \\ q(\mu) &= -3\mu + \underset{0 \leq x_1 \leq 2}{\text{minimum}} \{(-1 + 2\mu)x_1\} + \underset{0 \leq x_2 \leq 1}{\text{minimum}} \{(-1 + 4\mu)x_2\} \\ &= \begin{cases} -3 + 5\mu, & 0 \leq \mu \leq 1/4, \\ -2 + \mu, & 1/4 \leq \mu \leq 1/2, \\ -3\mu, & 1/2 \leq \mu. \end{cases} \end{aligned}$$

Check that $\mu^* = 1/2$, and hence that $q(\mu^*) = -3/2$. For linear programs, we have strong duality, but how do we obtain the optimal primal solution from μ^* ? It is clear that q is non-differentiable at μ^* . Let us utilize the characterization given in the system (6.12).

First, at μ^* , it is clear that $X(\mu^*)$ is the set $\{(2\alpha, 0)^T \mid 0 \leq \alpha \leq 1\}$. Among the subproblem solutions, we next have to find one that is primal feasible as well as complementary.

Primal feasibility means that $2 \cdot 2\alpha + 4 \cdot 0 \leq 3 \iff \alpha \leq 3/4$.

Further, complementarity means that $\mu^* \cdot (2x_1^* + 4x_2^* - 3) = 0 \iff \alpha = 3/4$, since $\mu^* \neq 0$. We conclude that the only primal vector that satisfies the system (6.12) together with the dual optimal solution $\mu^* = 1/2$ is $\mathbf{x}^* = (3/2, 0)^T$. ■

In the first example, the Lagrangian dual function is differentiable since $\mathbf{x}(\mu)$ is unique. The second one shows that otherwise, there may

be kinks in the function q where there are alternative solutions $\mathbf{x}(\mu)$; as a result, to obtain a primal optimal solution becomes more complex. The Dantzig–Wolfe algorithm, for example, represents a means by which to automatize the process that we have just shown; the algorithm generates extreme points of $X(\mu)$ algorithmically, and constructs the best feasible convex combination thereof, obtaining a primal–dual optimal solution in a finite number of iterations for linear programs.

The above examples motivate a deeper study of the differentiability properties of convex (or, concave) functions in general, and the Lagrangian dual objective function in particular.

6.3 Differentiability properties of the dual function

We have established that the Lagrangian dual problem (6.10) is a convex one, and further that under some circumstances the primal and dual optimal values are the same. We now turn to study the Lagrangian dual problem in detail, and in particular how it can be solved efficiently. First, we will establish when the dual function q is differentiable. We will see that differentiability holds only in some special cases, in which we can recognize the workings of the *Lagrange multiplier method*; this classic method was illustrated in Example 6.14. Most often, the function q will however be non-differentiable, and then this method will fail. This means that we must devise a more general numerical method which is not based on gradients but rather *subgradients*. This type of algorithm is the topic of the next section; we begin by studying the topic of subgradients of convex functions in general.

6.3.1 Subdifferentiability of convex functions

Throughout this section we suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a convex function, and study its subdifferentiability properties. We will later on apply our findings to the Lagrangian dual function q , or, rather, its negative $-q$. We first remark that a finite convex function is automatically continuous (cf. Theorem 4.27).

Definition 6.16 (subgradient) *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function. We say that a vector $\mathbf{g} \in \mathbb{R}^n$ is a subgradient of f at $\mathbf{x} \in \mathbb{R}^n$ if*

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^T(\mathbf{y} - \mathbf{x}), \quad \mathbf{y} \in \mathbb{R}^n. \quad (6.33)$$

The set of such vectors \mathbf{g} defines the subdifferential of f at \mathbf{x} , and is denoted $\partial f(\mathbf{x})$. ■

For concave functions, the reverse inequality of course holds; for simplicity we will refer also to such vectors \mathbf{g} as subgradients.

Notice the connection to the characterization of a convex function in C^1 in Theorem 3.40(a). The difference between them is that \mathbf{g} is not unique at a non-differentiable point. (Just as the gradient has a role in supporting hyperplanes to the graph of a convex function in C^1 , the role of a subgradient is the same; at a non-differentiable point there are more than one supporting hyperplane to the graph of f .)

We illustrate this in Figure 6.2.

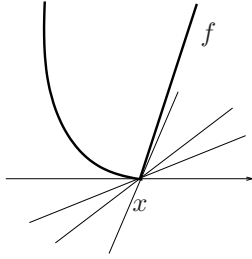


Figure 6.2: Three possible slopes of the convex function f at x .

Notice that a minimum \mathbf{x}^* of f over \mathbb{R}^n is characterized by the inclusion $\mathbf{0}^n \in \partial f(\mathbf{x}^*)$; recognize, again, the similarity to the C^1 case.

We list some additional basic results for convex functions next. Proofs will not be given here; we refer instead to the convex analysis text by Rockafellar [Roc70].

Proposition 6.17 (properties of a convex function) *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function.*

(a) [boundedness of $\partial f(\mathbf{x})$] *For every $\mathbf{x} \in \mathbb{R}^n$, $\partial f(\mathbf{x})$ is a nonempty, convex, and compact set. If X is bounded then $\cup_{\mathbf{x} \in X} \partial f(\mathbf{x})$ is bounded.*

(b) [closedness of ∂f] *The subdifferential mapping $\mathbf{x} \mapsto \partial f(\mathbf{x})$ is closed; in other words, if $\{\mathbf{x}_k\}$ is a sequence of vectors in \mathbb{R}^n converging to \mathbf{x} , and $\mathbf{g}_k \in \partial f(\mathbf{x}_k)$ holds for every k , then the sequence $\{\mathbf{g}_k\}$ of subgradients is bounded and every limit point thereof belongs to $\partial f(\mathbf{x})$.*

(c) [directional derivative and differentiability] *For every $\mathbf{x} \in \mathbb{R}^n$, the directional derivative of f at \mathbf{x} in the direction of $\mathbf{p} \in \mathbb{R}^n$ satisfies*

$$f'(\mathbf{x}; \mathbf{p}) = \text{maximum}_{\mathbf{g} \in \partial f(\mathbf{x})} \mathbf{g}^T \mathbf{p}. \quad (6.34)$$

In particular, f is differentiable at \mathbf{x} with gradient $\nabla f(\mathbf{x})$ if and only if it has $\nabla f(\mathbf{x})$ as its unique subgradient at \mathbf{x} ; in that case, $f'(\mathbf{x}; \mathbf{p}) = \nabla f(\mathbf{x})^T \mathbf{p}$.

Lagrangian duality

(d) [Danskin's Theorem—directional derivatives of a convex max function] Let Z be a compact subset of \mathbb{R}^m , and let $\phi : \mathbb{R}^n \times Z \rightarrow \mathbb{R}$ be continuous and such that $\phi(\cdot, \mathbf{z}) : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex for each $\mathbf{z} \in Z$. Let the function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be given by

$$f(\mathbf{x}) := \max_{\mathbf{z} \in Z} \phi(\mathbf{x}, \mathbf{z}), \quad \mathbf{x} \in \mathbb{R}^n. \quad (6.35)$$

The function f then is convex on \mathbb{R}^n and has a directional derivative at \mathbf{x} in the direction of \mathbf{p} equal to

$$f'(\mathbf{x}; \mathbf{p}) = \max_{\mathbf{z} \in Z(\mathbf{x})} \phi'(\mathbf{x}, \mathbf{z}; \mathbf{p}), \quad (6.36)$$

where $\phi'(\mathbf{x}, \mathbf{z}; \mathbf{p})$ is the directional derivative of $\phi(\cdot, \mathbf{z})$ at \mathbf{x} in the direction of \mathbf{p} , and $Z(\mathbf{x}) := \{\mathbf{z} \in \mathbb{R}^m \mid \phi(\mathbf{x}, \mathbf{z}) = f(\mathbf{x})\}$.

In particular, if $Z(\mathbf{x})$ contains a single point $\bar{\mathbf{z}}$ and $\phi(\cdot, \bar{\mathbf{z}})$ is differentiable at \mathbf{x} , then f is differentiable at \mathbf{x} , and $\nabla f(\mathbf{x}) = \nabla_{\mathbf{x}} \phi(\mathbf{x}, \bar{\mathbf{z}})$, where $\nabla_{\mathbf{x}} \phi(\mathbf{x}, \bar{\mathbf{z}})$ is the vector with components $\frac{\partial \phi(\mathbf{x}, \bar{\mathbf{z}})}{\partial x_i}$, $i = 1, \dots, n$.

If further $\phi(\cdot, \mathbf{z})$ is differentiable for all $\mathbf{z} \in Z$ and $\nabla_{\mathbf{x}} \phi(\mathbf{x}, \cdot)$ is continuous on Z for each \mathbf{x} , then

$$\partial f(\mathbf{x}) = \text{conv} \{ \nabla_{\mathbf{x}} \phi(\mathbf{x}, \mathbf{z}) \mid \mathbf{z} \in Z(\mathbf{x}) \}, \quad \mathbf{x} \in \mathbb{R}^n.$$

Proof. (a) This is a special case of [Roc70, Theorem 24.7].

(b) This is [Roc70, Theorem 24.5].

(c) This is [Roc70, Theorem 23.4 and 25.1].

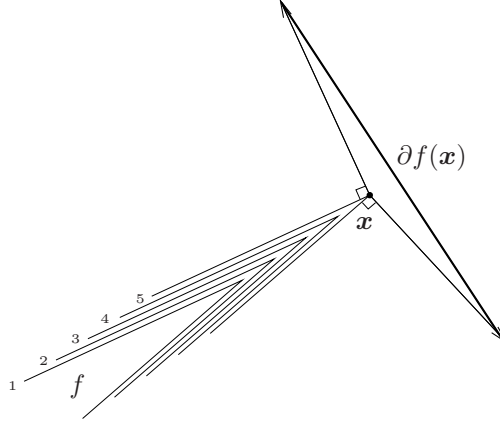
(d) This is [Ber99, Proposition B.25]. ■

Figure 6.3 illustrates the subdifferential of a convex function.

We apply parts of the above results in order to characterize a minimum of a convex function on \mathbb{R}^n .

Proposition 6.18 (optimality of a convex function over \mathbb{R}^n) *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function. The following three statements are equivalent:*

1. f is globally minimized at $\mathbf{x}^* \in \mathbb{R}^n$;
2. $\mathbf{0}^n \in \partial f(\mathbf{x}^*)$;
3. $f'(\mathbf{x}^*; \mathbf{p}) \geq 0$ for all $\mathbf{p} \in \mathbb{R}^n$.


 Figure 6.3: The subdifferential of a convex function f at \mathbf{x} .

Proof. We establish the result thus: $1 \implies 2 \implies 3 \implies 1$.

[$1 \implies 2$]: By the statement 1., we have that $f(\mathbf{y}) \geq f(\mathbf{x}^*)$ for every $\mathbf{y} \in \mathbb{R}^n$. This implies that for $\mathbf{g} = \mathbf{0}^n$, we satisfy the subgradient inequality (6.33). This establishes the statement 2.

[$2 \implies 3$]: We can equivalently write

$$\partial f(\mathbf{x}) = \{ \mathbf{g} \in \mathbb{R}^n \mid \mathbf{g}^T \mathbf{p} \leq f'(\mathbf{x}; \mathbf{p}), \quad \mathbf{p} \in \mathbb{R}^n \}.$$

With $\mathbf{g} = \mathbf{0}^n$ this definition immediately yields the statement 3.

[$3 \implies 1$]: By the compactness of the subdifferential [cf. Proposition 6.17(a)] and Weierstrass' Theorem 4.7 the maximum in the expression (6.34) is attained at some $\mathbf{g} \in \partial f(\mathbf{x}^*)$. It follows that, in the subgradient inequality (6.33), we get that

$$f(\mathbf{x}^* + \mathbf{p}) \geq f(\mathbf{x}^*) + \mathbf{g}^T \mathbf{p} \geq f(\mathbf{x}^*), \quad \mathbf{p} \in \mathbb{R}^n,$$

which is equivalent to the statement 1. ■

This result implies that a direction $\mathbf{p} \in \mathbb{R}^n$ is a descent direction with respect to f at \mathbf{x} if and only if $f'(\mathbf{x}; \mathbf{p}) < 0$ holds. This result cannot be extended to non-convex functions, even when the function f is in C^1 or even C^2 . [Take $f(x) := x^3$; $x = 0$; $p = -1$; see also the discussions following Proposition 4.16 and on saddle points in Example 11.2(b).] (A related result for possibly non-convex but differentiable functions is found in Proposition 4.16.)

6.3.2 Differentiability of the Lagrangian dual function

We consider the inequality constrained problem (6.4), where we make the following standing assumption:

$$f, g_i (i = 1, \dots, m) \in C^0, X \text{ is nonempty and compact.} \quad (6.37)$$

Under this assumption, the set of solutions to the Lagrangian subproblem,

$$X(\boldsymbol{\mu}) := \arg \min_{\mathbf{x} \in X} L(\mathbf{x}, \boldsymbol{\mu}), \quad \boldsymbol{\mu} \in \mathbb{R}^m, \quad (6.38)$$

is nonempty and compact for any choice of dual vector $\boldsymbol{\mu}$ by Weierstrass' Theorem 4.7. We first develop the subdifferentiability properties of the associated dual function q , stated in (6.9). The first result strengthens Theorem 6.4 under these additional assumptions.

Proposition 6.19 (subdifferentiability of the dual function) *Suppose that, in the problem (6.4), the compactness condition (6.37) holds.*

(a) *The dual function (6.9) is finite, continuous and concave on \mathbb{R}^m . If its supremum over \mathbb{R}_+^m is attained, then the optimal solution set therefore is closed and convex.*

(b) *The mapping $\boldsymbol{\mu} \mapsto X(\boldsymbol{\mu})$ is closed on \mathbb{R}^m . If $X(\bar{\boldsymbol{\mu}})$ is the singleton set $\{\bar{\mathbf{x}}\}$ for some $\bar{\boldsymbol{\mu}} \in \mathbb{R}^m$, and for some sequence $\{\boldsymbol{\mu}_k\} \subset \mathbb{R}^m$ with $\boldsymbol{\mu}_k \rightarrow \bar{\boldsymbol{\mu}}$, $\mathbf{x}_k \in X(\boldsymbol{\mu}_k)$ for all k , then $\mathbf{x}_k \rightarrow \bar{\mathbf{x}}$.*

(c) *Let $\boldsymbol{\mu} \in \mathbb{R}^m$. If $\mathbf{x} \in X(\boldsymbol{\mu})$, then $\mathbf{g}(\mathbf{x})$ is a subgradient to q at $\boldsymbol{\mu}$, that is, $\mathbf{g}(\mathbf{x}) \in \partial q(\boldsymbol{\mu})$.*

(d) *Let $\boldsymbol{\mu} \in \mathbb{R}^m$. Then,*

$$\partial q(\boldsymbol{\mu}) = \text{conv} \{ \mathbf{g}(\mathbf{x}) \mid \mathbf{x} \in X(\boldsymbol{\mu}) \}.$$

The set $\partial q(\boldsymbol{\mu})$ is convex and compact. Moreover, if U is a bounded set, then $\cup_{\boldsymbol{\mu} \in U} \partial q(\boldsymbol{\mu})$ is also bounded.

(e) *The directional derivative of q at $\boldsymbol{\mu} \in \mathbb{R}^m$ in the direction of $\mathbf{p} \in \mathbb{R}^m$ is*

$$q'(\boldsymbol{\mu}; \mathbf{p}) = \min_{\mathbf{g} \in \partial q(\boldsymbol{\mu})} \mathbf{g}^T \mathbf{p}.$$

Proof. (a) Theorem 6.4 stated the concavity of q on its effective domain. Weierstrass' Theorem 4.7 states that q is finite on \mathbb{R}^m , which is then also its effective domain. The continuity of q follows from that of any finite concave function, as we have already seen in Theorem 4.27. The closedness property of the solution set is a direct consequence of the continuity of q (the upper level set then automatically is closed), and complements the result of Theorem 6.9(a).

(b) Let $\{\boldsymbol{\mu}_k\}$ be a sequence of vectors in \mathbb{R}^m converging to $\bar{\boldsymbol{\mu}}$, and let $\mathbf{x}_k \in X(\boldsymbol{\mu}_k)$ be arbitrary. Let \mathbf{x} be arbitrary in X , and let further $\bar{\mathbf{x}} \in X$ be an arbitrary limit point of $\{\mathbf{x}_k\}$ (at least one exists by the compactness of X). From the property that for all k ,

$$L(\mathbf{x}_k, \boldsymbol{\mu}_k) \leq L(\mathbf{x}, \boldsymbol{\mu}_k),$$

follows, by the continuity of L , that, in the limit of k in the subsequence in which \mathbf{x}_k converges to $\bar{\mathbf{x}}$,

$$L(\bar{\mathbf{x}}, \bar{\boldsymbol{\mu}}) \leq L(\mathbf{x}, \bar{\boldsymbol{\mu}}),$$

so that $\bar{\mathbf{x}} \in X(\bar{\boldsymbol{\mu}})$, as desired. The special case of a singleton set $X(\bar{\boldsymbol{\mu}})$ follows.

(c) Let $\bar{\boldsymbol{\mu}} \in \mathbb{R}^m$ be arbitrary and let $\bar{\mathbf{x}} \in X(\bar{\boldsymbol{\mu}})$. We have that

$$\begin{aligned} q(\bar{\boldsymbol{\mu}}) &= \inf_{\mathbf{y} \in X} L(\mathbf{y}, \bar{\boldsymbol{\mu}}) = f(\bar{\mathbf{x}}) + \bar{\boldsymbol{\mu}}^T \mathbf{g}(\bar{\mathbf{x}}) \\ &= f(\bar{\mathbf{x}}) + \boldsymbol{\mu}^T \mathbf{g}(\bar{\mathbf{x}}) + (\bar{\boldsymbol{\mu}} - \boldsymbol{\mu})^T \mathbf{g}(\bar{\mathbf{x}}) \geq q(\boldsymbol{\mu}) + (\bar{\boldsymbol{\mu}} - \boldsymbol{\mu})^T \mathbf{g}(\bar{\mathbf{x}}), \end{aligned}$$

which implies that $\mathbf{g}(\bar{\mathbf{x}}) \in \partial q(\boldsymbol{\mu})$.

(d) The inclusion $\partial q(\boldsymbol{\mu}) \subseteq \text{conv} \{ \mathbf{g}(\mathbf{x}) \mid \mathbf{x} \in X(\boldsymbol{\mu}) \}$ follows from (c) and the convexity of $\partial q(\boldsymbol{\mu})$. The opposite inclusion follows by applying the Separation Theorem 3.24.⁵

(e) See Proposition 6.17(c). ■

The result in (c) is an independent proof of the concavity of q on \mathbb{R}^m .

The result (d) is particularly interesting, because by Carathéodory's Theorem 3.8 every subgradient of q at any point $\boldsymbol{\mu}$ is the convex combination of a finite number (in fact, at most $m+1$) of vectors of the form $\mathbf{g}(\mathbf{x}^s)$ with $\mathbf{x}^s \in X(\boldsymbol{\mu})$. Computationally, this has been utilized to devise efficient (proximal) bundle methods for the Lagrangian dual problem as well as to devise methods to recover primal optimal solutions.

Next, we establish the differentiability of the dual function under additional assumptions.

Proposition 6.20 (differentiability of the dual function) *Suppose that, in the problem (6.4), the compactness condition (6.37) holds.*

(a) *Let $\boldsymbol{\mu} \in \mathbb{R}^m$. The dual function q is differentiable at $\boldsymbol{\mu}$ if and only if $\{ \mathbf{g}(\mathbf{x}) \mid \mathbf{x} \in X(\boldsymbol{\mu}) \}$ is a singleton set, that is, if the value of the vector of constraint functions is invariant over the set of solutions $X(\boldsymbol{\mu})$ to the Lagrangian subproblem. Then, we have that*

$$\nabla q(\boldsymbol{\mu}) = \mathbf{g}(\mathbf{x}),$$

⁵See [BSS93, Theorem 6.3.7] for a detailed proof.

for every $\mathbf{x} \in X(\boldsymbol{\mu})$.

(b) The result in (a) holds in particular if the Lagrangian subproblem has a unique solution, that is, $X(\boldsymbol{\mu})$ is a singleton set. In particular, this property is satisfied for $\boldsymbol{\mu} \geq \mathbf{0}^m$ if further X is a convex set, f is strictly convex on X , and g_i ($i = 1, \dots, m$) are convex, in which case $q \in C^1$.

Proof. (a) The concave function q is differentiable at the point $\boldsymbol{\mu}$ (where it is finite) if and only if its subdifferential $\partial q(\boldsymbol{\mu})$ there is a singleton, cf. Proposition 6.17(c).

(b) Under either one of the assumptions stated, $X(\boldsymbol{\mu})$ is a singleton, whence the result follows from (a). Uniqueness follows from the convexity of the feasible set and strict convexity of the objective function, according to Proposition 4.11. That $q \in C^1$ follows from the continuity of \mathbf{g} and Proposition 6.19(b). ■

Proposition 6.21 (twice differentiability of the dual objective function) Suppose that, in the problem (6.4), $X = \mathbb{R}^n$, and f and g_i ($i = 1, \dots, m$) are convex functions in C^2 . Suppose that, at $\boldsymbol{\mu} \in \mathbb{R}^m$, the solution \mathbf{x} to the Lagrangian subproblem not only is unique, but also that the partial Hessian of the Lagrangian is positive definite at the pair $(\mathbf{x}, \boldsymbol{\mu})$, that is,

$$\nabla_{xx}^2 L(\mathbf{x}, \boldsymbol{\mu}) \text{ is positive definite.}$$

Then, the dual function q is twice differentiable at $\boldsymbol{\mu}$, with

$$\nabla^2 q(\boldsymbol{\mu}) = -\nabla \mathbf{g}(\mathbf{x})^T [\nabla_{xx}^2 L(\mathbf{x}, \boldsymbol{\mu})]^{-1} \nabla \mathbf{g}(\mathbf{x}).$$

Proof. The result follows from the Implicit Function Theorem, which is stated in Chapter 2, applied to the Lagrangian subproblem.⁶ ■

6.4 *Subgradient optimization methods

We begin by establishing the convergence of classic subgradient optimization methods as applied to a general convex optimization problem.

6.4.1 Convex problems

Consider the convex optimization problem to

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}), \tag{6.39a}$$

$$\text{subject to} \quad \mathbf{x} \in X, \tag{6.39b}$$

⁶See [Ber99, Pages 596–598] for a detailed analysis.

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex and the set $X \subseteq \mathbb{R}^n$ is nonempty, closed and convex.

The subgradient projection algorithm is as follows: select $\mathbf{x}_0 \in X$, and for $k = 0, 1, \dots$ generate

$$\mathbf{g}_k \in \partial f(\mathbf{x}_k), \quad (6.40a)$$

$$\mathbf{x}_{k+1} = \text{Proj}_X(\mathbf{x}_k - \alpha_k \mathbf{g}_k), \quad (6.40b)$$

where the sequence $\{\alpha_k\}$ is generated from one of the following three rules:

The first rule is termed the *divergent series* step length rule, and requires that

$$\alpha_k > 0, \quad k = 0, 1, \dots; \quad \lim_{k \rightarrow \infty} \alpha_k = 0; \quad \sum_{k=0}^{\infty} \alpha_k = +\infty. \quad (6.41)$$

The second rule adds to the requirements in (6.41) the square-summable restriction

$$\sum_{k=0}^{\infty} \alpha_k^2 < +\infty. \quad (6.42)$$

The conditions in (6.41) allow for convergence to *any* point from *any* starting point, since the total step is infinite, but convergence is therefore also quite slow; the additional condition in (6.42) means fast sequences are selected. An instance of the step length formulas which satisfies both (6.41) and (6.42) is the following:

$$\alpha_k = \gamma + \beta/(k+1), \quad k = 0, 1, \dots,$$

where $\beta > 0, \gamma \geq 0$.

The third step length rule is

$$\alpha_k = \theta_k \frac{f(\mathbf{x}_k) - f^*}{\|\mathbf{g}_k\|^2}, \quad 0 < \sigma_1 \leq \theta_k \leq 2 - \sigma_2 < 2, \quad (6.43)$$

where f^* is the optimal value of (6.39). We refer to this step length formula as the *Polyak step*, after the Russian mathematician Boris Polyak who invented the subgradient method in the 1960s together with Ermol'ev and Shor.

How is convergence established for subgradient optimization methods? As shall be demonstrated in Chapters 11 and 12 convergence of algorithms for problems with a *differentiable* objective function is typically based on generating descent directions, and step length rules that result in the sequence $\{\mathbf{x}_k\}$ of iterates being strictly descending in the

value of f . For the non-differentiable problem at hand, generating descent directions is a difficult task, since it is not true that the negative of an arbitrarily chosen subgradient of f at a non-optimal vector \mathbf{x} defines a descent direction.

In *bundle methods* one gathers information from more than one subgradient (hence the term *bundle*) around a current iteration point so that a descent direction can be generated, followed by an inexact line search. We concentrate here on the simpler methodology of subgradient optimization methods, in which we apply the formula (6.40) where the step length α_k is chosen based on very simple rules.

We establish below that if the step length is small enough, an iteration of the subgradient projection method leads to a vector that is closer to the set of optimal solutions. This technical result also motivates the construction of the Polyak step length rule, and hence shows that the convergence of subgradient methods is based on the reduction of the Euclidean distance to the optimal solutions rather than on the reduction of the value of the objective function f .

Proposition 6.22 (decreasing distance to the optimal solution set) *Suppose that $\mathbf{x}_k \in X$ is not optimal in (6.39), and that \mathbf{x}_{k+1} is given by (6.40) for some step length $\alpha_k > 0$.*

Then, for every optimal solution \mathbf{x}^ in (6.39),*

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| < \|\mathbf{x}_k - \mathbf{x}^*\|$$

holds for every step length α_k in the interval

$$\alpha_k \in (0, 2[f(\mathbf{x}_k) - f^*]/\|\mathbf{g}_k\|^2). \quad (6.44)$$

Proof. We have that

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 &= \|\text{Proj}_X(\mathbf{x}_k - \alpha_k \mathbf{g}_k) - \mathbf{x}^*\|^2 \\ &= \|\text{Proj}_X(\mathbf{x}_k - \alpha_k \mathbf{g}_k) - \text{Proj}_X(\mathbf{x}^*)\|^2 \\ &\leq \|\mathbf{x}_k - \alpha_k \mathbf{g}_k - \mathbf{x}^*\|^2 \\ &= \|\mathbf{x}_k - \mathbf{x}^*\|^2 - 2\alpha_k(\mathbf{x}_k - \mathbf{x}^*)^\top \mathbf{g}_k + \alpha_k^2 \|\mathbf{g}_k\|^2 \\ &\leq \|\mathbf{x}_k - \mathbf{x}^*\|^2 - 2\alpha_k[f(\mathbf{x}_k) - f^*] + \alpha_k^2 \|\mathbf{g}_k\|^2 \\ &< \|\mathbf{x}_k - \mathbf{x}^*\|^2, \end{aligned}$$

where we have utilized the property that the Euclidean projection is non-expansive (Theorem 4.32), the subgradient inequality (6.33) for convex

functions, and the bounds on α_k given by (6.44). ■

Our first convergence result is based on the divergent series step length formula (6.41), and establishes convergence to the optimal solution set X^* under an assumption on its boundedness. With the other two step length formulas, this condition will be possible to remove.

Recall the definition (3.12) of the minimum distance from a vector to a closed and convex set; our interest is in the distance from an arbitrary vector $\mathbf{x} \in \mathbb{R}^n$ to the solution set X^* :

$$\text{dist}_{X^*}(\mathbf{x}) := \text{minimum}_{\mathbf{y} \in X^*} \|\mathbf{y} - \mathbf{x}\|.$$

Theorem 6.23 (convergence of subgradient optimization methods, I) *Let $\{\mathbf{x}_k\}$ be generated by the method (6.40), (6.41). If X^* is bounded and the sequence $\{\mathbf{g}_k\}$ is bounded, then $f(\mathbf{x}_k) \rightarrow f^*$ and $\text{dist}_{X^*}(\mathbf{x}_k) \rightarrow 0$ holds.*

Proof. We show that the iterates will eventually belong to an arbitrarily small neighbourhood of the set of optimal solutions to (6.39).

Let $\delta > 0$ and $B^\delta := \{\mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x}\| \leq \delta\}$. Since f is convex, X is nonempty, closed and convex, and X^* is bounded, it follows from [Roc70, Theorem 27.2], applied to the lower semi-continuous, proper⁷ and convex function $f + \chi_X$ ⁸ that there exists an $\varepsilon = \varepsilon(\delta) > 0$ such that the level set $\{\mathbf{x} \in X \mid f(\mathbf{x}) \leq f^* + \varepsilon\} \subseteq X^* + B^{\delta/2}$; this level set is denoted by X^ε . Moreover, since for all k , $\|\mathbf{g}_k\| \leq \sup_s \{\|\mathbf{g}_s\|\} < \infty$, and $\alpha_k \rightarrow 0$, there exists an $N(\delta)$ such that $\alpha_k \|\mathbf{g}_k\|^2 \leq \varepsilon$ and $\alpha_k \|\mathbf{g}_k\| \leq \delta/2$ for all $k \geq N(\delta)$.

The sequel of the proof is based on induction and is organized as follows. In the first part, we show that there exists a finite $k(\delta) \geq N(\delta)$ such that $\mathbf{x}_{k(\delta)} \in X^* + B^\delta$. In the second part, we establish that if \mathbf{x}_k belongs to $X^* + B^\delta$ for some $k \geq N(\delta)$ then so does \mathbf{x}_{k+1} , by showing that either $\text{dist}_{X^*}(\mathbf{x}_{k+1}) < \text{dist}_{X^*}(\mathbf{x}_k)$ holds, or $\mathbf{x}_k \in X^\varepsilon$ so that $\mathbf{x}_{k+1} \in X^* + B^\delta$ since the step taken is not longer than $\delta/2$.

Let $\mathbf{x}^* \in X^*$ be arbitrary. In every iteration k we then have

$$\|\mathbf{x}^* - \mathbf{x}_{k+1}\|^2 = \|\mathbf{x}^* - \text{Proj}_X(\mathbf{x}_k - \alpha_k \mathbf{g}_k)\|^2 \quad (6.45a)$$

$$\leq \|\mathbf{x}^* - \mathbf{x}_k + \alpha_k \mathbf{g}_k\|^2 \quad (6.45b)$$

$$= \|\mathbf{x}^* - \mathbf{x}_k\|^2 + \alpha_k \left(2\mathbf{g}_k^\top (\mathbf{x}^* - \mathbf{x}_k) + \alpha_k \|\mathbf{g}_k\|^2 \right), \quad (6.45c)$$

⁷A proper function is a function which is finite at least at some vector and nowhere attains the value $-\infty$. See also Section 1.4.

⁸For any set $S \subset \mathbb{R}^n$ the function χ_S is the indicator function of the set S , that is, $\chi_S(\mathbf{x}) = 0$ if $\mathbf{x} \in S$; and $\chi_S(\mathbf{x}) = +\infty$ if $\mathbf{x} \notin S$. See also Section 13.1.

Lagrangian duality

where the inequality follows from the projection property. Now, suppose

$$2\mathbf{g}_s^T(\mathbf{x}^* - \mathbf{x}_s) + \alpha_s \|\mathbf{g}_s\|^2 < -\varepsilon \quad (6.46)$$

for all $s \geq N(\delta)$. Then, using (6.45) repeatedly, we obtain that for any $k \geq N(\delta)$,

$$\|\mathbf{x}^* - \mathbf{x}_{k+1}\|^2 < \|\mathbf{x}^* - \mathbf{x}_{N(\delta)}\|^2 - \varepsilon \sum_{s=N(\delta)}^k \alpha_s,$$

and from (6.40) it follows that the right-hand side of this inequality tends to minus infinity as $k \rightarrow \infty$, which clearly is impossible. Therefore,

$$2\mathbf{g}_k^T(\mathbf{x}^* - \mathbf{x}_k) + \alpha_k \|\mathbf{g}_k\|^2 \geq -\varepsilon \quad (6.47)$$

for at least one $k \geq N(\delta)$, say $k = k(\delta)$. From the definition of $N(\delta)$, it follows that $\mathbf{g}_{k(\delta)}^T(\mathbf{x}^* - \mathbf{x}_{k(\delta)}) \geq -\varepsilon$. From the definition of a subgradient (cf. Definition 6.16) we have that $f(\mathbf{x}^*) - f(\mathbf{x}_{k(\delta)}) \geq \mathbf{g}_{k(\delta)}^T(\mathbf{x}^* - \mathbf{x}_{k(\delta)})$, since $\mathbf{x}^*, \mathbf{x}_{k(\delta)} \in X$. Hence, $f(\mathbf{x}_{k(\delta)}) \leq f^* + \varepsilon$, that is, $\mathbf{x}_{k(\delta)} \in X^\varepsilon \subseteq X^* + B^{\delta/2} \subset X^* + B^\delta$.

Now, suppose that $\mathbf{x}_k \in X^* + B^\delta$ for some $k \geq N(\delta)$. If (6.46) holds for $s = k$, then, by using (6.45), we have that $\|\mathbf{x}^* - \mathbf{x}_{k+1}\| < \|\mathbf{x}^* - \mathbf{x}_k\|$ for any $\mathbf{x}^* \in X^*$. Hence,

$$\begin{aligned} \text{dist}_{X^*}(\mathbf{x}_{k+1}) &\leq \|\text{Proj}_{X^*}(\mathbf{x}_k) - \mathbf{x}_{k+1}\| < \|\text{Proj}_{X^*}(\mathbf{x}_k) - \mathbf{x}_k\| \\ &= \text{dist}_{X^*}(\mathbf{x}_k) \leq \delta. \end{aligned}$$

Thus, $\mathbf{x}_{k+1} \in X^* + B^\delta$. Otherwise, (6.47) must hold and, using the same arguments as above, we obtain that $f(\mathbf{x}_k) \leq f^* + \varepsilon$, i.e., $\mathbf{x}_k \in X^\varepsilon \subseteq X^* + B^{\delta/2}$. As

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}_k\| &= \|\text{Proj}_X(\mathbf{x}_k - \alpha_k \mathbf{g}_k) - \mathbf{x}_k\| \leq \|\mathbf{x}_k - \alpha_k \mathbf{g}_k - \mathbf{x}_k\| \\ &= \alpha_k \|\mathbf{g}_k\| \leq \delta/2 \end{aligned}$$

whenever $k \geq N(\delta)$, it follows that $\mathbf{x}_{k+1} \in X^* + B^{\delta/2} + B^{\delta/2} = X^* + B^\delta$.

By induction with respect to $k \geq k(\delta)$, it follows that $\mathbf{x}_k \in X^* + B^\delta$ for all $k \geq k(\delta)$. Since this holds for arbitrarily small values of $\delta > 0$ and f is continuous, the theorem follows. \blacksquare

We next introduce the additional requirement (6.42); the resulting algorithm's convergence behaviour is now much more favourable, and the proof is at the same time less technical.

Theorem 6.24 (convergence of subgradient optimization methods, II) *Let $\{\mathbf{x}_k\}$ be generated by the method (6.40), (6.41), (6.42). If X^* is nonempty and the sequence $\{\mathbf{g}_k\}$ is bounded, then $f(\mathbf{x}_k) \rightarrow f^*$ and $\mathbf{x}_k \rightarrow \mathbf{x}^* \in X^*$ holds.*

Proof. Let $\mathbf{x}^* \in X^*$ and $k \geq 1$. Repeated application of (6.45) yields

$$\|\mathbf{x}^* - \mathbf{x}_k\|^2 \leq \|\mathbf{x}^* - \mathbf{x}_0\|^2 + 2 \sum_{s=0}^{k-1} \alpha_s \mathbf{g}_s^T (\mathbf{x}^* - \mathbf{x}_s) + \sum_{s=0}^{k-1} \alpha_s^2 \|\mathbf{g}_s\|^2. \quad (6.48)$$

Since $\mathbf{x}^* \in X^*$ and $\mathbf{g}_s \in \partial f(\mathbf{x}_s)$ for all $s \geq 0$ we obtain that

$$f(\mathbf{x}_s) \geq f^* \geq f(\mathbf{x}_s) + \mathbf{g}_s^T (\mathbf{x}^* - \mathbf{x}_s), \quad s \geq 0, \quad (6.49)$$

and hence that $\mathbf{g}_s^T (\mathbf{x}^* - \mathbf{x}_s) \leq 0$ for all $s \geq 0$. Define $c := \sup_k \{\|\mathbf{g}_k\|\}$ and $p = \sum_{k=0}^{\infty} \alpha_k^2$, so that $\|\mathbf{g}_s\| \leq c$ for any $s \geq 0$ and $\sum_{s=0}^{k-1} \alpha_s^2 < p$. From (6.48) we then conclude that $\|\mathbf{x}^* - \mathbf{x}_k\|^2 < \|\mathbf{x}^* - \mathbf{x}_0\|^2 + pc^2$ for any $k \geq 1$, and thus that the sequence $\{\mathbf{x}_k\}$ is bounded.

Assume now that there is no subsequence $\{\mathbf{x}_{k_i}\}$ of $\{\mathbf{x}_k\}$ with $\mathbf{g}_{k_i}^T (\mathbf{x}^* - \mathbf{x}_{k_i}) \rightarrow 0$. Then there must exist an $\varepsilon > 0$ with $\mathbf{g}_s^T (\mathbf{x}^* - \mathbf{x}_s) \leq -\varepsilon$ for all sufficiently large values of s . From (6.48) and the conditions on the step lengths it follows that $\|\mathbf{x}^* - \mathbf{x}_s\| \rightarrow -\infty$, which clearly is impossible. The sequence $\{\mathbf{x}_k\}$ must therefore contain a subsequence $\{\mathbf{x}_{k_i}\}$ such that $\mathbf{g}_{k_i}^T (\mathbf{x}^* - \mathbf{x}_{k_i}) \rightarrow 0$. From (6.49) it follows that $f(\mathbf{x}_{k_i}) \rightarrow f^*$. The boundedness of $\{\mathbf{x}_k\}$ implies the existence of a limit point of the subsequence $\{\mathbf{x}_{k_i}\}$, say \mathbf{x}^∞ . From the continuity of f it follows that $\mathbf{x}^\infty \in X^*$.

To show that \mathbf{x}^∞ is the only limit point of $\{\mathbf{x}_k\}$, let $\delta > 0$ and choose an $M(\delta)$ such that $\|\mathbf{x}^\infty - \mathbf{x}_{M(\delta)}\|^2 \leq \delta/2$ and $\sum_{s=M(\delta)}^{\infty} \alpha_s^2 \leq \delta/(2c^2)$. Consider any $k > M(\delta)$. Analogously to the derivation of (6.48), and using (6.49), we then obtain that

$$\|\mathbf{x}^\infty - \mathbf{x}_k\|^2 \leq \|\mathbf{x}^\infty - \mathbf{x}_{M(\delta)}\|^2 + \sum_{s=M(\delta)}^{k-1} \alpha_s^2 \|\mathbf{g}_s\|^2 < \frac{\delta}{2} + \frac{\delta}{2c^2} c^2 = \delta.$$

Since this holds for arbitrarily small values of $\delta > 0$, we are done. \blacksquare

Note that the boundedness condition on $\{\mathbf{g}_k\}$ is fulfilled whenever we know before-hand that the sequence $\{\mathbf{x}_k\}$ is bounded, such as in the case when X itself is bounded; cf. Proposition 6.17(a).

We finally present the even stronger convergence properties of the subgradient projection method using the Polyak step.

Theorem 6.25 (convergence of subgradient optimization methods, III)
 Let $\{\mathbf{x}_k\}$ be generated by the method (6.40), (6.43). If X^* is nonempty then $f(\mathbf{x}_k) \rightarrow f^*$ and $\mathbf{x}_k \rightarrow \mathbf{x}^* \in X^*$ holds.

Proof. From Proposition 6.22 follows that the sequence $\{\|\mathbf{x}_k - \mathbf{x}^*\|\}$ is strictly decreasing for every $\mathbf{x}^* \in X^*$, and therefore has a limit. By construction of the step length, in which the step lengths are bounded away from zero and $2[f(\mathbf{x}_k) - f^*]/\|\mathbf{g}_k\|^2$, it follows from the proof of Proposition 6.22 that $[f(\mathbf{x}_k) - f^*]^2/\|\mathbf{g}_k\|^2 \rightarrow 0$ must hold. Since $\{\mathbf{g}_k\}$ must be bounded due to the boundedness of $\{\mathbf{x}_k\}$ [Proposition 6.17(a)], we have that $f(\mathbf{x}_k) \rightarrow f^*$. Further, \mathbf{x}_k is bounded, and due to the continuity property of f every limit point must then belong to X^* .

It remains to show that there can be only one limit point. This property follows from the monotone decrease of the distance $\|\mathbf{x}_k - \mathbf{x}^*\|$. In detail, the proof is as follows. Suppose two subsequences of $\{\mathbf{x}_k\}$ exist, such that they converge to two different vectors in X^* :

$$\mathbf{x}_{m_i} \rightarrow \mathbf{x}_1^*; \quad \mathbf{x}_{l_i} \rightarrow \mathbf{x}_2^*; \quad \mathbf{x}_1^* \neq \mathbf{x}_2^*.$$

We must then have $\|\mathbf{x}_{l_i} - \mathbf{x}_1^*\| \rightarrow \rho > 0$. Since $\mathbf{x}_1^* \in X^*$ and the distance to X^* is decreasing, $\|\mathbf{x}_k - \mathbf{x}_1^*\| \rightarrow \rho$ holds, and in particular $\|\mathbf{x}_{m_i} - \mathbf{x}_1^*\| \rightarrow \rho$, which is a contradiction. ■

Contrary to the slow convergence of the subgradient projection algorithms that rely on the divergent series step length rule, under additional conditions on the function f a subgradient algorithm based on the Polyak step length (6.43) is *geometrically convergent*, in the sense that there exist $c > 0$ and $\eta \in (0, 1)$ with

$$\|\mathbf{x}_k - \mathbf{x}^*\| \leq c\eta^k, \quad k = 0, 1, \dots$$

See Section 6.8 for references to other subgradient algorithms than those presented here.

6.4.2 Application to the Lagrangian dual problem

We remind ourselves that the Lagrangian dual problem is a concave maximization problem, and that the appearance of the dual function is similar to that of the following example:

Let $h(x) := \text{minimum}\{h_1(x), h_2(x)\}$, where $h_1(x) := 4 - |x|$ and $h_2(x) := 4 - (x - 2)^2$. Then,

$$h(x) = \begin{cases} 4 - x, & \text{if } 1 \leq x \leq 4; \\ 4 - (x - 2)^2 & \text{if } x \leq 1, x \geq 4; \end{cases}$$

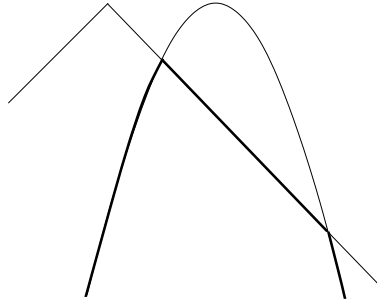


Figure 6.4: A convex min-function with three pieces.

cf. Figure 6.4.

The function h is non-differentiable at $x = 1$ and $x = 4$, since its graph has non-unique supporting hyperplanes there:

$$\partial h(x) = \begin{cases} \{4 - 2x\}, & \text{if } x < 1; \\ [-1, 2], & \text{if } x = 1; \\ \{-1\}, & \text{if } 1 < x < 4; \\ [-4, -1], & \text{if } x = 4; \\ \{4 - 2x\}, & \text{if } x > 4. \end{cases}$$

the subdifferential is here either a singleton (at differentiable points) or a closed interval (at non-differentiable points).

Note the monotonically decreasing nature of the relation $x \mapsto \partial h(x)$. Note also that $0 \in \partial h(1)$, whence $x^* = 1$ defines a maximum over \mathbb{R} .

Now, let $\mathbf{g} \in \partial q(\bar{\boldsymbol{\mu}})$, and let U^* be the set of optimal solutions to (6.10). Then,

$$U^* \subseteq \{ \boldsymbol{\mu} \in \mathbb{R}^m \mid \mathbf{g}^T(\boldsymbol{\mu} - \bar{\boldsymbol{\mu}}) \geq 0 \}.$$

In other words, any subgradient defines a half-space that contains the set of optimal solutions; cf. Figure 6.5. We therefore know that a small enough step in the direction of a subgradient gets us closer to the set of optimal solutions; cf. Proposition 6.22. But again consider Figure 6.5: an arbitrary subgradient, like the one depicted, may not define an ascent direction! As we saw in the previous section, convergence must be based on other arguments, like the decreasing distance to U^* alluded to above and in the previous section. In the next subsection we discuss in brief the generation of ascent directions.

We consider the Lagrangian dual problem (6.10). We suppose, as in the previous section, that X is compact so that the infimum in (6.9)

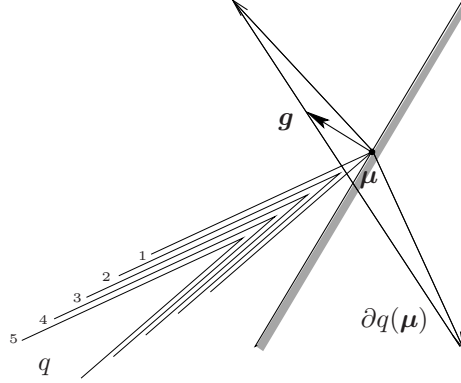


Figure 6.5: The half-space defined by the subgradient \mathbf{g} of q at $\boldsymbol{\mu}$. Note that the subgradient is not an ascent direction.

is attained for every $\boldsymbol{\mu} \geq \mathbf{0}^m$ (which is the set over which we wish to maximize q) and q is real-valued over \mathbb{R}_+^m .

In the case of our special concave maximization problem, the iteration has the form

$$\begin{aligned} \boldsymbol{\mu}_{k+1} &= \text{Proj}_{\mathbb{R}_+^m} [\boldsymbol{\mu}_k + \alpha_k \mathbf{g}_k] = [\boldsymbol{\mu}_k + \alpha_k \mathbf{g}_k]_+ \\ &= (\text{maximum} \{0, (\boldsymbol{\mu}_k)_i + \alpha_k (\mathbf{g}_k)_i\})_{i=1}^m, \end{aligned} \quad (6.50)$$

where $\mathbf{g}_k \in \partial q(\boldsymbol{\mu}_k)$ is arbitrarily chosen; we would typically use $\mathbf{g}_k = \mathbf{g}(\mathbf{x}_k)$, where $\mathbf{x}_k \in \arg \min_{\mathbf{x} \in X} L(\mathbf{x}, \boldsymbol{\mu}_k)$. The projection operation onto the first orthant is, as we can see, very simple.

Replacing the Polyak step (6.43) with the corresponding dual form

$$\alpha_k = \theta_k \frac{q^* - q(\boldsymbol{\mu}_k)}{\|\mathbf{g}_k\|^2}, \quad 0 < \sigma_1 \leq \theta_k \leq 2 - \sigma_2 < 2, \quad (6.51)$$

convergence will now be a simple consequence of the above theorems.

The compactness condition (6.37) and the fact that the feasible set of (6.4) is nonempty ensure that the problem (6.4) has an optimal solution; in particular, the feasibility condition (6.5) then holds. Further, if we introduce the Slater condition (6.16), we are ensured that there is no duality gap, and that the dual problem (6.10) has a compact set U^* of optimal solutions. Under these assumptions, we have the following results for subgradient optimization methods.

Theorem 6.26 (convergence of subgradient optimization methods) *Sup-*

pose that the problem (6.4) is feasible, and that the compactness condition (6.37) and the Slater condition (6.16) hold.

(a) Let $\{\mu_k\}$ be generated by the method (6.50), (6.41). Then, $q(\mu_k) \rightarrow q^*$, and $\text{dist}_{U^*}(\mu_k) \rightarrow 0$.

(b) Let $\{\mu_k\}$ be generated by the method (6.50), (6.41), (6.42). Then, $\{\mu_k\}$ converges to an optimal solution to (6.10).

(c) Let $\{\mu_k\}$ be generated by the method (6.50), (6.51). Then, $\{\mu_k\}$ converges to an optimal solution to (6.10).

Proof. The results follow from Theorems 6.23, 6.24, and 6.25, respectively. Note that in the first two cases, boundedness conditions were assumed for X^* and the sequence of subgradients. The corresponding conditions for the Lagrangian dual problem are fulfilled under the CQs imposed, since they imply that the search for an optimal solution is done over a compact set; cf. Theorem 6.9(a) and its proof. ■

6.4.3 The generation of ascent directions

Proposition 6.18 shows that the existence of a descent direction with respect to the convex function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at some $\bar{x} \in \mathbb{R}^n$ hinges on the existence of some vector $\bar{p} \in \mathbb{R}^n$ such that $f'(\bar{x}; \bar{p}) < 0$. According to the definition of the directional derivative and the compactness of $\partial f(\bar{x})$, this is equivalent to the statement that $\mathbf{g}^T \bar{p} \leq \varepsilon < 0$ for every $\mathbf{g} \in \partial f(\bar{x})$. In the context of Lagrangian duality we show below how we can generate an ascent direction for q at some $\mu \in \mathbb{R}^m$.

Definition 6.27 (steepest ascent direction) Suppose that the problem (6.4) is feasible, and that the compactness condition (6.37) holds. Consider the Lagrangian dual problem (6.10), and let $\mu \in \mathbb{R}^m$. A vector $\bar{p} \in \mathbb{R}^m$ with $\|\bar{p}\| \leq 1$ is a steepest ascent direction if

$$q'(\mu; \bar{p}) = \max_{\|p\| \leq 1} q'(\mu; p)$$

holds. ■

Proposition 6.28 (the shortest subgradient yields the steepest ascent direction) Suppose that the problem (6.4) is feasible, and that the compactness condition (6.37) holds. Consider the Lagrangian dual problem (6.10). The direction \bar{p} of steepest ascent with respect to q at μ is given below, where $\bar{g} \in \partial q(\mu)$ is the shortest subgradient in $\partial q(\mu)$ with respect to the Euclidean norm:

$$\bar{p} = \begin{cases} \mathbf{0}^m, & \text{if } \bar{g} = \mathbf{0}^m, \\ \frac{\bar{g}}{\|\bar{g}\|}, & \text{if } \bar{g} \neq \mathbf{0}^m. \end{cases}$$

Proof. By Definition 6.27 and Proposition 6.19(e), the following string of equalities and inequalities can easily be verified:

$$\begin{aligned}
 \max_{\|\mathbf{p}\| \leq 1} q'(\boldsymbol{\mu}; \mathbf{p}) &= \max_{\|\mathbf{p}\| \leq 1} \inf_{\mathbf{g} \in \partial q(\boldsymbol{\mu})} \mathbf{g}^T \mathbf{p} \\
 &\leq \inf_{\mathbf{g} \in \partial q(\boldsymbol{\mu})} \max_{\|\mathbf{p}\| \leq 1} \mathbf{g}^T \mathbf{p} \\
 &= \inf_{\mathbf{g} \in \partial q(\boldsymbol{\mu})} \|\mathbf{g}\| \\
 &= \|\bar{\mathbf{g}}\|.
 \end{aligned} \tag{6.52}$$

If we can construct a direction $\bar{\mathbf{p}}$ such that $q'(\boldsymbol{\mu}; \bar{\mathbf{p}}) = \|\bar{\mathbf{g}}\|$ then by (6.52) $\bar{\mathbf{p}}$ is the steepest ascent direction. If $\bar{\mathbf{g}} = \mathbf{0}^m$ then for $\bar{\mathbf{p}} = \mathbf{0}^m$ we obviously have that $q'(\boldsymbol{\mu}; \bar{\mathbf{p}}) = \|\bar{\mathbf{g}}\|$. Suppose then that $\bar{\mathbf{g}} \neq \mathbf{0}^m$, and let $\bar{\mathbf{p}} := \bar{\mathbf{g}}/\|\bar{\mathbf{g}}\|$. Note that

$$\begin{aligned}
 q'(\boldsymbol{\mu}; \bar{\mathbf{p}}) &= \inf_{\mathbf{g} \in \partial q(\boldsymbol{\mu})} \mathbf{g}^T \bar{\mathbf{p}} = \inf_{\mathbf{g} \in \partial q(\boldsymbol{\mu})} \frac{\bar{\mathbf{g}}^T \mathbf{g}}{\|\bar{\mathbf{g}}\|} \\
 &= \frac{1}{\|\bar{\mathbf{g}}\|} \inf_{\mathbf{g} \in \partial q(\boldsymbol{\mu})} \{\|\bar{\mathbf{g}}\|^2 + \bar{\mathbf{g}}^T(\mathbf{g} - \bar{\mathbf{g}})\} \\
 &= \|\bar{\mathbf{g}}\| + \frac{1}{\|\bar{\mathbf{g}}\|} \inf_{\mathbf{g} \in \partial q(\boldsymbol{\mu})} \bar{\mathbf{g}}^T(\mathbf{g} - \bar{\mathbf{g}}).
 \end{aligned} \tag{6.53}$$

Since $\bar{\mathbf{g}}$ is the shortest vector in $\partial q(\boldsymbol{\mu})$, then, by the variational inequality characterization of the projection of $\mathbf{0}^m$ onto $\partial q(\boldsymbol{\mu})$ established in Theorem 4.24, we obtain that $\bar{\mathbf{g}}^T(\mathbf{g} - \bar{\mathbf{g}}) \geq 0$ for every $\mathbf{g} \in \partial q(\boldsymbol{\mu})$. Hence, $\inf_{\mathbf{g} \in \partial q(\boldsymbol{\mu})} \bar{\mathbf{g}}^T(\mathbf{g} - \bar{\mathbf{g}}) = 0$ is achieved at $\bar{\mathbf{g}}$. From (6.53) it then follows that $q'(\boldsymbol{\mu}; \bar{\mathbf{p}}) = \|\bar{\mathbf{g}}\|$. We are done. ■

6.5 *Obtaining a primal solution

It remains for us to show how an optimal dual solution $\boldsymbol{\mu}^*$ can be *translated* into an optimal primal solution \mathbf{x}^* . Obviously, convexity and strong duality will be needed in general, if we are to be able to utilize the primal–dual optimality characterization in Theorem 6.7. It turns out that the generation of a primal optimum is automatic if q is differentiable at $\boldsymbol{\mu}^*$, which is also the condition under which the famous *Lagrange multiplier method* works. Unfortunately, in many cases, such as for most non-strictly convex optimization problems (like linear programming), this will not be the case, and then the translation work becomes more complex.

We start with the ideal case.

6.5.1 Differentiability at the optimal solution

The following results summarize the optimality conditions for the Lagrangian dual problem (6.10), and their consequences for the availability of a primal optimal solution in the absence of a duality gap.

Proposition 6.29 (optimality conditions for the dual problem) *Suppose that the compactness condition (6.37) holds in the problem (6.4). Suppose further that the vector μ^* solves the Lagrangian dual problem.*

(a) *The dual optimal solution is characterized by the inclusion*

$$\mathbf{0}^m \in -\partial q(\mu^*) + N_{\mathbb{R}_+^m}(\mu^*). \quad (6.54)$$

In other words, there then exists $\gamma^ \in \partial q(\mu^*)$ —an optimality-characterizing subgradient of q at μ^* —such that*

$$\mathbf{0}^m \leq \mu^* \perp \gamma^* \leq \mathbf{0}^m. \quad (6.55)$$

There exists a finite set of solutions $\mathbf{x}^i \in X(\mu^)$, $i = 1, \dots, k$, where $k \leq m + 1$ such that*

$$\gamma^* = \sum_{i=1}^k \alpha_i \mathbf{g}(\mathbf{x}^i); \quad \sum_{i=1}^k \alpha_i = 1; \quad \alpha_i \geq 0, \quad i = 1, \dots, k. \quad (6.56)$$

Hence, we have that

$$\sum_{i=1}^k \alpha_i \mu_i^* \mathbf{g}_i(\mathbf{x}^i) = \mathbf{0}. \quad (6.57)$$

(b) *If there is a duality gap, then q is non-differentiable at μ^* .*

(c) *If q is differentiable at μ^* , then there is no duality gap. Further, any vector in $X(\mu^*)$ then solves the primal problem (6.4).*

Proof. (a) The first result is a direct statement of the optimality conditions of the convex and subdifferentiable program (6.10); the complementarity conditions in (6.55) are an equivalent statement of the inclusion in (6.54).

The second result is an application of Carathéodory's Theorem 3.8 to the compact and convex set $\partial q(\mu^*)$.

(b) The result is established once (c) is.

(c) Let $\bar{\mathbf{x}}$ be any vector in $X(\mu^*)$ for which $\nabla q(\mu^*) = \mathbf{g}(\bar{\mathbf{x}})$ holds, cf. Proposition 6.20(a). We obtain from (6.55) that

$$\mathbf{0}^m \leq \mu^* \perp \mathbf{g}(\bar{\mathbf{x}}) \leq \mathbf{0}^m.$$

Hence, the pair $(\mu^*, \bar{\mathbf{x}})$ fulfills all the conditions stated in (6.12), so that, by Theorem 6.7, $\bar{\mathbf{x}}$ is an optimal solution to (6.4). ■

Remark 6.30 (the non-coordinability phenomenon and decomposition algorithms) Many interesting problems do not comply with the conditions in (c); for example, linear programming is one where the Lagrangian dual problem often is non-differentiable at every dual optimal solution.⁹ This is sometimes called the *non-coordinability phenomenon* (cf. [Las70, DiJ79]). It was in order to cope with this phenomenon that Dantzig–Wolfe decomposition ([DaW60, Las70]) and other column generation algorithms, Benders decomposition ([Ben62, Las70]) and generalized linear programming were developed; noticing that the convex combination of a finite number of candidate primal solutions are sufficient to verify an optimal primal–dual solution [cf. (6.57)], methodologies were developed to generate those vectors algorithmically. See also [LPS99] for overviews on the subject of generating primal optimal solutions from dual optimal ones, and [BSS93, Theorem 6.5.2] for an LP procedure that provides primal feasible solutions for convex programs.

Note that the equation (6.57) in (a) reduces to the complementarity condition that $\mu_i^* g_i(\bar{\mathbf{x}}) = 0$ holds, for the *averaged* solution, $\bar{\mathbf{x}} := \sum_{i=1}^k \alpha_i \mathbf{x}^i$, whenever all the functions g_i are affine. ■

6.5.2 Everett’s Theorem

The next result shows that the solution to the Lagrangian subproblem solves a perturbed version of the original problem. We state the result for the general problem to find

$$\begin{aligned} f^* &:= \infimum_{\mathbf{x}} f(\mathbf{x}), \\ \text{subject to } & \mathbf{x} \in X, \\ & g_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m, \\ & h_j(\mathbf{x}) = 0, \quad j = 1, \dots, \ell, \end{aligned} \tag{6.58}$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$, $i = 1, 2, \dots, m$, and $h_j : \mathbb{R}^n \rightarrow \mathbb{R}$, $j = 1, 2, \dots, \ell$, are given functions, and $X \subseteq \mathbb{R}^n$.

Theorem 6.31 (Everett’s Theorem) *Let $(\boldsymbol{\mu}, \boldsymbol{\lambda}) \in \mathbb{R}_+^m \times \mathbb{R}^\ell$. Consider the Lagrangian subproblem to*

$$\underset{\mathbf{x} \in X}{\text{minimize}} \left\{ f(\mathbf{x}) + \boldsymbol{\mu}^T \mathbf{g}(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{h}(\mathbf{x}) \right\}. \tag{6.59}$$

Suppose that $\bar{\mathbf{x}}$ is an optimal solution to this problem, and let $\mathcal{I}(\boldsymbol{\mu}) \subseteq \{1, \dots, m\}$ denote the set of indices i for which $\mu_i > 0$.

⁹In other words, even if a Lagrange multiplier vector is known, the Lagrangian subproblem may not identify a primal optimal solution.

(a) $\bar{\mathbf{x}}$ is an optimal solution to the perturbed primal problem to

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}), \quad (6.60)$$

subject to $\mathbf{x} \in X$,

$$\begin{aligned} g_i(\mathbf{x}) &\leq g_i(\bar{\mathbf{x}}), & i \in \mathcal{I}(\bar{\mathbf{x}}), \\ h_j(\mathbf{x}) &= h_j(\bar{\mathbf{x}}), & j = 1, \dots, \ell. \end{aligned}$$

(b) If $\bar{\mathbf{x}}$ is feasible in (6.58) and $\boldsymbol{\mu}^T \mathbf{g}(\bar{\mathbf{x}}) = 0$ holds, then $\bar{\mathbf{x}}$ solves (6.58), and the pair $(\boldsymbol{\mu}, \boldsymbol{\lambda})$ then solves the Lagrangian dual problem.

Proof. (a) Let \mathbf{x} satisfy the constraints of (6.60). Since we have that $\mathbf{h}(\mathbf{x}) = \mathbf{h}(\bar{\mathbf{x}})$ and $\boldsymbol{\mu}^T \mathbf{g}(\mathbf{x}) \leq \boldsymbol{\mu}^T \mathbf{g}(\bar{\mathbf{x}})$, the optimality of $\bar{\mathbf{x}}$ in (6.59) yields

$$\begin{aligned} f(\mathbf{x}) + \boldsymbol{\mu}^T \mathbf{g}(\bar{\mathbf{x}}) + \boldsymbol{\lambda}^T \mathbf{h}(\bar{\mathbf{x}}) &\geq f(\mathbf{x}) + \boldsymbol{\mu}^T \mathbf{g}(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{h}(\mathbf{x}) \\ &\geq f(\bar{\mathbf{x}}) + \boldsymbol{\mu}^T \mathbf{g}(\bar{\mathbf{x}}) + \boldsymbol{\lambda}^T \mathbf{h}(\bar{\mathbf{x}}), \end{aligned}$$

which shows that $f(\mathbf{x}) \geq f(\bar{\mathbf{x}})$. We are done.

(b) $\boldsymbol{\mu}^T \mathbf{g}(\bar{\mathbf{x}}) = 0$ implies that $g_i(\bar{\mathbf{x}}) = 0$ for $i \in \mathcal{I}(\boldsymbol{\mu})$; from (a) $\bar{\mathbf{x}}$ solves the problem to

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}), \quad (6.61)$$

subject to $\mathbf{x} \in X$,

$$\begin{aligned} g_i(\mathbf{x}) &\leq 0, & i \in \mathcal{I}(\bar{\mathbf{x}}), \\ h_j(\mathbf{x}) &= 0, & j = 1, \dots, \ell. \end{aligned}$$

In particular, then, since the feasible set of (6.58) is contained in that of (6.61) and $\bar{\mathbf{x}}$ is feasible in the former, $\bar{\mathbf{x}}$ must also solve (6.58). That the pair $(\boldsymbol{\mu}, \boldsymbol{\lambda})$ solves the dual problem follows by the equality between the primal and dual objective functions at $(\bar{\mathbf{x}}, \boldsymbol{\mu}, \boldsymbol{\lambda})$, and weak duality. ■

One important consequence of the result is that if the right-hand side perturbations $g_i(\bar{\mathbf{x}})$ and $h_i(\bar{\mathbf{x}})$ all are close to zero, the vector $\bar{\mathbf{x}}$ being near-feasible might mean that it is in fact acceptable as an approximate solution to the original problem. (This interpretation hinges on the dualized constraints being *soft* constraints, in the sense that a small violation is acceptable. See Section 1.8 for an introduction to the topic of soft constraints.)

6.6 *Sensitivity analysis

6.6.1 Analysis for convex problems

Consider the inequality constrained convex program (6.4), where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and g_i , $i = 1, \dots, m$, are convex functions and $X \subseteq \mathbb{R}^n$ is

a convex set. Suppose that the problem (6.4) is feasible, and that the compactness condition (6.37) and Slater condition (6.16) hold. This is the classic case where there exist multiplier vectors $\boldsymbol{\mu}^*$, according to Theorem 6.9, and strong duality holds.

For certain types of problems where the duality gap is zero and where there exist primal–dual optimal solutions, we have access to a beautiful theory of *sensitivity analysis*. The classic meaning of the term is the answer to the following question: what is the rate of change in f^* when a constraint right-hand side changes? This question answers important practical questions, like the following in manufacturing: If we buy one unit of additional resource at a given price, or if the demand of a product that we sell increases by a certain amount, then how much additional profit do we make?

We will here provide a basic result which states when this sensitivity analysis of the optimal objective value can be performed for the problem (6.4), and establish that the answer is determined precisely by the value of the Lagrange multiplier vector $\boldsymbol{\mu}^*$, provided that it is unique.

Definition 6.32 (perturbation function) *Consider the function $p : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{\pm\infty\}$ defined by*

$$\begin{aligned} p(\mathbf{u}) &:= \infimum_x f(\mathbf{x}), \\ \text{subject to } \quad &\mathbf{x} \in X, \\ &g_i(\mathbf{x}) \leq u_i, \quad i = 1, \dots, m, \quad \mathbf{u} \in \mathbb{R}^m; \end{aligned} \tag{6.62}$$

it is called the perturbation function, or primal function, associated with the problem (6.4). Its effective domain is the set $P := \{\mathbf{u} \in \mathbb{R}^m \mid p(\mathbf{u}) < +\infty\}$. ■

Under the above convexity conditions, we can establish that p is a convex function. Indeed, it holds that for any value of the Lagrange multiplier vector $\boldsymbol{\mu}^*$ for the problem (6.4) that

$$\begin{aligned} q(\boldsymbol{\mu}^*) &= \infimum_{\mathbf{x} \in X} \{f(\mathbf{x}) + (\boldsymbol{\mu}^*)^T \mathbf{g}(\mathbf{x})\} \\ &= \infimum_{\{(u, \mathbf{x}) \in P \times X \mid g(\mathbf{x}) \leq \mathbf{u}\}} \{f(\mathbf{x}) + (\boldsymbol{\mu}^*)^T \mathbf{g}(\mathbf{x})\} \\ &= \infimum_{\{(u, \mathbf{x}) \in P \times X \mid g(\mathbf{x}) \leq \mathbf{u}\}} \{f(\mathbf{x}) + (\boldsymbol{\mu}^*)^T \mathbf{u}\} \\ &= \infimum_{\mathbf{u} \in P} \infimum_{\{\mathbf{x} \in X \mid g(\mathbf{x}) \leq \mathbf{u}\}} \{f(\mathbf{x}) + (\boldsymbol{\mu}^*)^T \mathbf{u}\}. \end{aligned}$$

Since $\boldsymbol{\mu}^*$ is assumed to be a Lagrange multiplier vector, we have that $q(\boldsymbol{\mu}^*) = f^* = p(\mathbf{0}^m)$. By the definition of infimum, then, we have that

$$p(\mathbf{0}^m) \leq p(\mathbf{u}) + (\boldsymbol{\mu}^*)^T \mathbf{u}, \quad \mathbf{u} \in \mathbb{R}^m,$$

that is, $-\boldsymbol{\mu}^*$ (notice the sign!) is a subgradient of p at $\mathbf{u} = \mathbf{0}^m$ (see Definition 6.16). Moreover, by the result in Proposition 6.17(c), p is differentiable at $\mathbf{0}^m$ if and only if p is finite in a neighbourhood of $\mathbf{0}^m$ and $\boldsymbol{\mu}^*$ is a *unique* Lagrange multiplier vector, that is, the Lagrangian dual problem (6.10) has a unique optimal solution. We have therefore proved the following result:

Proposition 6.33 (a sensitivity analysis result) *Suppose that in the inequality constrained problem (6.4), $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$, $i = 1, \dots, m$, are convex functions and $X \subseteq \mathbb{R}^n$ is a convex set. Suppose that the problem (6.4) is feasible, and that the compactness assumption (6.37) and Slater condition (6.16) hold. Suppose further that the perturbed problem defined in (6.62) has an optimal solution in a neighbourhood of $\mathbf{u} = \mathbf{0}^m$, and that on the set of primal–dual optimal solutions to (6.4), (6.10), the dual optimal solution $\boldsymbol{\mu}^*$ is unique. Then, the perturbation function p is differentiable at $\mathbf{u} = \mathbf{0}^m$, and*

$$\nabla p(\mathbf{0}^m) = -\boldsymbol{\mu}^*$$

holds. ■

It is intuitive that the sign of $\nabla p(\mathbf{0}^m)$ should be non-positive; if a right-hand side of the (less-than) inequality constraints in (6.4) increases, then the feasible set becomes larger. [This means that we might be able to find feasible vectors \mathbf{x} in the new problem with $f(\mathbf{x}) < f^*$, where $f^* = p(\mathbf{0}^m)$ is the optimal value of the minimization problem (6.4).]

The result specializes immediately to linear programming problems, which is the problem type where this type of analysis is most often utilized. The proof of differentiability of the perturbation function at zero for that special case can however be done much more simply. (See Section 10.3.1.)

6.6.2 Analysis for differentiable problems

There exist local versions of the analysis valid also for non-convex problems, where we are interested in the effect of a problem perturbation on a KKT point. A special such analysis was recently performed by Bertsekas [Ber04], in which he shows that even when the problem is non-convex and the Lagrange multipliers are not unique, a sensitivity analysis is available as long as the functions defining the problem are differentiable. Suppose then that in the problem (6.4) the functions f and g_i , $i = 1, \dots, m$ are in C^1 and that X is nonempty. We generalize the concept of a *Lagrange multiplier vector* to here mean that it is a

Lagrangian duality

vector $\boldsymbol{\mu}^*$ associated with a *local* minimum \mathbf{x}^* such that

$$\left(\nabla f(\mathbf{x}^*) + \sum_{i=1}^m \mu_i^* \nabla g_i(\mathbf{x}^*) \right)^T \mathbf{p} \geq 0, \quad \mathbf{p} \in T_X(\mathbf{x}^*), \quad (6.63a)$$

$$\mu_i^* \geq 0, \quad i = 1, \dots, m, \quad (6.63b)$$

$$\mu_i^* = 0, \quad i \notin \mathcal{I}(\mathbf{x}^*), \quad (6.63c)$$

where $T_X(\mathbf{x}^*)$ is the tangent cone to X at \mathbf{x}^* (cf. Definition 5.2). Note that under an appropriate CQ this is equivalent to the KKT conditions, in which case we are simply requiring here that \mathbf{x}^* is a local minimum.

In the below result we utilize the notation

$$g_i^+(\mathbf{x}) := \text{maximum} \{0, g_i(\mathbf{x})\}, \quad i = 1, \dots, m,$$

and let $\mathbf{g}^+(\mathbf{x})$ be the m -vector of elements $g_i^+(\mathbf{x})$, $i = 1, \dots, m$.

Theorem 6.34 (sensitivity from the minimum norm multiplier vector) *Suppose that \mathbf{x}^* is a local minimum in the problem (6.4), and that the set of Lagrange multipliers is nonempty. Let $\boldsymbol{\mu}^*$ denote the Lagrange multiplier vector of minimum Euclidean norm. Then, for every sequence $\{\mathbf{x}_k\} \subset X$ of infeasible vectors such that $\mathbf{x}_k \rightarrow \mathbf{x}^*$ we have that*

$$f(\mathbf{x}^*) - f(\mathbf{x}_k) \leq \|\boldsymbol{\mu}^*\| \cdot \|\mathbf{g}^+(\mathbf{x}_k)\| + o(\|\mathbf{x}_k - \mathbf{x}^*\|). \quad (6.64)$$

Furthermore, if $\boldsymbol{\mu}^* \neq \mathbf{0}^m$ and $T_X(\mathbf{x}^*)$ is convex, the above inequality is sharp in the sense that there exists a sequence of infeasible vectors $\{\mathbf{x}_k\} \subset X$ such that

$$\lim_{k \rightarrow \infty} \frac{f(\mathbf{x}^*) - f(\mathbf{x}_k)}{\|\mathbf{g}^+(\mathbf{x}_k)\|} = \|\boldsymbol{\mu}^*\|,$$

and for this sequence

$$\lim_{k \rightarrow \infty} \frac{g_i^+(\mathbf{x}_k)}{\|\mathbf{g}^+(\mathbf{x}_k)\|} = \frac{\mu_i^*}{\|\boldsymbol{\mu}^*\|}, \quad i = 1, \dots, m,$$

holds. ■

Theorem 6.34 establishes the optimal rate of cost improvement with respect to infeasible constraint perturbations (in effect, those that imply an enlargement of the feasible set).

We finally remark that under stronger conditions still, the operator $\mathbf{u} \mapsto \mathbf{x}^*(\mathbf{u})$ assigning the (unique) optimal solution \mathbf{x}^* to each perturbation vector $\mathbf{u} \in \mathbb{R}^m$ is differentiable at $\mathbf{u} = \mathbf{0}^m$. Such a result is

reminiscent to the Implicit Function Theorem, which however only covers equality systems. If we are to study the sensitivity of \mathbf{x}^* to changes in the right-hand sides of inequality constraints as well, then the analysis becomes complicated due to the fact that we must be able to predict if some active constraints may become inactive in the process. In some circumstances, different directions of change in the right-hand sides may cause different subsets of the active constraints $\mathcal{I}(\mathbf{x}^*)$ at \mathbf{x}^* to become inactive, and this would most probably then be a non-differentiable point of the operator $\mathbf{u} \mapsto \mathbf{x}^*(\mathbf{u})$. A sufficient condition (but not necessary, at least in the case of linear constraints) for this to not happen is when \mathbf{x}^* is *strictly complementary*, that is, when there exists a multiplier vector $\boldsymbol{\mu}^*$ with $\mu_i^* > 0$ for every $i \in \mathcal{I}(\mathbf{x}^*)$.

6.7 Applications

We provide two example applications of Lagrangian duality. The first describes the primal–dual relationship between currents and voltages in an electrical network of devices (voltage sources, diodes, and resistors); this application illustrates that Lagrange multipliers often have direct interpretations. The second application concerns a classic combinatorial optimization problem: the traveling salesman problem. We show how to approximately solve this problem through Lagrangian relaxation and subgradient optimization.

6.7.1 Electrical networks

An *electrical network* (or, *circuit*) is an interconnection of analog electrical elements such as resistors, inductors, capacitors, diodes, and transistors. Its size varies from the smallest integrated circuit to an entire electricity distribution network. A circuit is a network that has at least one closed loop. A network is a connection of 2 or more simple circuit elements, and may not be a circuit. The goal when designing electrical networks for signal processing is to apply a predefined operation on *potential differences* (measured in *volts*) or *currents* (measured in *amperes*). Typical functions for these electrical networks are amplification, oscillation and analog linear algorithmic operations such as addition, subtraction, multiplication, and division. In the case of power distribution networks, engineers design the circuit to transport energy as efficiently as possible while at the same time taking into account economic factors, network safety and redundancy. These networks use components such as power lines, cables, circuit breakers, switches and transformers.

To design any electrical circuits, electrical engineers need to be able

to predict the voltages and currents in the circuit. Linear circuits (that is, an electrical network where all elements have a linear current–voltage relation) can be quite easily analyzed through the use of complex numbers and systems of linear equations,¹⁰ while nonlinear elements require a more sophisticated analysis. The classic electrical laws describing the equilibrium state of an electrical network are due to G. Kirchhoff [Kir1847]; referred to as *Kirchhoff’s circuit laws* they express in a mathematical form the conservation of charge and energy.¹¹

Formally, we let an electrical circuit be described by branches (or, links) connecting nodes. We present a simple example where the only devices are voltage sources, resistors, and diodes. The resulting equilibrium conditions will be shown to be represented as the solution to a strictly convex quadratic program. In general, devices such as resistors can be non-linear, but linearity is assumed throughout this section.

- A *voltage source* maintains a constant branch voltage v_s irrespective of the branch current c_s . The power absorbed by the device is $-v_s c_s$.
- A *diode* permits the branch current c_d to flow in one direction only, but consumes no power regardless of the current or voltage on the branch. Denoting the branch voltage by v_d , the direction condition can be stated as a complementarity condition:

$$c_d \geq 0; \quad v_d \geq 0; \quad v_d c_d = 0. \quad (6.65)$$

- A *resistor* consumes power in relation with its resistance, denoted by R_r . We recognize the following law describing the relationship between the branch current and voltage in a linear resistor:

$$v_r = -R_r c_r. \quad (6.66)$$

The power consumed is given by

$$-v_r c_r = \frac{v_r^2}{R_r} = R_r c_r^2, \quad (6.67)$$

where we have utilized (6.66) to derive two alternative relations.

We must be careful about the direction of flow of currents and voltages, and thus define, for each type of device, a node–branch incidence

¹⁰For such networks already Maxwell [Max1865] had stated equilibrium conditions.

¹¹These laws can be derived from Maxwell’s equations, but Kirchhoff preceded Maxwell and derived his equations from work done by G. Ohm.

matrix of the form

$$n_{ij} := \begin{cases} -1, & \text{if branch } j \text{ has node } i \text{ as its origin,} \\ 1, & \text{if branch } j \text{ ends in node } i, \\ 0, & \text{otherwise.} \end{cases}$$

The interpretation of a current flow variable is that the direction is from the negative to the positive terminal of the device, that is, from the origin to the ending node of the branch; a negative variable value will therefore correspond to a flow in the opposite direction. Note that for the diodes, the latter is not allowed, as seen in (6.65).

For the three types of devices we hence yield incidence matrices denoted by \mathbf{N}_S , \mathbf{N}_R , and \mathbf{N}_D , creating a partitioned matrix $\mathbf{N} = [\mathbf{N}_S \mathbf{N}_D \mathbf{N}_R]$. Similarly, we let $\mathbf{c} = (\mathbf{c}_S^T, \mathbf{c}_D^T, \mathbf{c}_R^T)^T$ and $\mathbf{v} = (\mathbf{v}_S^T, \mathbf{v}_D^T, \mathbf{v}_R^T)^T$ represent the vectors of branch currents and voltages. We also let $\mathbf{p} = (\mathbf{p}_S^T, \mathbf{p}_D^T, \mathbf{p}_R^T)^T$ denote the vector of node potentials. Before stating the optimization problem whose minimum describes the equilibrium of the system, we recall the two fundamental equilibrium laws:

Kirchhoff's current law: The sum of all currents entering a node is equal to the sum of all currents leaving the node. In other words, $\mathbf{N}\mathbf{c} = \mathbf{0}$.¹²

$$\mathbf{N}_S \mathbf{c}_S + \mathbf{N}_D \mathbf{c}_D + \mathbf{N}_R \mathbf{c}_R = \mathbf{0}. \quad (6.68)$$

Kirchhoff's voltage law: The difference between the node potentials at the ends of each branch is equal to the branch voltage. In other words, $\mathbf{N}^T \mathbf{p} = \mathbf{v}$.¹³

$$\mathbf{N}_S^T \mathbf{p} = \mathbf{v}_S, \quad (6.69a)$$

$$\mathbf{N}_D^T \mathbf{p} = \mathbf{v}_D, \quad (6.69b)$$

$$\mathbf{N}_R^T \mathbf{p} = \mathbf{v}_R. \quad (6.69c)$$

We summarize the equations representing the characteristics of the electrical devices as follows: For the diodes, (6.65) yields

$$\mathbf{v}_D \geq \mathbf{0}; \quad \mathbf{c}_D \geq \mathbf{0}; \quad \mathbf{v}_D^T \mathbf{c}_D = 0. \quad (6.70)$$

For the resistors, (6.66) yields

$$\mathbf{v}_R = -\mathbf{R}\mathbf{c}_R, \quad (6.71)$$

¹²This law is also referred to as the first law, the point rule, the junction rule, and the node law.

¹³This law is a corollary to Ohm's law, and is also referred to as the loop law.

Lagrangian duality

\mathbf{R} being the diagonal matrix with elements equal to the values R_r .

Hence, (6.68)–(6.71) represent the equilibrium conditions of the circuit. We will now describe the optimization problem whose optimality conditions are, precisely, (6.68)–(6.71) [note that \mathbf{v}_S is fixed]:

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} \mathbf{c}_R^T \mathbf{R} \mathbf{c}_R - \mathbf{v}_S^T \mathbf{c}_S, \\ & \text{subject to} \quad \mathbf{N}_S \mathbf{c}_S + \mathbf{N}_D \mathbf{c}_D + \mathbf{N}_R \mathbf{c}_R = \mathbf{0}, \\ & \quad \quad \quad -\mathbf{c}_D \leq \mathbf{0}. \end{aligned} \tag{6.72}$$

In the problem (6.72) we wish to determine branch currents \mathbf{c}_S , \mathbf{c}_D , and \mathbf{c}_R so as to minimize the sum of half the energy absorbed in the resistors and the energy loss of the voltage source. Note the sign condition on the diode currents.

Note that this is a convex program with linear constraints, and thus the KKT conditions are both necessary and sufficient for the global optimality of the currents. It is instrumental to check that the KKT conditions for (6.72) are given by (6.68)–(6.71), where the Lagrange multipliers are given by $(\mathbf{p}^T, \mathbf{v}_D^T)^T$.

In the discussion terminating in the Strong Duality Theorem 6.13, we showed that the Lagrangian dual of a strictly convex quadratic optimization problem is yet another convex quadratic optimization problem. In our case, following that development, we can derive the following dual optimization problem in terms of the node potentials \mathbf{p} (notice, again, that \mathbf{v}_S is fixed):

$$\begin{aligned} & \text{maximize} \quad -\frac{1}{2} \mathbf{v}_R^T \mathbf{R}^{-1} \mathbf{v}_R, \\ & \text{subject to} \quad \mathbf{N}_S^T \mathbf{p} = \mathbf{v}_S, \\ & \quad \quad \quad \mathbf{N}_D^T \mathbf{p} - \mathbf{v}_D = \mathbf{0}, \\ & \quad \quad \quad \mathbf{N}_R^T \mathbf{p} - \mathbf{v}_R = \mathbf{0}, \\ & \quad \quad \quad \mathbf{v}_D \geq \mathbf{0}. \end{aligned} \tag{6.73}$$

In the dual problem (6.73) the matrix \mathbf{R}^{-1} is the diagonal matrix of conductances. The objective function is equivalent to the minimization of the power absorbed by the resistors, and we wish to determine the branch voltages \mathbf{v}_D and \mathbf{v}_R , and the potential vector \mathbf{p} .

Verify that the KKT conditions for this problem, again, reduce to the equilibrium conditions (6.68)–(6.71). In other words, the Lagrange multipliers for the dual problem (6.73) are the (primal) branch currents.

Finally, let us note that by Theorem 6.13(a) the two problems (6.72) and (6.73) have the same objective value at optimality. That is,

$$\frac{1}{2}\mathbf{c}_R^T \mathbf{R} \mathbf{c}_R + \frac{1}{2}\mathbf{v}_R^T \mathbf{R}^{-1} \mathbf{v}_R - \mathbf{v}_S^T \mathbf{c}_S = 0.$$

By (6.70)–(6.71), the above equation reduces to

$$\mathbf{v}_S^T \mathbf{c}_S + \mathbf{v}_D^T \mathbf{c}_D + \mathbf{v}_R^T \mathbf{c}_R = 0,$$

which is precisely the principle of energy conservation.

6.7.2 A Lagrangian relaxation of the traveling salesman problem

Lagrangian relaxation has shown to be remarkably efficient for some combinatorial optimization problems. This is surprising when taking into account that such problems are integer or mixed-integer problems, which suffer from non-zero duality gaps in general. What then lies behind their popularity?

- One can show that Lagrangian relaxation of an integer program is always at least as good as that of a *continuous relaxation*¹⁴ (in the sense that the value of f_R is higher for Lagrangian relaxation than for a continuous relaxation);
- Together with heuristics for finding primal feasible solution, good feasible solutions are often found;
- The Lagrangian relaxed problems can be made computationally much simpler than the original problem, while still keeping a lot of the structure of the original problem.

6.7.2.1 The traveling salesman problem

Let the graph $\mathcal{G} = (\mathcal{N}, \mathcal{L})$ be defined by a number of cities (or, nodes) $i \in \mathcal{N}$ and undirected links in between subsets of pairs of them: $(i, j) \in \mathcal{L} \subseteq \mathcal{N} \times \mathcal{N}$. Notice that the links (i, j) and (j, i) are identical, and are in \mathcal{L} represented by one non-directed link only.

Let c_{ij} denote the distance between the cities i and j , $\{i, j\} \subset \mathcal{N}$. We introduce the following binary variables:

$$x_{ij} := \begin{cases} 1, & \text{if link } (i, j) \text{ is part of the TSP tour,} \\ 0, & \text{otherwise,} \end{cases} \quad (i, j) \in \mathcal{L}.$$

¹⁴The continuous relaxation amounts to removing the integrality conditions, replacing, for example, $x_j \in \{0, 1\}$ by $x_j \in [0, 1]$.

Lagrangian duality

With these definitions, the undirected traveling salesman problem (TSP) is to

$$\text{minimize}_x \quad \sum_{(i,j) \in \mathcal{L}} c_{ij} x_{ij}, \quad (6.74a)$$

$$\text{subject to} \quad \sum_{(i,j) \in \mathcal{L}: \{i,j\} \subseteq \mathcal{S}} x_{ij} \leq |\mathcal{S}| - 1, \quad \mathcal{S} \subset \mathcal{N}, \quad (6.74b)$$

$$\sum_{(i,j) \in \mathcal{L}} x_{ij} = n, \quad (6.74c)$$

$$\sum_{i \in \mathcal{N}: (i,j) \in \mathcal{L}} x_{ij} = 2, \quad j \in \mathcal{N}, \quad (6.74d)$$

$$x_{ij} \in \{0, 1\}, \quad (i, j) \in \mathcal{L}. \quad (6.74e)$$

The constraints have the following interpretation: (6.74b) implies that there can be no *sub-tours*, that is, a tour where fewer than n cities are visited (if $\mathcal{S} \subset \mathcal{N}$ then there can be at most $|\mathcal{S}| - 1$ links between nodes in the set \mathcal{S} , where $|\mathcal{S}|$ is the cardinality—number of members—of the set \mathcal{S}); (6.74c) implies that in total n cities must be visited; and (6.74d) implies that each city is connected to two others, such that we make sure to arrive from one city and leave for the next.

This problem is NP-hard, which implies that there is no known polynomial algorithm for solving it. We resort therefore to the use of relaxation techniques, in particular Lagrangian relaxation. We have more than one alternative relaxation to perform: If we Lagrangian relax the tree constraints (6.74b) and (6.74c) the remaining problem is a *2-matching* problem; it can be solved in polynomial time. If we instead Lagrangian relax the degree constraints (6.74d) for every node except for one node the remaining problem is a *1-MST* problem, that is, a special type of minimum spanning tree (MST) problem.

The following definition is classic: a *Hamiltonian path* (respectively, *cycle*) is a path (respectively, cycle) which passes every node in the graph exactly once. Every Hamiltonian cycle is a Hamiltonian path from a node s to another node, t , followed by a link (t, s) ; a subgraph which consists of a spanning tree plus an extra link such that all nodes have degree two. This is then a feasible solution to the TSP.

A *1-MST* problem is the problem to find an MST in the graph that excludes node s , followed by the addition of the two least expensive links from node s to that tree. If all nodes happen to get degree two, then the 1-MST solution is a traveling salesman tour (that is, a Hamiltonian cycle). The idea behind solving the Lagrangian dual problem is then to find proper multiplier values such that the Lagrangian relaxation will produce feasible solutions.

6.7.2.2 Lagrangian relaxation of the traveling salesman problem

Suppose that we Lagrangian relax the degree constraints (6.74d), except for that for node 1. We assume that the starting node for the trip, node $s \in \mathcal{N}$, and all the links in \mathcal{L} connected to it, have been removed temporarily (in the 1-MST, this data is re-introduced later), but without changing the notation to reflect this.

The subproblem is the following: a 1-MST defined by

$$\begin{aligned} q(\boldsymbol{\lambda}) &= \underset{\mathbf{x}}{\text{minimum}} \sum_{(i,j) \in \mathcal{L}} c_{ij} x_{ij} + \sum_{j \in \mathcal{N}} \lambda_j \left(2 - \sum_{i \in \mathcal{N}: (i,j) \in \mathcal{L}} x_{ij} \right) \\ &= 2 \sum_{j \in \mathcal{N}} \lambda_j + \underset{\mathbf{x}}{\text{minimum}} \sum_{(i,j) \in \mathcal{L}} (c_{ij} - \lambda_i - \lambda_j) x_{ij}. \end{aligned}$$

We see immediately the role of the Lagrange multipliers: a high (low) value of the multiplier λ_j makes node j attractive (unattractive) in the above 1-MST problem, and will therefore lead to more (less) links being attached to it.

When solving the Lagrangian dual problem, we will use the class of subgradient optimization methods, an overview of which is found in Section 6.4.

What is the updating step in the subgradient method, and what is its interpretation? It is as usual an update in the direction of a subgradient, that is, the direction of

$$h_i(\mathbf{x}(\boldsymbol{\lambda})) := 2 - \sum_{i \in \mathcal{N}: (i,j) \in \mathcal{L}} x_{ij}(\boldsymbol{\lambda}), \quad i \in \mathcal{N},$$

where the value of $x_{ij} \in \{0,1\}$ is the solution to the 1-MST solution with link costs $c_{ij} - \lambda_i - \lambda_j$. We see from the direction formula that

$$\lambda_j^{\text{new}} := \lambda_j + \alpha \left(2 - \sum_{i \in \mathcal{N}: (i,j) \in \mathcal{L}} x_{ij}(\boldsymbol{\lambda}) \right), \quad j \in \mathcal{N},$$

where $\alpha > 0$ is a step length. It is interesting to investigate what the update means:

$$\text{current degree at node } j : \begin{cases} > 2 \implies \lambda_j \downarrow (\text{link cost } \uparrow) \\ = 2 \implies \lambda_j - (\text{link cost constant}) \\ < 2 \implies \lambda_j \uparrow (\text{link cost } \downarrow) \end{cases}$$

In other words, the updating formula in a subgradient method is such that the link cost in the 1-MST subproblem is shifted upwards

(downwards) if there are too many (too few) links connected to node j in the 1-MST. We are hence adjusting the *node prices* of the nodes in such a way as to try to influence the 1-MST problem to always choose 2 links per node to connect to.

6.7.2.3 A feasibility heuristic

A feasibility heuristic takes the optimal solution from the Lagrangian minimization problem over \mathbf{x} and adjusts it such that a feasible solution to the original problem is constructed. Since one cannot predict if, or when, a primal feasible solution will be found directly from the subproblem, the heuristic will provide a solution that can be used in place of an optimal one, if one is not found. Moreover, as we know from Lagrangian duality theory, we then have access to both lower and upper bounds on the optimal value f^* of the original problem, and so we have a quality measure of the feasible solutions found.

A feasibility heuristic which can be used together with our Lagrangian relaxation is as follows.

Identify a path in the 1-MST with many links. Then form a subgraph with the remaining nodes and find a path that passes all of them. Put the two paths together in the best way. The resulting path is a Hamiltonian cycle, that is, a feasible solution.

6.7.2.4 The Philips example

In 1987–1988 an M.Sc. project was performed at the department of mathematics at Linköping University, in cooperation with the company Philips, Norrköping. The project was initiated with the goal to improve the current practice of solving a production planning problem.

The problem was as follows: Philips produce circuit boards, perhaps several hundreds or thousands of the same type. There is a new batch of patterns (holes) to be drilled every day, and perhaps even several such batches per day.

To speed up the production process the drilling machine is connected to a microcomputer that selects the ordering of the holes to be drilled, given their coordinates. The algorithm for performing the sorting used to be a simple sorting operation that found, for every fixed x -coordinate, the corresponding y -coordinates and sorted them in increasing order. The movement of the drill was therefore from left to right, and for each fixed x -coordinate the movement was vertical. The time it took to drill the holes on one circuit board was, however, far too long, simply because the drill traveled around a lot without performing any tasks, following a path that was too long. (On the other hand, the actual ordering was

very fast to produce!) All in all, the complete batch production took too long because of the poorly planned drill movement.

It was observed that the production planning problem is a traveling salesman problem, where the cities are the holes to be drilled, and the distances between them correspond to the Euclidean distances between them. Therefore, an efficient TSP heuristic was devised and implemented, for use in conjunction with the microcomputer. In fact, it was based on precisely the above Lagrangian relaxation, a subgradient optimization method, and a graph-search type heuristic of the form discussed above.

A typical run with the algorithm took a few minutes, and was always stopped after a fixed number of subgradient iterations; the generation of feasible solutions with the above graph search technique was performed at every K^{th} iteration, where $K > 1$ is an integer. (Moreover, feasible solutions were not generated during the first iterations of the dual procedure, because of the poor quality of λ_k for low values of k ; often the traveling salesman tour resulting from the heuristic is better when the multipliers are near-optimal in the Lagrangian dual problem.)

In one of the examples implemented it was found that the optimal path length was in the order to 2 meters, and that the upper and lower bounds on f^* produced lead to the conclusion that the relative error of the path length of the best feasible solution found was *less than 7 %*, a quite good result, also showing that the duality gap for the problem at hand (together with the Lagrangian relaxation chosen) is quite small.

After implementing the new procedure, Philips could report an increase in production by some 70 %. Hence, the slightly longer time it took to provide a better production plan, that is, the traveling salesman tour for the drill to follow, was more than well compensated by the fact that the drilling could be done much faster.

This is a case where Lagrangian relaxation helped to solve a large-scale, complex and difficult problem by utilizing problem structure.

6.8 Notes and further reading

Lagrangian duality has been developed in many sources, including early developments by Arrow, Hurwicz, and Uzawa [AHU58], Everett [Eve63], and Falk [Fal67], and later on by Rockafellar [Roc70]. Our development follows to a large extent that of portions of the text books by Bertsekas [Ber99], Bazaraa et al. [BSS93], and Rockafellar [Roc70].

The Relaxation Theorem 6.1 can almost be considered to be folklore, and can be found in a slightly different form in [Wol98, Proposition 2.3].

The differentiability properties of convex functions were developed largely by Rockafellar [Roc70], whose text we mostly follow.

Subgradient methods were developed in the Soviet Union in the 1960s, predominantly by Ermol'ev, Polyak, and Shor. Text book treatments of subgradient methods are found, for example, in [Sho85, HiL93, Ber99]. Theorem 6.23 is essentially due to Ermol'ev [Erm66]; the proof stems from [LPS96]. Theorem 6.24 is due to Shepilov [She76]; finally, Theorem 6.25 is due to Polyak [Pol69].

Everett's Theorem 6.31 is due to Everett [Eve63].

Theorem 6.34 stems from [Ber04, Proposition 1.1].

That the equilibrium conditions of an electrical or hydraulic network are attained as the minimum of the total energy loss were known more than a century ago. Mathematical programming models for the electrical network equilibrium problems described in Section 6.7.1 date at least as far back as to Duffin [Duf46, Duf47] and d'Auriac [dAu47]. Duffin constructs his objective function as a sum of integrals of resistance functions. The possibility of viewing the equilibrium problem in at least two related, dual, ways as that of either finding the optimal flows of currents or the optimal potentials was also known early in the analysis of electrical networks; these two principles are written out in [Cro36] in work on pipe networks, and explicitly stated as a pair of primal–dual quadratic programming problems in [Den59]; we followed his development, as represented in [BSS93, Section 1.2.D].

The traveling salesman problem is an essential model problem in combinatorial optimization. Excellent introductions to the field can be found in [Law76, PaS82, NeW88, Wol98, Sch03]. It was the work in [HWC74, Geo74, Fis81, Fis85], among others, in the 1970s and 1980s on the traveling salesman problem and its relatives that made Lagrangian relaxation and subgradient optimization popular, and it remains most popular within the combinatorial optimization field.

6.9 Exercises

Exercise 6.1 (numerical example of Lagrangian relaxation) Consider the convex problem to

$$\begin{aligned} &\text{minimize} \quad \frac{1}{x_1} + \frac{4}{x_2}, \\ &\text{subject to} \quad x_1 + x_2 \leq 4, \\ &\quad \quad \quad x_1, x_2 \geq 0. \end{aligned}$$

(a) Lagrangian relax the first constraint, and write down the resulting implicit dual objective function and the dual problem. Motivate why the

relaxed problem always has a unique optimum, whence the dual objective function is everywhere differentiable.

(b) Solve the implicit Lagrangian dual problem by utilizing that the gradient to a differentiable dual objective function can be expressed by using the functions that are involved in the relaxed constraints and the unique solution to the relaxed problem.

(c) Give an explicit dual problem (a dual problem only in terms of the Lagrange multipliers). Solve it to confirm the results in (b).

(d) Find the original problem's optimal solution.

(e) Show that strong duality holds.

Exercise 6.2 (global optimality conditions) Consider the problem to

$$\begin{aligned} & \text{minimize } f(\mathbf{x}) := x_1 + 2x_2^2 + 3x_3^3, \\ & \text{subject to } x_1 + 2x_2 + x_3 \leq 3, \\ & \quad 2x_1^2 + x_2 \geq 2, \\ & \quad 2x_1 + x_3 = 2, \\ & \quad x_j \geq 0, \quad j = 1, 2, 3. \end{aligned}$$

(a) Formulate the Lagrangian dual problem that results from Lagrangian relaxing all but the sign constraints.

(b) State the global primal–dual optimality conditions.

Exercise 6.3 (Lagrangian relaxation) Consider the problem to

$$\begin{aligned} & \text{minimize } f(\mathbf{x}) := x_1^2 + 2x_2^2, \\ & \text{subject to } x_1 + x_2 \geq 2, \\ & \quad x_1^2 + x_2^2 \leq 5. \end{aligned}$$

Find an optimal solution through Lagrangian duality.

Exercise 6.4 (Lagrangian relaxation) In many circumstances it is of interest to calculate the Euclidean projection of a vector onto a subspace. Especially, consider the problem to find the Euclidean projection of the vector $\mathbf{y} \in \mathbb{R}^n$ onto the null space of the matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$, that is, to find an $\mathbf{x} \in \mathbb{R}^n$ that solves the problem to

$$\begin{aligned} & \text{minimize } f(\mathbf{x}) := \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|^2, \\ & \text{subject to } \mathbf{A}\mathbf{x} = \mathbf{0}^m, \end{aligned}$$

where \mathbf{A} is such that $\text{rank } \mathbf{A} = m$.

The solution to this problem is classic: the projection is given by

$$\mathbf{x}^* = \mathbf{y} - \mathbf{A}^T(\mathbf{A}\mathbf{A}^T)^{-1}\mathbf{A}\mathbf{y}.$$

If we let $\mathbf{P} := \mathbf{I}^n - \mathbf{A}^T(\mathbf{A}\mathbf{A}^T)^{-1}\mathbf{A}$, where $\mathbf{I}^n \in \mathbb{R}^{n \times n}$ is the unit matrix, be the *projection matrix*, the formula is simply $\mathbf{x}^* = \mathbf{P}\mathbf{y}$.

Lagrangian duality

Derive this formula by utilizing Lagrangian duality. Motivate every step by showing that the necessary properties are fulfilled.

[Note: This exercise is similar to that in Example 5.51, but utilizes Lagrangian duality rather than the KKT conditions to derive the projection formula.]

Exercise 6.5 (Lagrangian relaxation, exam 040823) Consider the following linear optimization problem:

$$\begin{aligned} & \text{minimize} && f(x, y) := x - 0.5y, \\ & \text{subject to} && -x + y \leq -1, \\ & && -2x + y \leq -2, \\ & && (x, y)^T \in \mathbb{R}_+^2. \end{aligned}$$

(a) Show that the problem satisfies Slater's constraint qualification. Derive the Lagrangian dual problem corresponding to the Lagrangian relaxation of the two linear inequality constraints, and show that its set of optimal solutions is convex and bounded.

(b) Calculate the set of subgradients of the Lagrangian dual function at the dual points $(1/4, 1/3)^T$ and $(1, 0)^T$.

Exercise 6.6 (Lagrangian relaxation) Provide an explicit form of the Lagrangian dual problem for the problem to

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^m \sum_{j=1}^n x_{ij} \ln x_{ij} \\ & \text{subject to} && \sum_{i=1}^m x_{ij} = b_j, && j = 1, \dots, n, \\ & && \sum_{j=1}^n x_{ij} = a_i, && i = 1, \dots, m, \\ & && x_{ij} \geq 0, && i = 1, \dots, m, \quad j = 1, \dots, n, \end{aligned}$$

where $a_i > 0$, $b_j > 0$ for all i, j , and where the linear equalities are Lagrangian relaxed.

Exercise 6.7 (Lagrangian relaxation) Given is the problem to

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) := 2x_1^2 + x_2^2 + x_1 - 3x_2, \quad (6.75a)$$

$$\text{subject to} \quad x_1^2 + x_2 \geq 8, \quad (6.75b)$$

$$x_1 \in [1, 3], \quad (6.75c)$$

$$x_2 \in [2, 5]. \quad (6.75d)$$

Lagrangian relax the constraint (6.75b) with a multiplier μ . Formulate the Lagrangian dual problem and calculate the dual function's value at $\mu = 1$, $\mu = 2$, and $\mu = 3$. Within which interval lies the optimal value f^* ? Also, draw the dual function.

Exercise 6.8 (Lagrangian duality for integer problems) Consider the primal problem to

$$\begin{aligned} & \text{minimize } f(\mathbf{x}), \\ & \text{subject to } \mathbf{g}(\mathbf{x}) \leq \mathbf{0}^m, \\ & \mathbf{x} \in X, \end{aligned}$$

where $X \subseteq \mathbb{R}^n$, $f : \mathbb{R}^n \rightarrow \mathbb{R}$, and $\mathbf{g} : \mathbb{R}^n \rightarrow \mathbb{R}^m$. If the restrictions $\mathbf{g}(\mathbf{x}) \leq \mathbf{0}^m$ are complicating side constraints which are Lagrangian relaxed, we obtain the Lagrangian dual problem to

$$\text{maximize}_{\boldsymbol{\mu} \geq \mathbf{0}^m} q(\boldsymbol{\mu}),$$

where

$$q(\boldsymbol{\mu}) := \text{minimum}_{\mathbf{x} \in X} \{f(\mathbf{x}) + \boldsymbol{\mu}^T \mathbf{g}(\mathbf{x})\}, \quad \boldsymbol{\mu} \in \mathbb{R}^m.$$

(a) Suppose that the set X is finite (for example, consisting of a finite number of integer vectors). Denote the elements of X by \mathbf{x}^p , $p = 1, \dots, P$. Show that the dual objective function is piece-wise linear. How many linear segments can it have, at most? Why is it *not* always built up by that many segments?

[Note: This property holds regardless of any properties of f and \mathbf{g} .]

(b) Illustrate the result in (a) on the linear 0/1 problem to find

$$\begin{aligned} z^* = & \text{maximum } z = 5x_1 + 8x_2 + 7x_3 + 9x_4, \\ & \text{subject to } \quad 3x_1 + 2x_2 + 2x_3 + 4x_4 \leq 5, \\ & \quad 2x_1 + x_2 + 2x_3 + x_4 = 3, \\ & \quad x_1, x_2, x_3, x_4 \in \{0, 1\}, \end{aligned}$$

where the first constraint is considered complicating.

(c) Suppose that the function f and all components of \mathbf{g} are linear, and that the set X is a polytope (that is, a bounded polyhedron). Show that the dual objective function is also in this case piece-wise linear. How many linear pieces can it be built from, at most?

Exercise 6.9 (Lagrangian relaxation) Consider the problem to

$$\begin{aligned} & \text{minimize } z = 2x_1 + x_2, \\ & \text{subject to } \quad x_1 + x_2 \geq 5, \\ & \quad x_1 \leq 4, \\ & \quad x_2 \leq 4, \\ & \quad x_1, x_2 \geq 0, \text{ integer.} \end{aligned}$$

Lagrangian relax the first constraint. Describe the Lagrangian function and the dual problem. Calculate the Lagrangian dual function at these four points: $\mu = 0, 1, 2, 3$. Give the best lower and upper bounds on the optimal value of the original problem that you have found.

Lagrangian duality

Exercise 6.10 (surrogate relaxation) Consider an optimization problem of the form

$$\begin{aligned} & \text{minimize } f(\mathbf{x}), \\ & \text{subject to } g_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m, \quad (P) \\ & \quad \mathbf{x} \in X, \end{aligned}$$

where the functions $f, g_i : \mathbb{R}^n \rightarrow \mathbb{R}$ are continuous and the set $X \subset \mathbb{R}^n$ is closed and bounded. The problem is assumed to have an optimal solution, \mathbf{x}^* . Introduce parameters $\mu_i \geq 0, i = 1, \dots, m$, and define

$$\begin{aligned} s(\boldsymbol{\mu}) &:= \text{minimum } f(\mathbf{x}), \\ & \text{subject to } \boldsymbol{\mu}^T \mathbf{g}(\mathbf{x}) \leq 0, \quad (S) \\ & \quad \mathbf{x} \in X. \end{aligned}$$

This problem therefore has exactly one explicit constraint.

(a) [weak duality] Show that \mathbf{x}^* is a feasible solution to the problem (S) and that $s(\boldsymbol{\mu}) \leq f^*$ therefore always holds, that is, the problem (S) is a *relaxation* of the original one. Motivate also why $\text{maximum}_{\boldsymbol{\mu} \geq \mathbf{0}^m} s(\boldsymbol{\mu}) \leq f^*$ must hold. Explain the potential usefulness of this result!

(b) [example] Consider the linear 0/1 problem

$$\begin{aligned} z^* &= \text{maximum } z = 5x_1 + 8x_2 + 7x_3 + 9x_4, \\ & \text{subject to } \begin{aligned} 3x_1 + 2x_2 + 3x_3 + 3x_4 &\leq 6, & (1) \\ 2x_1 + 3x_2 + 3x_3 + 4x_4 &\leq 5, & (2) \\ 2x_1 + x_2 + 2x_3 + x_4 &= 3, \\ x_1, x_2, x_3, x_4 &\in \{0, 1\}. \end{aligned} \end{aligned}$$

Surrogate relax the constraints (1) and (2) with multipliers $\mu_1, \mu_2 \geq 0$ and formulate the problem (S). Let $\bar{\boldsymbol{\mu}} = (1, 2)^T$. Calculate $s(\bar{\boldsymbol{\mu}})$.

Consider again the original problem and Lagrangian relax the constraints (1) and (2) with multipliers $\mu_1, \mu_2 \geq 0$. Calculate the Lagrangian dual objective value at $\boldsymbol{\mu} = \bar{\boldsymbol{\mu}}$.

Compare the two results!

(c) [comparison with Lagrangian duality] Let $\boldsymbol{\mu} \geq \mathbf{0}^m$ and

$$q(\boldsymbol{\mu}) := \text{minimum}_{\mathbf{x} \in X} \{f(\mathbf{x}) + \boldsymbol{\mu}^T \mathbf{g}(\mathbf{x})\}.$$

Show that $q(\boldsymbol{\mu}) \leq s(\boldsymbol{\mu})$, and that

$$\text{maximum}_{\boldsymbol{\mu} \geq \mathbf{0}^m} q(\boldsymbol{\mu}) \leq \text{maximum}_{\boldsymbol{\mu} \geq \mathbf{0}^m} s(\boldsymbol{\mu}) \leq f^*$$

holds.

Part IV

Linear Programming

Linear programming: An introduction

VII

Linear programming (LP) models, that is, the collection of optimization models with linear objective functions and polyhedral feasible regions, are very useful in practice. The reason for this is that many real world problems can be described by LP models (even if several approximations must be made first) and, perhaps as importantly, there exist efficient algorithms for solving linear programs; the most famous of them is the simplex method, which will be presented in Chapter 9. Often, LP models deal with situations where a number of resources (materials, machines, people, land, etc.) are available and are to be combined to yield several products.

To introduce the concept of linear programming we use a simplified manufacturing problem. In Section 7.1 we describe the problem. From the problem description we develop an LP model in Section 7.2. It turns out that the LP model only contains two variables. Hence it is possible to solve the problem graphically, which is done in Section 7.3. In Section 7.4 we discuss what happens if the data of the problem is modified. Namely, we see how the optimal solution changes if the supply of raw material or the prices of the products are modified. Finally, in Section 7.5 we develop what we call the linear programming *dual* problem to the manufacturing problem.

7.1 The manufacturing problem

A manufacturer produces two pieces of furniture: tables and chairs. The production of the furniture requires the use of two different pieces of raw material: large and small pieces. A table is assembled by putting together two pieces of each, while a chair is assembled from one of the larger pieces and two of the smaller pieces (see Figure 7.1).

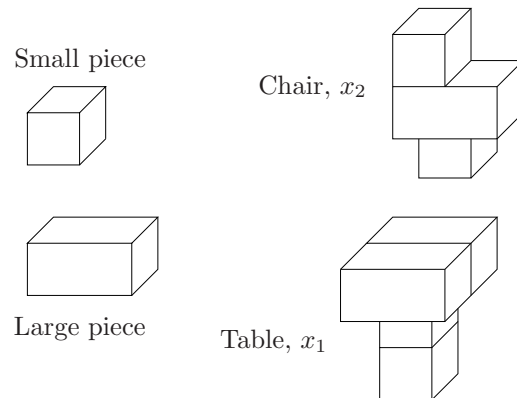


Figure 7.1: Illustration of the manufacturing problem.

When determining the optimal production plan, the manufacturer must take into account that only 6 large and 8 small pieces are available. A table is sold for 1600 SEK, while a chair sells for 1000 SEK. Under the assumption that all items produced can be sold, and that the raw material has already been paid for, the problem is to determine the production plan that maximizes the total income, within the limited resources.

7.2 A linear programming model

In order to develop a linear programming model for the manufacturing problem we introduce the following variables:

x_1 = number of tables manufactured and sold,
 x_2 = number of chairs manufactured and sold,
 z = total income.

The variable z is, strictly speaking, not a variable, but will be defined by the variables x_1 and x_2 .

The income from each product is given by the price of the product multiplied by the number of products sold. Hence, the total income is

$$z = 1600x_1 + 1000x_2. \quad (7.1)$$

Given that we produce x_1 tables and x_2 chairs the required number of large pieces is $2x_1 + x_2$ and the required number of small peaces is

$2x_1 + 2x_2$. But only 6 large pieces and 8 small pieces are available, so we must have that

$$2x_1 + x_2 \leq 6, \quad (7.2)$$

$$2x_1 + 2x_2 \leq 8. \quad (7.3)$$

Further, it is impossible to produce a negative number of chairs or tables, which gives that

$$x_1, x_2 \geq 0. \quad (7.4)$$

(Also, the number of chairs and tables produced must be integers, but we will not take that into account here.)

Now the objective is to maximize the total income, so if we combine the income function (7.1) and the constraints (7.2)–(7.4) we get the following linear programming model:

$$\begin{aligned} &\text{maximize} && z = 1600x_1 + 1000x_2, && (7.5) \\ &\text{subject to} && 2x_1 && + x_2 \leq 6, \\ &&& 2x_1 && + 2x_2 \leq 8, \\ &&& x_1, && x_2 \geq 0. \end{aligned}$$

7.3 Graphical solution

The feasible region of the linear programming formulation (7.5) is shown in Figure 7.2. The figure also includes lines corresponding to different values of the cost function. For example, the line $z = 0 = 1600x_1 + 1000x_2$ passes through the origin, and the line $z = 2600 = 1600x_1 + 1000x_2$ passes through the point $(1, 1)^T$. We see that the value of the cost function increases as these lines move upward, and it follows that the optimal solution is $\mathbf{x}^* = (2, 2)^T$ and $z^* = 5200$. Observe that the optimal solution is an extreme point, which is in accordance with Theorem 4.12. This fact will be very important in the development of the simplex method in Chapter 9.

7.4 Sensitivity analysis

In this section we investigate how the optimal solution changes if the data of the problem is changed. We consider three different changes (made independently of each other), namely

1. an increase in the number of large pieces available;

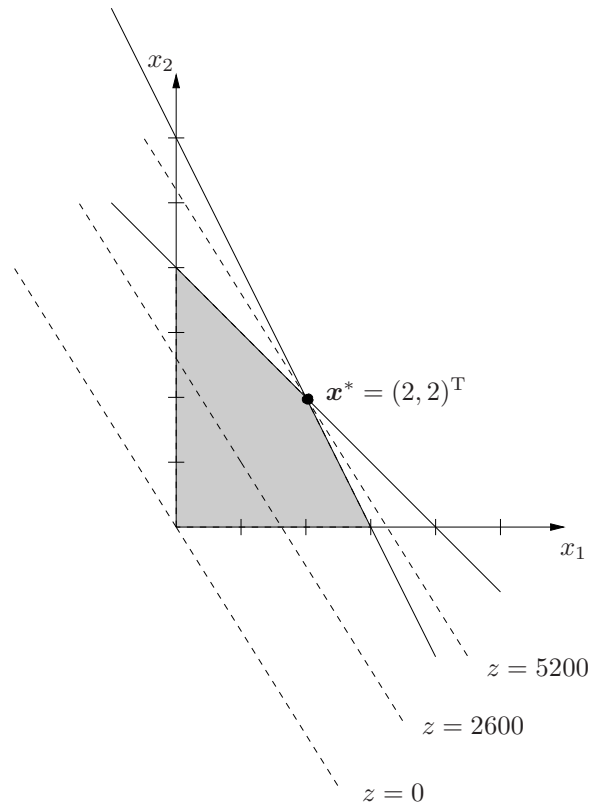


Figure 7.2: Graphical solution of the manufacturing problem.

2. an increase in the number of small pieces available; and
3. a decrease in the price of the tables.

7.4.1 An increase in the number of large pieces available

Assume that the number of large pieces available increases from 6 to 7. Then the linear program becomes

$$\begin{aligned}
 &\text{maximize} && z = 1600x_1 + 1000x_2, \\
 &\text{subject to} && 2x_1 + x_2 \leq 7, \\
 & && 2x_1 + 2x_2 \leq 8, \\
 & && x_1, \quad x_2 \geq 0.
 \end{aligned}$$

The feasible region is shown in Figure 7.3.

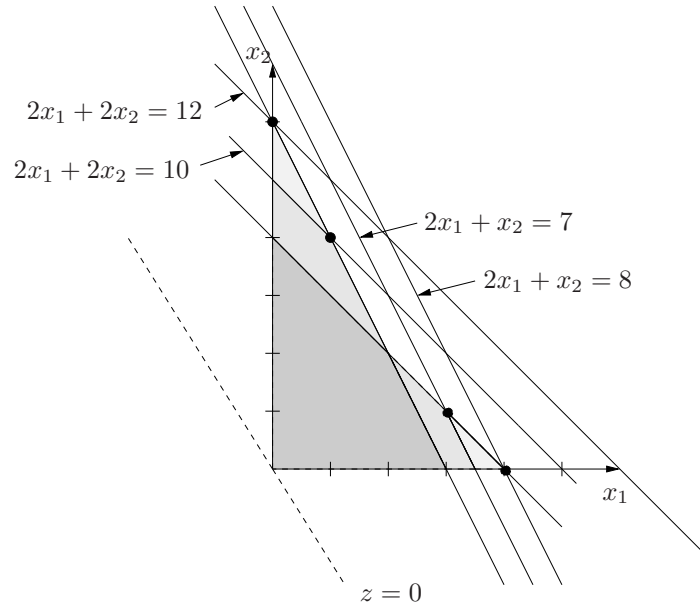


Figure 7.3: An increase in the number of large and small pieces available.

We see that the optimal solution becomes $(3,1)^T$ and $z^* = 5800$, which means that an additional large piece increases the income by $5800 - 5200 = 600$. Hence the *shadow price* of the large pieces is 600. The figure also illustrates what happens if the number of large pieces is 8. Then the optimal solution becomes $(4,0)^T$ and $z^* = 6400$. But what happens if we increase the number of large pieces further? From the figure it follows that the optimal solution will not change (since $x_2 \geq 0$ must apply), so an increase larger than 2 in the number of large pieces gives no further income. This illustrates that the validity of the shadow price depends on the actual increment; exactly when the shadow price is valid is investigated in Theorem 10.8 and Remark 10.9.

7.4.2 An increase in the number of small pieces available

Starting from the original setup, in the same manner as for the large pieces it follows from Figure 7.3 that two additional small pieces give the new optimal solution $\mathbf{x}^* = (1,4)^T$ and $z^* = 5600$, so the income per

additional small piece is $(5600 - 5200)/2 = 200$. Hence the shadow price of the small pieces is 200. However, no more than 4 small pieces are worth this price, since $x_1 \geq 0$ must apply.

7.4.3 A decrease in the price of the tables

Now assume that the price of tables is decreased from 1600 to 800. The new linear program becomes

$$\begin{aligned} \text{maximize} \quad & z = 800x_1 + 1000x_2, \\ \text{subject to} \quad & 2x_1 + x_2 \leq 6, \\ & 2x_1 + 2x_2 \leq 8, \\ & x_1, \quad x_2 \geq 0. \end{aligned}$$

This new situation is illustrated in Figure 7.4, from which we see that

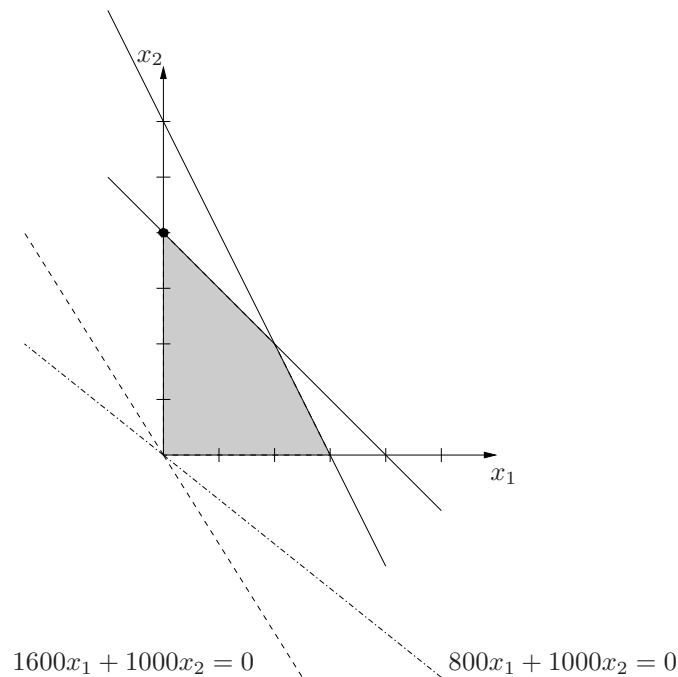


Figure 7.4: A decrease in the price of the tables.

the optimal solution is $(0, 4)^T$, that is, we will not produce any tables. This is natural, since it takes the same number of small pieces to produce

a table and a chair but the table requires one more large piece, and in addition the price of a table is now lower than that of a chair.

7.5 The dual of the manufacturing problem

7.5.1 A competitor

Suppose that another manufacturer (let us call them Billy) produce bookshelves whose raw material is identical to those used for the tables and chairs, that is, the small and large pieces. Billy wish to expand their production, and are interested in acquiring the resources that “our” factory sits on. Let us ask ourselves two questions, which (as we shall see) have identical answers: (1) what is the lowest bid (price) for the total capacity at which we are willing to sell?; (2) what is the highest bid (price) that Billy are prepared to offer for the resources? The answer to those two questions is a measure of the wealth of the company in terms of their resources.

7.5.2 A dual problem

To study the problem, we introduce the variables

y_1 = the price which Billy offers for each large piece,
 y_2 = the price which Billy offers for each small piece,
 w = the total bid which Billy offers.

In order to accept to sell our resources, it is reasonable to require that the price offered is at least as high as the value that the resources represent in our optimal production plan, as otherwise we would earn more by using the resources ourselves. Consider, for example, the net income on a table sold. It is 1600 SEK, and for that we use two large and two small pieces. The bid would therefore clearly be too low unless $2y_1 + 2y_2 \geq 1600$. The corresponding requirement for the chairs is that $y_1 + 2y_2 \geq 1000$.

Billy is interested in minimizing the total bid, under the condition that the offer is accepted. Observing that y_1 and y_2 are prices and therefore non-negative, we have the following mathematical model for Billy’s problem:

$$\begin{array}{ll} \text{minimize} & w = 6y_1 + 8y_2, \\ \text{subject to} & 2y_1 + 2y_2 \geq 1600, \\ & y_1 + 2y_2 \geq 1000, \\ & y_1, \quad y_2 \geq 0. \end{array} \tag{7.6}$$

This is usually called the *dual problem* of our production planning problem (which would then be the *primal problem*).

The optimal solution to this problem is $\mathbf{y}^* = (600, 200)^T$. (Check this!) The total offer is $w^* = 5200$.

Remark 7.1 (the linear programming dual) Observe that the dual problem (7.6) is in accordance with the Lagrangian duality theory of Section 6.2.4. The linear programming dual will be discussed further in Chapter 10. ■

7.5.3 Interpretations of the dual optimal solution

From the above we see that the dual optimal solution is identical to the shadow prices for the resource (capacity) constraints. (This is indeed a general conclusion in linear programming.) To motivate that this is reasonable in our setting, we may consider Billy as a fictitious competitor only, which we use together with the dual problem to measure the value of our resources. This (fictitious) measure can be used to create internal prices in a company in order to utilize limited resources as efficiently as possible, especially if the resource is common to several independent sub-units. The price that the dual optimal solution provides will then be a price directive for the sub-units, that will make them utilize the scarce resources in a manner which is optimal for the overall goal.

We note that the optimal value $z^* = 5200$ of the production agrees with the total value $w^* = 5200$ of the resources in our company. (This is also a general result in linear programming; see the Strong Duality Theorem 10.6.) Billy will of course not pay more than what the resource is worth, but can at the same time not offer less than the profit that our company can make ourselves, since we would then not agree to sell. It follows immediately that for each feasible production plan \mathbf{x} and price \mathbf{y} , it holds that $z \leq w$, since

$$\begin{aligned} z &= 1600x_1 + 1000x_2 \leq (2y_1 + 2y_2)x_1 + (y_1 + 2y_2)x_2 \\ &= y_1(2x_1 + x_2) + y_2(2x_1 + 2x_2) \leq 6y_1 + 8y_2 = w, \end{aligned}$$

where in the inequalities we utilize all the constraints of the primal and dual problems. (Also this fact is general in linear programming; see the Weak Duality Theorem 10.4.) So, each offer accepted (from our point of view) must necessarily be an upper bound on our own possible profit, and this upper bound is what Billy wish to minimize in the dual problem.

Linear programming models

VIII

We begin this chapter with a presentation of the axioms underlying the use of linear programming models and discuss the modelling process. Then, in Section 8.2, the geometry of linear programming is studied. It is shown that every linear program can be transformed into the *standard form* which is the form that the simplex method requires. Further, we introduce the concept of *basic feasible solution* and discuss its connection to extreme points. A version of the Representation Theorem adapted to the standard form is presented, and we show that if there exists an optimal solution to a linear program in standard form, then there exists an optimal solution among the basic feasible solutions. Finally, we define the term *adjacent extreme point* and give an algebraic characterization of adjacency which actually proves that the simplex method at each iteration step moves from one extreme point to an adjacent one.

8.1 Linear programming modelling

Many real world situations can be modelled as linear programs. However, the applicability of a linear program requires certain axioms to be fulfilled. Hence, often approximations of the real world problem must be made prior to the formulation of a linear program. The axioms underlying the use of linear programming models are:

- proportionality (linearity: no economies-of-scale; no fixed costs);
- additivity (no substitute-time-effects);
- divisibility (continuity); and
- determinism (no randomness).

George Dantzig presented the linear programming model and the simplex method for solving it at an econometrics conference in Wisconsin in the late 40s. The economist Hotelling stood up, devastatingly smiling, and stated that “But we all know the world is nonlinear.” The young graduate student George Dantzig could not respond, but was defended by John von Neumann, who stood up and concluded that “The speaker titled his talk ‘linear programming’ and carefully stated his axioms. If you have an application that satisfies the axioms, well use it. If it does not, then don’t”; he sat down, and Hotelling was silenced. (See Dantzig’s account of the early history of linear programming in [LRS91, pp. 19–31].)

Now if the problem considered (perhaps after approximations) fulfills the axioms above, then it can be formulated as a linear program. However, in practical modelling situations we usually do not talk about the axioms; they naturally appear when a linear program is formulated.

To formulate a real world problem as a linear program is an art in itself, and unfortunately there is little theory to help in formulating the problem in this way. The general approach can however be described by two steps:

1. Prepare a list of all the decision variables in the problem. This list must be complete in the sense that if an optimal solution providing the values of each of the variables is obtained, then the decision maker should be able to translate it into an optimum policy that can be implemented.
2. Use the variables from step 1 to formulate all the constraints and the objective function of the problem.

We illustrate the two-step modelling process by an example.

Example 8.1 (the transportation problem) In the transportation problem we have a set of nodes or locations called *sources*, which have a commodity available for shipment, and another set of locations called *demand centers*, or *sinks*, which require this commodity. The amount of commodity available at each source and the amount required at each demand center are specified, as well as the cost per unit of transporting the commodity from each source to each demand center. The problem is to determine the quantity to be transported from each source to each demand center, so as to meet all the requirements at minimum total shipping cost.

Consider the problem where the commodity is iron ore, the sources are found at mines 1 and 2, where the ore is produced, and the demand

Table 8.1: Unit cost of shipping ore from mine to steel plant (KSEK per Mton).

Plant	1	2	3
Mine 1	9	16	28
Mine 2	14	29	19

centers are three steel plants. The unit costs of shipping ore from each mine to each steel plant are given in Table 8.1.

Further, the amount of ore available at the mines and the Mtons of ore required at each steel plant are given in the Tables 8.2 and 8.3.

Table 8.2: Amount of ore available at the mines (Mtons).

Mine 1	103
Mine 2	197

Table 8.3: Ore requirements at the steel plants (Mtons).

Plant 1	71
Plant 2	133
Plant 3	96

We use the two-step modelling process to formulate a linear programming model.

Step 1: The activities in the transportation model are to ship ore from mine i to steel plant j for $i = 1, 2$ and $j = 1, 2, 3$. It is convenient to represent the variables corresponding to the levels at which these activities are carried out by double subscripted symbols. Hence, for $i = 1, 2$ and $j = 1, 2, 3$, we introduce the following variables:

x_{ij} = amount of ore (in Mtons) shipped from mine i to steel plant j .

We also introduce a variable corresponding to the total cost of the shipping:

z = total shipping cost.

Step 2: The transportation problem considered is illustrated in Figure 8.1.

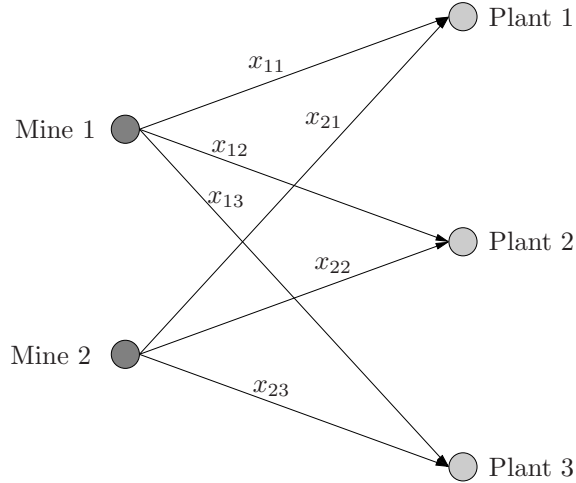


Figure 8.1: Illustration of the transportation problem.

The items in this problem are the ore at various locations. Consider the ore at mine 1. According to Table 8.2 there are only 103 Mtons of it available, and the amount of ore shipped out of mine 1, which is $x_{11} + x_{12} + x_{13}$, cannot exceed the amount available, leading to the constraint

$$x_{11} + x_{12} + x_{13} \leq 103.$$

Likewise, if we consider ore at mine 2 we get the constraint

$$x_{21} + x_{22} + x_{23} \leq 197.$$

Further, at steel plant 1 according to Table 8.3 there are at least 71 Mtons of ore required, so the amount of ore shipped to steel plant 1 has to be greater than or equal to this amount, leading to the constraint

$$x_{11} + x_{21} \geq 71.$$

In the same manner, for the steel plants 2 and 3 we get

$$x_{12} + x_{22} \geq 133,$$

$$x_{13} + x_{23} \geq 96.$$

Of course it is impossible to ship a negative amount of ore, yielding the constraints

$$x_{ij} \geq 0, \quad i = 1, 2, \quad j = 1, 2, 3.$$

From Table 8.1 follows that the total cost (in KSEK) of shipping is

$$z = 9x_{11} + 16x_{12} + 28x_{13} + 14x_{21} + 29x_{22} + 19x_{23}.$$

Finally, since the objective is to minimize the total cost we get the following linear programming model:

$$\begin{array}{ll} \text{minimize} & z = 9x_{11} + 16x_{12} + 28x_{13} + 14x_{21} + 29x_{22} + 19x_{23}, \\ \text{subject to} & \begin{array}{rcl} x_{11} & +x_{12} & +x_{13} & \leq 103, \\ & & x_{21} & +x_{22} & +x_{23} & \leq 197, \\ x_{11} & & & +x_{21} & & \geq 71, \\ & x_{12} & & & +x_{22} & \geq 133, \\ & & x_{13} & & & +x_{23} & \geq 96, \\ x_{11}, & x_{12}, & x_{13}, & x_{21}, & x_{22}, & x_{23} & \geq 0. \end{array} \end{array}$$

The transportation problem may be given in a compact general formulation. Assume that we have N sources and M demand centers. For $i = 1, \dots, N$ and $j = 1, \dots, M$, introduce the variables

x_{ij} = amount of commodity shipped from source i to demand center j ,
and let

$$z = \text{total shipping cost.}$$

Further for $i = 1, \dots, N$ and $j = 1, \dots, M$ introduce the shipping costs

c_{ij} = unit cost of shipping commodity from source i to demand center j .

Also, let

s_i = amount of commodity available at source i , $i = 1, \dots, N$,

d_j = amount of commodity required at demand center j , $j = 1, \dots, M$.

Consider source i . The amount of commodity available is given by s_i , which gives the constraint

$$\sum_{j=1}^M x_{ij} \leq s_i.$$

Linear programming models

Similarly, the amount of commodity required at demand center j is given by d_j , leading to the constraint

$$\sum_{i=1}^N x_{ij} \geq d_j.$$

It is impossible to ship a negative amount of commodity, which gives

$$x_{ij} \geq 0, \quad i = 1, \dots, N, \quad j = 1, \dots, M.$$

Finally, the total cost for shipping is

$$z = \sum_{i=1}^N \sum_{j=1}^M c_{ij} x_{ij},$$

leading to the following linear program:

$$\begin{aligned} \text{minimize} \quad & z = \sum_{i=1}^N \sum_{j=1}^M c_{ij} x_{ij}, \\ \text{subject to} \quad & \sum_{j=1}^M x_{ij} \leq s_i, \quad i = 1, \dots, N, \\ & \sum_{i=1}^N x_{ij} \geq d_j, \quad j = 1, \dots, M, \\ & x_{ij} \geq 0, \quad i = 1, \dots, N, \quad j = 1, \dots, M. \end{aligned}$$

If, for some reason, it is impossible to transport any commodities from a source i to a sink j , then we may either remove this variable altogether from the model, or, more simply, give it the unit price $c_{ij} = +\infty$.

Note, finally, that there exists a feasible solution to the transportation problem if and only if $\sum_{i=1}^N s_i \geq \sum_{j=1}^M d_j$. ■

8.2 The geometry of linear programming

In Section 3.2 we studied the class of feasible sets in linear programming, namely the sets of polyhedra; they are sets of the form

$$P := \{ \mathbf{x} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{x} \leq \mathbf{b} \},$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$. In particular, we proved the Representation Theorem 3.22 and promised that it would be important in the

development of the simplex method. In this section we consider polyhedra of the form

$$P := \{ \mathbf{x} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{x} = \mathbf{b}; \quad \mathbf{x} \geq \mathbf{0}^n \}, \quad (8.1)$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$, and $\mathbf{b} \in \mathbb{R}^m$ is such that $\mathbf{b} \geq \mathbf{0}^m$. The advantage of this form is that the constraints (except for the non-negativity constraints) are *equalities*, which admits pivot operations to be carried out. The simplex method uses pivot operations at each iteration step and hence it is necessary that the polyhedron (that is, the feasible region) is represented in the form (8.1). This is, however, not a restriction, as we will see in Section 8.2.1, since every polyhedron can be transformed into this form! We will use the term *standard form* when a polyhedron is represented in the form (8.1). In Section 8.2.2 we introduce the concept of *basic feasible solution* and show that each basic feasible solution corresponds to an extreme point. We also restate the Representation Theorem 3.22 for polyhedra in standard form and prove that if there exists an optimal solution to a linear program, then there exists an optimal solution among the extreme points. Finally, in Section 8.2.3, a strong connection between basic feasible solutions and adjacent extreme points is discussed. This connection shows that the simplex method at each iteration step moves from an extreme point to an adjacent extreme point.

8.2.1 Standard form

A linear programming problem in *standard form* is a problem of the form

$$\begin{aligned} &\text{minimize} && z = \mathbf{c}^T \mathbf{x}, \\ &\text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b}, \\ &&& \mathbf{x} \geq \mathbf{0}^n, \end{aligned}$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \geq \mathbf{0}^m$. The purpose of this section is to show that every linear program can be transformed into the standard form. In order to do that we must

- express the objective function in minimization form;
- transform all the constraints into equality constraints with non-negative right-hand sides; and
- transform any unrestricted and non-positive variables into non-negative ones.

Objective function

Constant terms in the objective function will not change the set of optimal solutions and can therefore be eliminated. If the objective is to

$$\text{maximize } z = \mathbf{c}^T \mathbf{x},$$

then change the objective function so that the objective becomes

$$\text{minimize } \tilde{z} := -z = -\mathbf{c}^T \mathbf{x}.$$

This change does not affect the set of feasible solutions to the problem and the equation

$$[\text{maximum value of } z] = -[\text{minimum value of } \tilde{z}]$$

can be used to get the maximum value of the original objective function.

Inequality constraints and negative right-hand sides

Consider the inequality constraint

$$\mathbf{a}^T \mathbf{x} \leq b,$$

where $\mathbf{a} \in \mathbb{R}^n$ and $b \in \mathbb{R}$. By introducing a non-negative *slack variable* s this constraint can be written as

$$\mathbf{a}^T \mathbf{x} + s = b, \tag{8.2a}$$

$$s \geq 0, \tag{8.2b}$$

which has the desired form of an equation. Suppose that $b < 0$. By multiplying both sides of (8.2a) by -1 the negativity in the right-hand side is eliminated and we are done. Similarly, a constraint of the form

$$\mathbf{a}^T \mathbf{x} \geq b,$$

can be written as

$$\mathbf{a}^T \mathbf{x} - s = b,$$

$$s \geq 0.$$

We call such variables s *surplus variables*.

Remark 8.2 (on the role of slack and surplus variables) Slack and surplus variables may appear to be only help variables, but they often have a clear interpretation as decision variables. Consider, for example, the

model (7.5) of a furniture production problem. The two inequality constraints are associated with the capacity of production stemming from the availability of raw material. Suppose then that we introduce slack variables in these constraints, which leads to the equivalent problem to

$$\begin{array}{llll} \text{maximize} & z = 1600x_1 + 1000x_2, & & \\ \text{subject to} & 2x_1 & +x_2 +s_1 & = 6, \\ & 2x_1 & +2x_2 & +s_2 = 8, \\ & x_1, & x_2, & s_1, s_2 \geq 0. \end{array}$$

The new variables s_1 and s_2 have the following interpretation: the value of s_i ($i = 1, 2$) is the level of inventory (or, remaining capacity of raw material of type i) that will be left unused when the production plan (x_1, x_2) has been implemented. This interpretation makes it clear that the values of s_1 and s_2 are clear consequences of our decision-making.

Surplus variables have a corresponding interpretation. In the case of the transportation model of the previous section, a demand constraint $(\sum_{i=1}^N x_{ij} \geq d_j, j = 1, \dots, M)$ may be fulfilled with equality (in which case the customer gets an amount exactly according to the demand) or it is fulfilled with strict inequality (in which case the customer gets a surplus of the product asked for). ■

Unrestricted and non-positive variables

Consider the linear program

$$\begin{array}{ll} \text{minimize} & z = \mathbf{c}^T \mathbf{x}, \\ \text{subject to} & \mathbf{A}\mathbf{x} = \mathbf{b}, \\ & x_2 \leq 0, \\ & x_j \geq 0, \quad j = 3, \dots, n, \end{array}$$

which is assumed to be in standard form except for the unrestricted variable x_1 and the non-positive variable x_2 . The x_2 -variable can be replaced by the non-negative variable $\tilde{x}_2 := -x_2$. The x_1 -variable can be transformed into the difference of two non-negative variables. Namely, introduce the variables $x_1^+ \geq 0$ and $x_1^- \geq 0$ and let $x_1 = x_1^+ - x_1^-$. Substituting x_1 with $x_1^+ - x_1^-$ wherever it occurs transforms the problem into standard form. The drawback of this method to handle unrestricted variables is that often the resulting problem is numerically unstable. However, there are other techniques to handle unrestricted variables that overcome this problem; one of them is discussed in Exercise 8.5.

Example 8.3 (standard form) Consider the linear program

$$\begin{aligned} \text{maximize } z &= 9x_1 - 7x_2 + 3y_1, \\ \text{subject to } 3x_1 + x_2 - y_1 &\leq 1, \\ 4x_1 - x_2 + 2y_1 &\geq 3, \\ x_1, x_2 &\geq 0. \end{aligned} \tag{8.4}$$

In order to transform the objective function into the minimization form, let

$$\tilde{z} := -z = -9x_1 + 7x_2 - 3y_1.$$

Further, by introducing the slack variable s_1 and the surplus variable s_2 the constraints can be transformed into an equality form by

$$\begin{aligned} 3x_1 + x_2 - y_1 + s_1 &= 1, \\ 4x_1 - x_2 + 2y_1 - s_2 &= 3, \\ x_1, x_2, s_1, s_2 &\geq 0. \end{aligned}$$

Finally, by introducing the variables y_1^+ and y_1^- we can handle the unrestricted variable y_1 by substituting it by $y_1^+ - y_1^-$ wherever it occurs. We arrive at the standard form to

$$\begin{aligned} \text{minimize } \tilde{z} &= -9x_1 + 7x_2 - 3y_1^+ + 3y_1^-, \\ \text{subject to } 3x_1 + x_2 - y_1^+ + y_1^- + s_1 &= 1, \\ 4x_1 - x_2 + 2y_1^+ - 2y_1^- - s_2 &= 3, \\ x_1, x_2, y_1^+, y_1^-, s_1, s_2 &\geq 0. \end{aligned} \tag{8.5}$$

Clearly, an optimal solution $(x_1, x_2, y_1^+, y_1^-, s_1, s_2)$ to (8.5) can be transformed into an optimal solution (x_1, x_2, y_1) to (8.4) by using the substitution $y_1 = y_1^+ - y_1^-$. ■

8.2.2 Basic feasible solutions and the Representation Theorem

In this section we introduce the concept of *basic feasible solution* and show the equivalence between extreme point and basic feasible solution. From this we can draw the conclusion that if there exists an optimal solution then there exists an optimal solution among the basic feasible solutions. This fact is crucial in the simplex method which searches for an optimal solution among the basic feasible solutions.

Consider a linear program in standard form,

$$\begin{aligned} & \text{minimize} && z = \mathbf{c}^T \mathbf{x}, \\ & \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b}, \\ & && \mathbf{x} \geq \mathbf{0}^n, \end{aligned} \tag{8.6}$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$ with $\text{rank } \mathbf{A} = \text{rank}(\mathbf{A}, \mathbf{b}) = m$ (otherwise, we can always delete rows), $n > m$, and $\mathbf{b} \geq \mathbf{0}^m$. A point $\tilde{\mathbf{x}}$ is a *basic solution* of (8.6) if

1. the equality constraints are satisfied at $\tilde{\mathbf{x}}$, that is, $\mathbf{A}\tilde{\mathbf{x}} = \mathbf{b}$; and
2. the columns of \mathbf{A} corresponding to the non-zero components of $\tilde{\mathbf{x}}$ are linearly independent.

A basic solution that also satisfies the non-negativity constraints is called a *basic feasible solution*, or, in short, a BFS.

Since $\text{rank } \mathbf{A} = m$, we can solve the system $\mathbf{A}\mathbf{x} = \mathbf{b}$ by selecting m variables of \mathbf{x} corresponding to m linearly independent columns of \mathbf{A} . Hence, we partition the columns of \mathbf{A} into two parts: one with $n - m$ columns of \mathbf{A} corresponding to components of \mathbf{x} that are set to 0; these are called the *non-basic* variables and are denoted by the subvector $\mathbf{x}_N \in \mathbb{R}^{n-m}$. The others represent the *basic variables*, and are denoted by $\mathbf{x}_B \in \mathbb{R}^m$. According to this partition,

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_B \\ \mathbf{x}_N \end{pmatrix}, \quad \mathbf{A} = (\mathbf{B}, \mathbf{N}),$$

which yields that

$$\mathbf{A}\mathbf{x} = \mathbf{B}\mathbf{x}_B + \mathbf{N}\mathbf{x}_N = \mathbf{b}.$$

Since $\mathbf{x}_N = \mathbf{0}^{n-m}$ by construction, we get the basic solution

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_B \\ \mathbf{x}_N \end{pmatrix} = \begin{pmatrix} \mathbf{B}^{-1}\mathbf{b} \\ \mathbf{0}^{n-m} \end{pmatrix}.$$

Further, if $\mathbf{B}^{-1}\mathbf{b} \geq \mathbf{0}^m$ then \mathbf{x} is a basic feasible solution.

Remark 8.4 (degenerate basic solution) If more than $n - m$ variables are zero at a basic solution \mathbf{x} , then the corresponding partition is not unique. Such a basic solution is called *degenerate*. ■

Example 8.5 (partitioning) Consider the linear program

$$\begin{aligned} \text{minimize } z &= 4x_1 + 3x_2 + 7x_3 - 2x_4, \\ \text{subject to } & \begin{aligned} x_1 & & -x_3 & & & & = 3, \\ x_1 & -x_2 & & -2x_4 & & & = 1, \\ 2x_1 & & & & +x_4 +x_5 & = 7, \\ x_1, & x_2, & x_3, & x_4, & x_5 & \geq 0. \end{aligned} \end{aligned}$$

The constraint matrix and the right-hand side vector are given by

$$\mathbf{A} = \begin{pmatrix} 1 & 0 & -1 & 0 & 0 \\ 1 & -1 & 0 & -2 & 0 \\ 2 & 0 & 0 & 1 & 1 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 3 \\ 1 \\ 7 \end{pmatrix}.$$

(a) The partition $\mathbf{x}_B = (x_2, x_3, x_4)^T$, $\mathbf{x}_N = (x_1, x_5)^T$,

$$\mathbf{B} = \begin{pmatrix} 0 & -1 & 0 \\ -1 & 0 & -2 \\ 0 & 0 & 1 \end{pmatrix}, \quad \mathbf{N} = \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{pmatrix},$$

corresponds to the basic solution

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_B \\ \mathbf{x}_N \end{pmatrix} = \begin{pmatrix} x_2 \\ x_3 \\ x_4 \\ x_1 \\ x_5 \end{pmatrix} = \begin{pmatrix} \mathbf{B}^{-1}\mathbf{b} \\ \mathbf{0}^2 \end{pmatrix} = \begin{pmatrix} -15 \\ -3 \\ 7 \\ 0 \\ 0 \end{pmatrix}.$$

This is, however, not a basic *feasible* solution (since x_2 and x_3 are negative).

(b) The partition, $\mathbf{x}_B = (x_1, x_2, x_5)^T$, $\mathbf{x}_N = (x_3, x_4)^T$,

$$\mathbf{B} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & -1 & 0 \\ 2 & 0 & 1 \end{pmatrix}, \quad \mathbf{N} = \begin{pmatrix} -1 & 0 \\ 0 & -2 \\ 0 & 1 \end{pmatrix},$$

corresponds to the basic solution

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_B \\ \mathbf{x}_N \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 \\ x_5 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} \mathbf{B}^{-1}\mathbf{b} \\ \mathbf{0}^2 \end{pmatrix} = \begin{pmatrix} 3 \\ 2 \\ 1 \\ 0 \\ 0 \end{pmatrix}.$$

This is clearly a basic feasible solution.

(c) The partition $\mathbf{x}_B = (x_2, x_4, x_5)^T$, $\mathbf{x}_N = (x_1, x_3)^T$,

$$\mathbf{B} = \begin{pmatrix} 0 & 0 & 0 \\ -1 & -2 & 0 \\ 0 & 1 & 1 \end{pmatrix}, \quad \mathbf{N} = \begin{pmatrix} 1 & -1 \\ 1 & 0 \\ 2 & 0 \end{pmatrix},$$

does not correspond to a basic solution since the system $\mathbf{B}\mathbf{x}_B = \mathbf{b}$ is infeasible.

(d) Finally, the partition $\mathbf{x}_B = (x_1, x_4, x_5)^T$, $\mathbf{x}_N = (x_2, x_3)^T$,

$$\mathbf{B} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & -2 & 0 \\ 2 & 1 & 1 \end{pmatrix}, \quad \mathbf{N} = \begin{pmatrix} 0 & -1 \\ -1 & 0 \\ 0 & 0 \end{pmatrix},$$

corresponds to the basic feasible solution

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_B \\ \mathbf{x}_N \end{pmatrix} = \begin{pmatrix} x_1 \\ x_4 \\ x_5 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} \mathbf{B}^{-1}\mathbf{b} \\ \mathbf{0}^2 \end{pmatrix} = \begin{pmatrix} 3 \\ 1 \\ 0 \\ 0 \\ 0 \end{pmatrix},$$

which is *degenerate* (since a basic variable, x_5 , has value zero). ■

Remark 8.6 (partitioning) The above partitioning technique will be used frequently in what follows and from now on when we say that $\mathbf{A} = (\mathbf{B}, \mathbf{N})$ is a *partition* of \mathbf{A} we will always mean that the columns of \mathbf{A} and the variables of \mathbf{x} have been rearranged so that \mathbf{B} corresponds to the basic variables \mathbf{x}_B and \mathbf{N} to the non-basic variables \mathbf{x}_N . ■

We are now ready to prove the equivalence between extreme point and basic feasible solution.

Theorem 8.7 (equivalence between extreme point and BFS) *Suppose $\mathbf{A} \in \mathbb{R}^{m \times n}$ with $\text{rank } \mathbf{A} = m$, and $\mathbf{b} \in \mathbb{R}^m$. Then, a vector $\mathbf{x} \in \mathbb{R}^n$ is an extreme point of the set $\{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{x} = \mathbf{b}; \mathbf{x} \geq \mathbf{0}^n\} \neq \emptyset$ if and only if it is a basic feasible solution.*

Proof. Let \mathbf{x} be a basic feasible solution with the corresponding partition $\mathbf{A} = (\mathbf{B}, \mathbf{N})$, where $\text{rank } \mathbf{B} = m$ (such a partition exists since $\text{rank } \mathbf{A} = m$). Then the equality subsystem (see Section 3.2.3) of

$$\begin{aligned} \mathbf{A}\mathbf{x} &= \mathbf{b}, \\ \mathbf{x} &\geq \mathbf{0}^n \end{aligned}$$

is given by

$$\begin{aligned} \mathbf{B}\mathbf{x}_B + \mathbf{N}\mathbf{x}_N &= \mathbf{b}, \\ \mathbf{x}_N &= \mathbf{0}^{n-m} \end{aligned}$$

(if some of the basic variables equal zero we get additional rows of the form “ $x_i = 0$ ” but these will not affect the proof). Since $\text{rank } \mathbf{B} = m$ it follows that

$$\text{rank} \begin{pmatrix} \mathbf{B} & \mathbf{N} \\ \mathbf{0}^{(n-m) \times m} & \mathbf{I}^{n-m} \end{pmatrix} = n.$$

The result then follows from Theorem 3.17. ■

Remark 8.8 (degenerate extreme point) An extreme point that corresponds to more than one BFS is said to be degenerate. This typically occurs when we have redundant constraints. ■

We present a reformulation of the Representation Theorem 3.22 that is adapted to the standard form. Consider the polyhedral cone $C := \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{x} = \mathbf{0}^m; \mathbf{x} \geq \mathbf{0}^n\}$. From Theorem 3.28 it follows that C is finitely generated, that is, there exist vectors $\mathbf{d}^1, \dots, \mathbf{d}^r \in \mathbb{R}^n$ such that

$$C = \text{cone}\{\mathbf{d}^1, \dots, \mathbf{d}^r\} := \left\{ \mathbf{x} \in \mathbb{R}^n \mid \mathbf{x} = \sum_{i=1}^r \alpha_i \mathbf{d}^i; \alpha_1, \dots, \alpha_r \geq 0 \right\}.$$

There are, of course, infinitely many ways to generate a certain polyhedral cone C . Assume that $C = \text{cone}\{\mathbf{d}^1, \dots, \mathbf{d}^r\}$. If there exists a vector $\mathbf{d}^i \in \{\mathbf{d}^1, \dots, \mathbf{d}^r\}$ such that

$$\mathbf{d}^i \in \text{cone}\{\{\mathbf{d}^1, \dots, \mathbf{d}^r\} \setminus \{\mathbf{d}^i\}\},$$

then \mathbf{d}^i is not necessary in the description of C . If we similarly continue to remove vectors from $\{\mathbf{d}^1, \dots, \mathbf{d}^r\}$, one at a time, we end up with a set generating C such that none of the vectors of the set can be written as a non-negative linear combination of the others. Such a set is naturally called the set of *extreme directions* of C (cf. Definition 3.11 of extreme point).

Theorem 8.9 (Representation Theorem) Let $P := \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{x} = \mathbf{b}; \mathbf{x} \geq \mathbf{0}^n\}$ and $V := \{\mathbf{v}^1, \dots, \mathbf{v}^k\}$ be the set of extreme points of P . Further, let $C := \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{x} = \mathbf{0}^m; \mathbf{x} \geq \mathbf{0}^n\}$ and $D := \{\mathbf{d}^1, \dots, \mathbf{d}^r\}$ be the set of extreme directions of C .

- (a) P is nonempty if and only if V is nonempty (and finite).
- (b) P is unbounded if and only if D is nonempty (and finite).
- (c) Every $\mathbf{x} \in P$ can be represented as the sum of a convex combination of the points in V and a non-negative linear combination of the points in D , that is,

$$\mathbf{x} = \sum_{i=1}^k \alpha_i \mathbf{v}^i + \sum_{j=1}^r \beta_j \mathbf{d}^j,$$

for some $\alpha_1, \dots, \alpha_k \geq 0$ such that $\sum_{i=1}^k \alpha_i = 1$, and $\beta_1, \dots, \beta_r \geq 0$. ■

We have arrived at the important result that if there exists an optimal solution to a linear program in standard form then there exists an optimal solution among the basic feasible solutions.

Theorem 8.10 (existence and properties of optimal solutions) *Let the sets P , V , and D be defined as in Theorem 8.9 and consider the linear program*

$$\begin{aligned} &\text{minimize} && z = \mathbf{c}^T \mathbf{x}, \\ &\text{subject to} && \mathbf{x} \in P. \end{aligned} \tag{8.7}$$

(a) *This problem has a finite optimal solution if and only if P is nonempty and z is lower bounded on P , that is, if P is nonempty and $\mathbf{c}^T \mathbf{d}^j \geq 0$ for all $\mathbf{d}^j \in D$.*

(b) *If the problem has a finite optimal solution, then there exists an optimal solution among the extreme points.*

Proof. (a) Let $\mathbf{x} \in P$. Then it follows from Theorem 8.9 that

$$\mathbf{x} = \sum_{i=1}^k \alpha_i \mathbf{v}^i + \sum_{j=1}^r \beta_j \mathbf{d}^j, \tag{8.8}$$

for some $\alpha_1, \dots, \alpha_k \geq 0$ such that $\sum_{i=1}^k \alpha_i = 1$, and $\beta_1, \dots, \beta_r \geq 0$. Hence

$$\mathbf{c}^T \mathbf{x} = \sum_{i=1}^k \alpha_i \mathbf{c}^T \mathbf{v}^i + \sum_{j=1}^r \beta_j \mathbf{c}^T \mathbf{d}^j. \tag{8.9}$$

As \mathbf{x} varies over P , the value of z clearly corresponds to variations of the weights α_i and β_j . The first term in the right-hand side of (8.9) is bounded as $\sum_{i=1}^k \alpha_i = 1$. The second term is lower bounded as \mathbf{x} varies

over P if and only if $\mathbf{c}^T \mathbf{d}^j \geq 0$ holds for all $\mathbf{d}^j \in D$, since otherwise we could let $\beta_j \rightarrow +\infty$ for an index j with $\mathbf{c}^T \mathbf{d}^j < 0$, and get that $z \rightarrow -\infty$. If $\mathbf{c}^T \mathbf{d}^j \geq 0$ for all $\mathbf{d}^j \in D$, then it is clearly optimal to choose $\beta_j = 0$ for $j = 1, \dots, r$. It remains to search for the optimal solution in the convex hull of V .

(b) Assume that $\mathbf{x} \in P$ is an optimal solution and let \mathbf{x} be represented as in (8.8). From the above we have that we can choose $\beta_1 = \dots = \beta_r = 0$, so we can assume that

$$\mathbf{x} = \sum_{i=1}^k \alpha_i \mathbf{v}^i.$$

Further, let

$$a \in \arg \min_{i \in \{1, \dots, k\}} \mathbf{c}^T \mathbf{v}^i.$$

Then,

$$\mathbf{c}^T \mathbf{v}^a = \mathbf{c}^T \mathbf{v}^a \sum_{i=1}^k \alpha_i = \sum_{i=1}^k \alpha_i \mathbf{c}^T \mathbf{v}^a \leq \sum_{i=1}^k \alpha_i \mathbf{c}^T \mathbf{v}^i = \mathbf{c}^T \mathbf{x},$$

that is, the extreme point \mathbf{v}^a is a global minimum. ■

Note that part (b) of the theorem implies that if there exists an optimal solution to (8.7), then there exists an optimal solution with no more than m positive variables. This interesting fact does not hold for a general optimization problem.

Remark 8.11 The bounded case of Theorem 8.10 was already given in Theorem 4.12. ■

8.2.3 Adjacent extreme points

Consider the polytope in Figure 8.2. Clearly, every point on the line segment joining the extreme points \mathbf{x} and \mathbf{u} cannot be written as a convex combination of any pair of points that are not on this line segment. However, this is not true for the points on the line segment between the extreme points \mathbf{x} and \mathbf{w} . The extreme points \mathbf{x} and \mathbf{u} are said to be *adjacent* (while \mathbf{x} and \mathbf{w} are not adjacent).

Definition 8.12 (adjacent extreme points) *Two extreme points \mathbf{x} and \mathbf{u} of a polyhedron P are adjacent if each point \mathbf{y} on the line segment between \mathbf{x} and \mathbf{u} has the property that if*

$$\mathbf{y} = \lambda \mathbf{v} + (1 - \lambda) \mathbf{w},$$

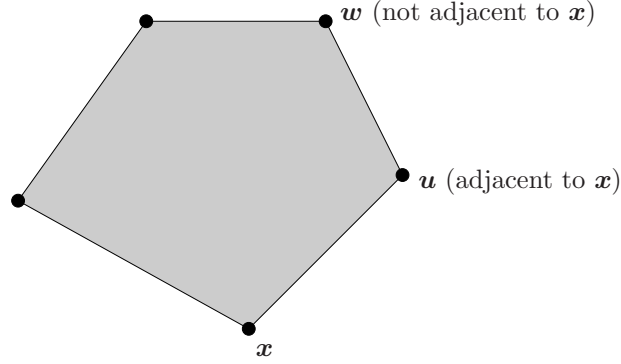


Figure 8.2: Illustration of adjacent extreme points.

where $\lambda \in (0, 1)$ and $v, w \in P$, then both v and w must be on the line segment between x and u . ■

Now, consider the polyhedron in standard form,

$$P := \{x \in \mathbb{R}^n \mid Ax = b; \quad x \geq 0^n\}.$$

Let $u \in P$ be a basic feasible solution (and hence an extreme point of P) corresponding to the partition $A = (B^1, N^1)$, where $\text{rank } B^1 = m$, that is,

$$u = \begin{pmatrix} u_{B^1} \\ u_{N^1} \end{pmatrix} = \begin{pmatrix} (B^1)^{-1}b \\ 0^{n-m} \end{pmatrix}.$$

Further, let $B^1 = (b^1, \dots, b^m)$ and $N^1 = (n^1, \dots, n^{n-m})$ (that is, $b^i \in \mathbb{R}^m$, $i = 1, \dots, m$, and $n^j \in \mathbb{R}^m$, $j = 1, \dots, n - m$, are columns of A). Construct a new partition (B^2, N^2) by replacing one column of B^1 , say b^1 , with one column of N^1 , say n^1 , that is,

$$\begin{aligned} B^2 &= (n^1, b^2, \dots, b^m), \\ N^2 &= (b^1, n^2, \dots, n^{n-m}). \end{aligned}$$

Assume that the partition (B^2, N^2) corresponds to a basic feasible solution v (i.e., v is an extreme point), and that the two extreme points u and v corresponding to (B^1, N^1) and (B^2, N^2) , respectively, are not equal. Then we have the following elegant result.

Proposition 8.13 (algebraic characterization of adjacency) *Let u and v be the extreme points that correspond to the partitions (B^1, N^1) and (B^2, N^2) described above. Then u and v are adjacent BFSs.*

Proof. If the variables of \mathbf{v} are ordered in the same way as the variables of \mathbf{u} , then the vectors must be of the form

$$\begin{aligned}\mathbf{u} &= (u_1, \dots, u_m, 0, 0, \dots, 0)^T, \\ \mathbf{v} &= (0, v_2, \dots, v_{m+1}, 0, \dots, 0)^T.\end{aligned}$$

Take a point \mathbf{x} on the line segment between \mathbf{u} and \mathbf{v} , that is,

$$\mathbf{x} = \lambda \mathbf{u} + (1 - \lambda) \mathbf{v}$$

for some $\lambda \in (0, 1)$. In order to prove the theorem we must show that if \mathbf{x} can be written as a convex combination of two feasible points, then these points must be on the line segment between \mathbf{u} and \mathbf{v} . So assume that

$$\mathbf{x} = \alpha \mathbf{y}^1 + (1 - \alpha) \mathbf{y}^2$$

for some feasible points \mathbf{y}^1 and \mathbf{y}^2 , and $\alpha \in (0, 1)$. Then it follows that \mathbf{y}^1 and \mathbf{y}^2 must be solutions to the system

$$\begin{aligned}y_1 \mathbf{b}^1 + \dots + y_m \mathbf{b}^m + y_{m+1} \mathbf{n}^1 &= \mathbf{b}, \\ y_{m+2} &= \dots = y_n = 0, \\ \mathbf{y} &\geq \mathbf{0}^n,\end{aligned}$$

or, equivalently [by multiplying both sides of the first row by $(\mathbf{B}^1)^{-1}$],

$$\begin{aligned}\mathbf{y} &= \begin{pmatrix} (\mathbf{B}^1)^{-1} \mathbf{b} \\ \mathbf{0}^{n-m} \end{pmatrix} + \begin{pmatrix} -y_{m+1} (\mathbf{B}^1)^{-1} \mathbf{n}^1 \\ y_{m+1} \\ \mathbf{0}^{n-m-1} \end{pmatrix}, \\ \mathbf{y} &\geq \mathbf{0}^n.\end{aligned}$$

But this is in fact the line segment between \mathbf{u} and \mathbf{v} (if $y_{m+1} = 0$ then $\mathbf{y} = \mathbf{u}$ and if $y_{m+1} = v_{m+1}$ then $\mathbf{y} = \mathbf{v}$). In other words, \mathbf{y}^1 and \mathbf{y}^2 are on the line segment between \mathbf{u} and \mathbf{v} , and we are done. ■

Since the simplex method at each iteration performs exactly the above replacement action the proposition actually shows that the simplex method at each non-degenerate iteration moves from one extreme point to an adjacent.

Remark 8.14 Actually a converse of the implication in Proposition 8.13 also holds. Namely, if two extreme points \mathbf{u} and \mathbf{v} are adjacent, then there exists a partition $(\mathbf{B}^1, \mathbf{N}^1)$ corresponding to \mathbf{u} and a partition $(\mathbf{B}^2, \mathbf{N}^2)$ corresponding to \mathbf{v} such that the columns of \mathbf{B}^1 and \mathbf{B}^2 are the same except for one. The proof is similar to that of Proposition 8.13. ■

8.3 Notes and further reading

The material in this chapter can be found in most books on linear programming, such as [Dan63, Chv83, Mur83, Sch86, DaT97, Pad99, Van01, DaT03, DaM05].

8.4 Exercises

Exercise 8.1 (LP modelling) Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$. Formulate the following problems as linear programming problems.

- (a) minimize $\sum_{i=1}^m |(\mathbf{A}\mathbf{x} - \mathbf{b})_i|$ subject to $\max_{i=1, \dots, n} |x_i| \leq 1$.
- (b) minimize $\sum_{i=1}^m |(\mathbf{A}\mathbf{x} - \mathbf{b})_i| + \max_{i=1, \dots, n} |x_i|$.

Exercise 8.2 (LP modelling) Consider the sets $V := \{\mathbf{v}^1, \dots, \mathbf{v}^k\} \subset \mathbb{R}^n$ and $W := \{\mathbf{w}^1, \dots, \mathbf{w}^l\} \subset \mathbb{R}^n$. Formulate the following problems as linear programming problems.

- (a) Construct, if possible, a hyperplane that separates the sets V and W , that is, find $\mathbf{a} \in \mathbb{R}^n$, with $\mathbf{a} \neq \mathbf{0}^n$, and $b \in \mathbb{R}$ such that

$$\begin{aligned} \mathbf{a}^T \mathbf{v} &\leq b, & \text{for all } \mathbf{v} \in V, \\ \mathbf{a}^T \mathbf{w} &\geq b, & \text{for all } \mathbf{w} \in W. \end{aligned}$$

- (b) Construct, if possible, a sphere that separates the sets V and W , that is, find a center $\mathbf{x}^c \in \mathbb{R}^n$ and a radius $R \geq 0$ such that

$$\begin{aligned} \|\mathbf{v} - \mathbf{x}^c\|_2 &\leq R, & \text{for all } \mathbf{v} \in V, \\ \|\mathbf{w} - \mathbf{x}^c\|_2 &\geq R, & \text{for all } \mathbf{w} \in W. \end{aligned}$$

Exercise 8.3 (linear-fractional programming) Consider the linear-fractional program

$$\begin{aligned} \text{minimize} \quad & f(\mathbf{x}) := (\mathbf{c}^T \mathbf{x} + \alpha) / (\mathbf{d}^T \mathbf{x} + \beta), \\ \text{subject to} \quad & \mathbf{A}\mathbf{x} \leq \mathbf{b}, \end{aligned} \tag{8.10}$$

where $\mathbf{c}, \mathbf{d} \in \mathbb{R}^n$, $\mathbf{A} \in \mathbb{R}^{m \times n}$, and $\mathbf{b} \in \mathbb{R}^m$. Further, assume that the polyhedron $P := \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{x} \leq \mathbf{b}\}$ is bounded and that $\mathbf{d}^T \mathbf{x} + \beta > 0$ for all $\mathbf{x} \in P$. Show that (8.10) can be solved by solving the linear program

$$\begin{aligned} \text{minimize} \quad & g(\mathbf{y}, z) := \mathbf{c}^T \mathbf{y} + \alpha z, \\ \text{subject to} \quad & \mathbf{A}\mathbf{y} - z\mathbf{b} \leq \mathbf{0}^m, \\ & \mathbf{d}^T \mathbf{y} + \beta z = 1, \\ & z \geq 0. \end{aligned} \tag{8.11}$$

[Hint: Suppose that \mathbf{y}^* together with z^* are a solution to (8.11), and show that $z^* > 0$ and that \mathbf{y}^*/z^* is a solution to (8.10).]

Linear programming models

Exercise 8.4 (standard form) Transform the linear program

$$\begin{aligned} \text{minimize } z &= x_1 - 5x_2 - 7x_3, \\ \text{subject to } &5x_1 - 2x_2 + 6x_3 \geq 5, \\ &3x_1 + 4x_2 - 9x_3 = 3, \\ &7x_1 + 3x_2 + 5x_3 \leq 9, \\ &x_1 \geq -2, \end{aligned}$$

into standard form.

Exercise 8.5 (standard form) Consider the linear program

$$\begin{aligned} \text{minimize } z &= 5x_1 + 3x_2 - 7x_3, \\ \text{subject to } &2x_1 + 4x_2 + 6x_3 = 11, \\ &3x_1 - 5x_2 + 3x_3 + x_4 = 11, \\ &x_1, \quad x_2, \quad x_4 \geq 0. \end{aligned}$$

(a) Show how to transform this problem into standard form by eliminating one constraint and the unrestricted variable x_3 .

(b) Why cannot this technique be used to eliminate variables with non-negativity restrictions?

Exercise 8.6 (basic feasible solutions) Suppose that a linear program includes a free variable x_j . When transforming this problem into standard form, x_j is replaced by

$$\begin{aligned} x_j &= x_j^+ - x_j^-, \\ x_j^+, x_j^- &\geq 0. \end{aligned}$$

Show that no basic feasible solution can include both x_j^+ and x_j^- as non-zero basic variables.

Exercise 8.7 (equivalent systems) Consider the system of equations

$$\sum_{j=1}^n a_{ij}x_j = b_i, \quad i = 1, \dots, m. \quad (8.12)$$

Show that this system is equivalent to the system

$$\sum_{j=1}^n a_{ij}x_j \leq b_i, \quad i = 1, \dots, m, \quad (8.13a)$$

$$\sum_{i=1}^m \sum_{j=1}^n a_{ij}x_j \geq \sum_{i=1}^m b_i. \quad (8.13b)$$

The simplex method

IX

This chapter presents the simplex method for solving linear programs. In Section 9.1 the algorithm is presented. First, we assume that a basic feasible solution is known at the start of the algorithm, and then we describe what to do when a BFS is not known from the beginning. In Section 9.2 we discuss termination characteristics of the algorithm. It turns out that if all the BFSs of the problem are non-degenerate, then the algorithm terminates. However, if there exist degenerate BFSs then there is a possibility that the algorithm cycles between degenerate BFSs and hence never terminates. Fortunately, the simple *Bland's rule* eliminates cycling, and which we describe. We close the chapter by discussing the computational complexity of the simplex algorithm. In the worst case, the algorithm visits all the extreme points of the problem, and since the number of extreme points may be exponential in the dimension of the problem, the simplex algorithm does not belong to the desirable polynomial complexity class. The simplex method is therefore not theoretically satisfactory, but in practice it works very well and thus it frequently appears in commercial linear programming codes.

9.1 The algorithm

Assume that we have a linear program in standard form:

$$\begin{aligned} \text{minimize} \quad & z = \mathbf{c}^T \mathbf{x}, \\ \text{subject to} \quad & \mathbf{A}\mathbf{x} = \mathbf{b}, \\ & \mathbf{x} \geq \mathbf{0}^n, \end{aligned}$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$, $n > m$, $\text{rank } \mathbf{A} = m$, $\mathbf{b} \geq \mathbf{0}^m$, and $\mathbf{c} \in \mathbb{R}^n$. (This is not a restriction, as was shown in Section 8.2.1.) At each iteration

The simplex method

the simplex algorithm starts at the current basic feasible solution (BFS) and moves to an adjacent BFS such that the objective function value decreases. It terminates with an optimal BFS (if there exists a finite optimal solution), or a *direction of unboundedness*, that is, a point in $C := \{\mathbf{p} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{p} = \mathbf{0}^m; \mathbf{p} \geq \mathbf{0}^n\}$ along which the objective function diverges to $-\infty$. (Observe that if $\mathbf{p} \in C$ is a direction of unboundedness and $\tilde{\mathbf{x}}$ is a feasible solution, then every solution $\mathbf{y}(\alpha)$ of the form

$$\mathbf{y}(\alpha) := \tilde{\mathbf{x}} + \alpha\mathbf{p}, \quad \alpha \geq 0,$$

is feasible. Hence if $\mathbf{c}^T\mathbf{p} < 0$ then $z = \mathbf{c}^T\mathbf{y}(\alpha) \rightarrow -\infty$ as $\alpha \rightarrow +\infty$.)

9.1.1 A BFS is known

Assume that a basic feasible solution $\mathbf{x} = (\mathbf{x}_B^T, \mathbf{x}_N^T)^T$ corresponding to the partition $\mathbf{A} = (\mathbf{B}, \mathbf{N})$ is known. Then we have that

$$\mathbf{A}\mathbf{x} = (\mathbf{B}, \mathbf{N}) \begin{pmatrix} \mathbf{x}_B \\ \mathbf{x}_N \end{pmatrix} = \mathbf{B}\mathbf{x}_B + \mathbf{N}\mathbf{x}_N = \mathbf{b},$$

or, equivalently,

$$\mathbf{x}_B = \mathbf{B}^{-1}\mathbf{b} - \mathbf{B}^{-1}\mathbf{N}\mathbf{x}_N. \quad (9.1)$$

Further, rearrange the components of \mathbf{c} such that $\mathbf{c} = (\mathbf{c}_B^T, \mathbf{c}_N^T)^T$ has the same ordering as $\mathbf{x} = (\mathbf{x}_B^T, \mathbf{x}_N^T)^T$. Then from (9.1) it follows that

$$\begin{aligned} \mathbf{c}^T\mathbf{x} &= \mathbf{c}_B^T\mathbf{x}_B + \mathbf{c}_N^T\mathbf{x}_N \\ &= \mathbf{c}_B^T(\mathbf{B}^{-1}\mathbf{b} - \mathbf{B}^{-1}\mathbf{N}\mathbf{x}_N) + \mathbf{c}_N^T\mathbf{x}_N \\ &= \mathbf{c}_B^T\mathbf{B}^{-1}\mathbf{b} + (\mathbf{c}_N^T - \mathbf{c}_B^T\mathbf{B}^{-1}\mathbf{N})\mathbf{x}_N. \end{aligned} \quad (9.2)$$

The principle of the simplex algorithm is now easy to describe. At the start of an iteration we are located at the BFS given by

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_B \\ \mathbf{x}_N \end{pmatrix} = \begin{pmatrix} \mathbf{B}^{-1}\mathbf{b} \\ \mathbf{0}^{n-m} \end{pmatrix},$$

which is an extreme point according to Theorem 8.7. Proposition 8.13 implies that if we construct a new partition by replacing one column of \mathbf{B} by one column of \mathbf{N} such that the new partition corresponds to a basic feasible solution, $\tilde{\mathbf{x}}$, not equal to \mathbf{x} , then $\tilde{\mathbf{x}}$ is adjacent to \mathbf{x} . The principle of the simplex algorithm is to move to an adjacent extreme point such that the objective function value decreases. From (9.2) follows that if

we increase the j^{th} component of the non-basic vector \mathbf{x}_N from 0 to 1, then the change in the objective function value becomes

$$(\tilde{\mathbf{c}}_N)_j := (\mathbf{c}_N^T - \mathbf{c}_B^T \mathbf{B}^{-1} \mathbf{N})_j,$$

that is, the change in the objective function value resulting from a unit increase of the non-basic variable $(\mathbf{x}_N)_j$ from zero is given by the j^{th} component of the vector $\tilde{\mathbf{c}}_N^T := \mathbf{c}_N^T - \mathbf{c}_B^T \mathbf{B}^{-1} \mathbf{N}$.

We call $(\tilde{\mathbf{c}}_N)_j$ the *reduced cost* of the non-basic variable $(\mathbf{x}_N)_j$ for $j = 1, \dots, n - m$. Actually, we can define the reduced cost, $\tilde{\mathbf{c}} = (\tilde{\mathbf{c}}_B^T, \tilde{\mathbf{c}}_N^T)^T$, of all the variables at the given BFS by

$$\begin{aligned} \tilde{\mathbf{c}}^T &:= \mathbf{c}^T - \mathbf{c}_B^T \mathbf{B}^{-1} \mathbf{A} = (\mathbf{c}_B^T, \mathbf{c}_N^T) - \mathbf{c}_B^T \mathbf{B}^{-1} (\mathbf{B}, \mathbf{N}) \\ &= ((\mathbf{0}^m)^T, \mathbf{c}_N^T - \mathbf{c}_B^T \mathbf{B}^{-1} \mathbf{N}); \end{aligned}$$

note that the reduced costs of the basic variables are $\tilde{\mathbf{c}}_B = \mathbf{0}^m$.

If $(\tilde{\mathbf{c}}_N)_j \geq 0$ for all $j = 1, \dots, n - m$, then there exists no adjacent extreme point such that the objective function value decreases and we stop; \mathbf{x} is then an optimal solution.

Proposition 9.1 (optimality in the simplex method) *Let \mathbf{x}^* be the basic feasible solution that corresponds to the partition $\mathbf{A} = (\mathbf{B}, \mathbf{N})$. If $(\tilde{\mathbf{c}}_N)_j \geq 0$ for all $j = 1, \dots, n - m$, then \mathbf{x}^* is an optimal solution.*

Proof. Since $\mathbf{c}_B^T \mathbf{B}^{-1} \mathbf{b}$ is constant, it follows from (9.2) that the original linear program is equivalent to

$$\begin{aligned} \text{minimize } z &= \tilde{\mathbf{c}}_N^T \mathbf{x}_N \\ \text{subject to } \mathbf{x}_B + \mathbf{B}^{-1} \mathbf{N} \mathbf{x}_N &= \mathbf{B}^{-1} \mathbf{b}, \\ \mathbf{x}_B &\geq \mathbf{0}^m, \\ \mathbf{x}_N &\geq \mathbf{0}^{n-m}, \end{aligned}$$

or, equivalently [by reducing the \mathbf{x}_B variables through (9.1)],

$$\begin{aligned} \text{minimize } z &= \tilde{\mathbf{c}}_N^T \mathbf{x}_N \\ \text{subject to } \mathbf{B}^{-1} \mathbf{N} \mathbf{x}_N &\leq \mathbf{B}^{-1} \mathbf{b}, \\ \mathbf{x}_N &\geq \mathbf{0}^{n-m}. \end{aligned} \tag{9.3}$$

Since \mathbf{x}^* is a BFS it follows that $\mathbf{x}_N^* := \mathbf{0}^{n-m}$ is feasible in (9.3). But $\tilde{\mathbf{c}}_N \geq \mathbf{0}^{n-m}$ so $\mathbf{x}_N^* = \mathbf{0}^{n-m}$ is in fact optimal in (9.3). (Why?) Hence

$$\mathbf{x}^* = \begin{pmatrix} \mathbf{B}^{-1} \mathbf{b} \\ \mathbf{0}^{n-m} \end{pmatrix}$$

is an optimal solution to the original problem. ■

Remark 9.2 (optimality condition) Proposition 9.1 states that if $(\tilde{c}_N)_j \geq 0$ for all $j = 1, \dots, n - m$, then \mathbf{x}^* is an optimal extreme point. But is it true that if \mathbf{x}^* is an optimal extreme point, then $(\tilde{c}_N)_j \geq 0$ for all $j = 1, \dots, n - m$? The answer to this question is *no*: if the optimal basic feasible solution \mathbf{x}^* is degenerate, then there may exist basis representations of \mathbf{x}^* such that $(\tilde{c}_N)_j < 0$ for some j . However, it holds that if \mathbf{x}^* is an optimal extreme point, then there exists at least one basis representation of \mathbf{x}^* such that $(\tilde{c}_N)_j \geq 0$ for all $j = 1, \dots, n - m$. Proposition 9.1 can hence be strengthened: *\mathbf{x}^* is an optimal extreme point if and only if there exists a basis representation of it such that $\tilde{c}_N \geq \mathbf{0}^{n-m}$.* ■

If some of the reduced costs are negative, then we choose a non-basic variable with the most negative reduced cost to enter the basis. We must also choose one variable from \mathbf{x}_B to leave the basis. Suppose that the variable $(\mathbf{x}_N)_j$ has been chosen to enter the basis. Then, according to (9.1), when the value of $(\mathbf{x}_N)_j$ is increased from zero we will move along the half-line (or, ray)

$$\mathbf{l}(\mu) := \begin{pmatrix} \mathbf{B}^{-1}\mathbf{b} \\ \mathbf{0}^{n-m} \end{pmatrix} + \mu \begin{pmatrix} -\mathbf{B}^{-1}\mathbf{N}_j \\ \mathbf{e}_j \end{pmatrix}, \quad \mu \geq 0,$$

where \mathbf{e}_j is the j^{th} unit vector. In order to maintain feasibility we must have that $\mathbf{l}(\mu) \geq \mathbf{0}^n$. If $\mathbf{l}(\mu) \geq \mathbf{0}^n$ for all $\mu \geq 0$, then $z \rightarrow -\infty$ as $\mu \rightarrow +\infty$, that is,

$$\mathbf{p} = \begin{pmatrix} -\mathbf{B}^{-1}\mathbf{N}_j \\ \mathbf{e}_j \end{pmatrix}$$

is a direction of unboundedness and $z \rightarrow -\infty$ along the half-line $\mathbf{l}(\mu)$, $\mu \geq 0$. Observe that this occurs if and only if

$$\mathbf{B}^{-1}\mathbf{N}_j \leq \mathbf{0}^m.$$

Otherwise, the maximal value of μ in order to maintain feasibility is

$$\mu^* = \underset{i \in \{k \mid (\mathbf{B}^{-1}\mathbf{N}_j)_k > 0\}}{\text{minimum}} \frac{(\mathbf{B}^{-1}\mathbf{b})_i}{(\mathbf{B}^{-1}\mathbf{N}_j)_i}.$$

If $\mu^* > 0$ it follows that $\mathbf{l}(\mu^*)$ is an extreme point adjacent to \mathbf{x} . Actually we move to $\mathbf{l}(\mu^*)$ by choosing the outgoing variable $(\mathbf{x}_B)_i$, where

$$i \in \arg \underset{i \in \{k \mid (\mathbf{B}^{-1}\mathbf{N}_j)_k > 0\}}{\text{minimum}} \frac{(\mathbf{B}^{-1}\mathbf{b})_i}{(\mathbf{B}^{-1}\mathbf{N}_j)_i},$$

to leave the basis.

We are now ready to state the simplex algorithm.

The Simplex Algorithm:

Step 0 (initialization: BFS). Let $\mathbf{x} = (\mathbf{x}_B^T, \mathbf{x}_N^T)^T$ be a BFS corresponding to the partition $\mathbf{A} = (\mathbf{B}, \mathbf{N})$.

Step 1 (descent direction or termination: entering variable, pricing). Calculate the reduced costs of the non-basic variables:

$$(\tilde{\mathbf{c}}_N)_j := (\mathbf{c}_N^T - \mathbf{c}_B^T \mathbf{B}^{-1} \mathbf{N})_j, \quad j = 1, \dots, n - m.$$

If $(\tilde{\mathbf{c}}_N)_j \geq 0$ for all $j = 1, \dots, n - m$ then stop; \mathbf{x} is then optimal. Otherwise choose $(\mathbf{x}_N)_j$, where

$$j \in \arg \min_{j \in \{1, \dots, n-m\}} (\tilde{\mathbf{c}}_N)_j,$$

to enter the basis.

Step 2 (line search or termination: leaving variable). If

$$\mathbf{B}^{-1} \mathbf{N}_j \leq \mathbf{0}^m,$$

then the problem is unbounded, stop; $\mathbf{p} := ((-\mathbf{B}^{-1} \mathbf{N}_j)^T, \mathbf{e}_j^T)^T$ is then a direction of unboundedness. Otherwise choose $(\mathbf{x}_B)_i$, where

$$i \in \arg \min_{i \in \{k \mid (\mathbf{B}^{-1} \mathbf{N}_j)_k > 0\}} \frac{(\mathbf{B}^{-1} \mathbf{b})_i}{(\mathbf{B}^{-1} \mathbf{N}_j)_i},$$

to leave the basis.

Step 3 (update: change basis). Construct a new partition by swapping $(\mathbf{x}_B)_i$ with $(\mathbf{x}_N)_j$. Go to Step 1.

Remark 9.3 (the simplex algorithm as a feasible descent method) In the above description, we have chosen to use terms similar to those that will be used for several descent methods in nonlinear optimization that are described in Part V; see, for example, the algorithm description in Section 11.1 for unconstrained nonlinear optimization problems. The simplex method is a very special type of descent algorithm: in order to remain feasible we generate feasible descent directions \mathbf{p} (Step 1) that follow the boundary of the polyhedron; because of the fact that the objective function is linear, a line search would yield an infinite step unless a new boundary makes such a step infeasible; this is the role of Step 2. Finally, termination at an optimal solution (Step 1) is based on a special property of linear programming which allows us to decide on global optimality based only on local information (that is, the current BFS's reduced costs). (The convexity of LP is a crucial property for this principle to be valid.) More on the characterization of this optimality criterion, and its relationships to the optimality principles in the Chapters 4–6 will be discussed in the next chapter. ■

Remark 9.4 (calculating the reduced costs) When calculating the reduced costs of the non-basic variables at the pricing Step 1 of the simplex algorithm, it is appropriate to first calculate

$$\mathbf{y}^T := \mathbf{c}_B^T \mathbf{B}^{-1}$$

through the system

$$\mathbf{B}^T \mathbf{y} = \mathbf{c}_B,$$

and then calculate the reduced costs by

$$\tilde{\mathbf{c}}_N^T = \mathbf{c}_N^T - \mathbf{y}^T \mathbf{N}.$$

By this procedure we avoid the matrix-matrix multiplication $\mathbf{B}^{-1} \mathbf{N}$. ■

Remark 9.5 (alternative pricing rules) If n is very large, it can be costly to compute the reduced costs at the pricing Step 1 of the simplex algorithm. A methodology which saves computations is *partial pricing*, in which only a subset of the elements $(\tilde{\mathbf{c}}_N)_j$ is calculated.

Another problem with the standard pricing rule is that the use of the criterion to minimize $_{j \in \{1, \dots, n-m\}} \{(\tilde{\mathbf{c}}_N)_j\}$ does not take into account the actual improvement that is made. In particular, a different scaling of the variables might mean that a unit change is a dramatic move in one variable, and a very small move in another. The *steepest-edge rule* eliminates this scaling problem somewhat: With $(\mathbf{x}_N)_j$ being the entering variable we have that

$$\begin{pmatrix} \mathbf{x}_B \\ \mathbf{x}_N \end{pmatrix}^{\text{new}} := \begin{pmatrix} \mathbf{x}_B \\ \mathbf{x}_N \end{pmatrix} + (\mathbf{x}_N)_j \mathbf{p}_j, \quad \mathbf{p}_j = \begin{pmatrix} -\mathbf{B}^{-1} \mathbf{N}_j \\ \mathbf{e}_j \end{pmatrix}.$$

Choose j in

$$\arg \min_{j \in \{1, \dots, n-m\}} \frac{\mathbf{c}^T \mathbf{p}_j}{\|\mathbf{p}_j\|},$$

that is, the usual pricing rule based on $\mathbf{c}^T \mathbf{p}_j = \mathbf{c}_B^T (-\mathbf{B}^{-1} \mathbf{N}_j) + (\mathbf{c}_N)_j = (\tilde{\mathbf{c}}_N)_j$ is replaced by a rule wherein the reduced costs are scaled by the length of the candidate search directions \mathbf{p}_j . (Other scaling factors can of course be used.) ■

Remark 9.6 (initial basic feasible solution) Consider the linear program

$$\begin{aligned} \text{minimize} \quad & z = \mathbf{c}^T \mathbf{x}, \\ \text{subject to} \quad & \mathbf{A} \mathbf{x} \leq \mathbf{b}, \\ & \mathbf{x} \geq \mathbf{0}^n, \end{aligned} \tag{9.4}$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \geq \mathbf{0}^m$, and $\mathbf{c} \in \mathbb{R}^n$. By introducing slack variables $\mathbf{s} \in \mathbb{R}^m$ we get

$$\begin{aligned} \text{minimize } z &= \mathbf{c}^T \mathbf{x}, \\ \text{subject to } \quad \mathbf{A}\mathbf{x} + \mathbf{I}^m \mathbf{s} &= \mathbf{b}, \\ \mathbf{x} &\geq \mathbf{0}^n, \\ \mathbf{s} &\geq \mathbf{0}^m. \end{aligned} \tag{9.5}$$

Since $\mathbf{b} \geq \mathbf{0}^m$ it then follows that the partition $(\mathbf{I}^m, \mathbf{A})$ corresponds to a basic feasible solution to (9.5), that is, the slack variables \mathbf{s} are the basic variables. (This corresponds to the origin in the problem (9.4), which is clearly feasible when $\mathbf{b} \geq \mathbf{0}^m$.)

Similarly, if we can identify an identity matrix among the columns of the constraint matrix, then (if the right-hand side is non-negative, which is the case if the problem is in standard form) we obtain a BFS by taking the variables that correspond to these columns as basic variables. ■

Example 9.7 (the simplex method) Consider the linear program

$$\begin{aligned} \text{minimize } z &= x_1 - 2x_2 - 4x_3 + 4x_4, \\ \text{subject to } \quad & -x_2 + 2x_3 + x_4 \leq 4, \\ & -2x_1 + x_2 + x_3 - 4x_4 \leq 5, \\ & x_1 - x_2 + 2x_4 \leq 3, \\ & x_1, x_2, x_3, x_4 \geq 0. \end{aligned}$$

By introducing the slack variables x_5, x_6 and x_7 we get the problem to

$$\begin{aligned} \text{minimize } z &= x_1 - 2x_2 - 4x_3 + 4x_4, \\ \text{subject to } \quad & -x_2 + 2x_3 + x_4 + x_5 = 4, \\ & -2x_1 + x_2 + x_3 - 4x_4 + x_6 = 5, \\ & x_1 - x_2 + 2x_4 + x_7 = 3, \\ & x_1, x_2, x_3, x_4, x_5, x_6, x_7 \geq 0. \end{aligned}$$

According to Remark 9.6 we can take $\mathbf{x}_B = (x_5, x_6, x_7)^T$ and $\mathbf{x}_N = (x_1, x_2, x_3, x_4)^T$ as the initial basic and non-basic vector, respectively. The reduced costs of the non-basic variables then become

$$\mathbf{c}_N^T - \mathbf{c}_B^T \mathbf{B}^{-1} \mathbf{N} = (1, -2, -4, 4),$$

and hence we choose x_3 as the entering variable. Further, from

$$\begin{aligned} \mathbf{B}^{-1} \mathbf{b} &= (4, 5, 3)^T, \\ \mathbf{B}^{-1} \mathbf{N}_3 &= (2, 1, 0)^T, \end{aligned}$$

The simplex method

$$\arg \min_{i \in \{k \mid (\mathbf{B}^{-1}\mathbf{N}_3)_k > 0\}} \frac{(\mathbf{B}^{-1}\mathbf{b})_i}{(\mathbf{B}^{-1}\mathbf{N}_3)_i} = \{1\},$$

so we choose x_5 to leave the basis. The new basic and non-basic vectors are $\mathbf{x}_B = (x_3, x_6, x_7)^T$ and $\mathbf{x}_N = (x_1, x_2, x_5, x_4)^T$, and the reduced costs of the non-basic variables become

$$\mathbf{c}_N^T - \mathbf{c}_B^T \mathbf{B}^{-1} \mathbf{N} = (1, -4, 2, 6),$$

so x_2 is the entering variable, and from

$$\begin{aligned} \mathbf{B}^{-1}\mathbf{b} &= (2, 3, 3)^T, \\ \mathbf{B}^{-1}\mathbf{N}_2 &= (-1/2, 3/2, -1)^T, \end{aligned}$$

$$\arg \min_{i \in \{k \mid (\mathbf{B}^{-1}\mathbf{N}_2)_k > 0\}} \frac{(\mathbf{B}^{-1}\mathbf{b})_i}{(\mathbf{B}^{-1}\mathbf{N}_2)_i} = \{2\},$$

and hence x_6 is the leaving variable. The new basic and non-basic vectors become $\mathbf{x}_B = (x_3, x_2, x_7)^T$ and $\mathbf{x}_N = (x_1, x_6, x_5, x_4)^T$, and the reduced costs of the non-basic variables are

$$\mathbf{c}_N^T - \mathbf{c}_B^T \mathbf{B}^{-1} \mathbf{N} = (-13/3, 8/3, 2/3, -6),$$

so x_4 is the entering variable and

$$\begin{aligned} \mathbf{B}^{-1}\mathbf{b} &= (3, 2, 5)^T, \\ \mathbf{B}^{-1}\mathbf{N}_4 &= (-1, -3, -1)^T. \end{aligned}$$

But since $\mathbf{B}^{-1}\mathbf{N}_4 \leq \mathbf{0}^3$ it follows that the objective function diverges to $-\infty$ along the half-line given by

$$\mathbf{l}(\mu) = (x_1, x_2, x_3, x_4)^T = (0, 2, 3, 0)^T + \mu(0, 3, 1, 1)^T, \quad \mu \geq 0.$$

We conclude that the problem is unbounded. ■

9.1.2 A BFS is not known: phase I & II

Often a basic feasible solution is not known initially. (In fact, only if the origin is feasible in (9.4) we know a BFS immediately.) However, an initial basic feasible solution can be found by solving a linear program that is a pure feasibility problem. We call this the *phase I problem*.

Consider the following linear program in standard form:

$$\begin{aligned} \text{minimize} \quad & z = \mathbf{c}^T \mathbf{x}, \\ \text{subject to} \quad & \mathbf{A}\mathbf{x} = \mathbf{b}, \\ & \mathbf{x} \geq \mathbf{0}^n. \end{aligned} \tag{9.6}$$

In order to find a basic feasible solution we introduce the *artificial variables* $\mathbf{a} \in \mathbb{R}^m$ and consider the *phase I problem* to

$$\begin{aligned} \text{minimize } w &= (\mathbf{1}^m)^T \mathbf{a}, \\ \text{subject to } \quad \mathbf{A}\mathbf{x} + \mathbf{I}^m \mathbf{a} &= \mathbf{b}, \\ \mathbf{x} &\geq \mathbf{0}^n, \\ \mathbf{a} &\geq \mathbf{0}^m. \end{aligned} \tag{9.7}$$

In other words, we introduce an additional (artificial) variable a_i for every linear constraint $i = 1, \dots, m$, and thus construct the unit matrix in $\mathbb{R}^{m \times m}$ sought.

We obtain a BFS to the phase I problem (9.7) by taking the artificial variables \mathbf{a} as the basic variables. (Remember that $\mathbf{b} \geq \mathbf{0}^m$; the simplicity of finding an initial BFS for the phase I problem is in fact the reason why we require this to hold!) Then the phase I problem (9.7) can be solved by the simplex method stated in the previous section. Note that the phase I problem is bounded from below $[(\mathbf{1}^m)^T \mathbf{a} \geq 0]$ which means that an optimal solution to (9.7) always exists by Theorem 8.10.

Assume that the optimal objective function value is w^* . We observe that if and only if the part \mathbf{x}^* of an optimal solution $((\mathbf{x}^*)^T, (\mathbf{a}^*)^T)^T$ to the problem (9.7) is a feasible solution to the original problem (9.6), then $((\mathbf{x}^*)^T, (\mathbf{0}^m)^T)^T$ is an optimal solution to the phase I problem and $w^* = 0$. Hence, if $w^* > 0$, then the original linear program is infeasible. We have the following cases:

1. If $w^* > 0$, then the original problem is infeasible.
2. If $w^* = 0$, then if the optimal basic feasible solution is $(\mathbf{x}^T, \mathbf{a}^T)^T$ we must have that $\mathbf{a} = \mathbf{0}^m$, and \mathbf{x} corresponds to a basic feasible solution to the original problem.¹

Therefore, if there exists a feasible solution to the original problem (9.6), then a BFS is found by solving the phase I problem (9.7). This BFS can then be used as the starting BFS when solving the original problem, which is called the *phase II problem*, with the simplex method.

Remark 9.8 (artificial variables) The purpose of introducing artificial variables is to get an identity matrix among the columns of the constraint matrix. If some of the columns of the constraint matrix of the original problem consists of only zeros except for one positive entry, then it is not

¹Notice that if the final BFS in the phase I problem is degenerate then one or several artificial variables a_i may remain in the basis with value zero; in order to remove them from the basis a number of degenerate pivots may have to be performed; this is naturally always possible.

The simplex method

necessary to introduce an artificial variable in the corresponding row. An example of a linear constraint for which an original variable naturally serves as a basic variable is a \leq -constraint with a positive right-hand side, in which case we can use the corresponding slack variable. ■

Example 9.9 (phase I & II) Consider the following linear program:

$$\begin{array}{llllll} \text{minimize} & z = & 2x_1, & & & \\ \text{subject to} & & x_1 & & -x_3 & = 3, \\ & & x_1 & -x_2 & & -2x_4 = 1, \\ & & 2x_1 & & & +x_4 \leq 7, \\ & & x_1, & x_2, & x_3, & x_4 \geq 0. \end{array}$$

By introducing a slack variable x_5 we get the equivalent linear program in standard form:

$$\begin{array}{llllllll} \text{minimize} & z = & 2x_1, & & & & & (9.8) \\ \text{subject to} & & x_1 & & -x_3 & & & = 3, \\ & & x_1 & -x_2 & & -2x_4 & & = 1, \\ & & 2x_1 & & & +x_4 & +x_5 & = 7, \\ & & x_1, & x_2, & x_3, & x_4, & x_5 & \geq 0. \end{array}$$

We cannot identify the identity matrix among the columns of the constraint matrix of the problem (9.8), but the third unit vector \mathbf{e}_3 is found in the column corresponding to the x_5 -variable. Therefore, we leave the problem (9.8) for a while, and instead introduce two artificial variables a_1 and a_2 and consider the phase I problem to

$$\begin{array}{llllllllll} \text{minimize} & w = & & & & & a_1 & +a_2 & & \\ \text{subject to} & & x_1 & & -x_3 & & +a_1 & & & = 3, \\ & & x_1 & -x_2 & & -2x_4 & & +a_2 & = 1, \\ & & 2x_1 & & & +x_4 & +x_5 & & = 7, \\ & & x_1, & x_2, & x_3, & x_4, & x_5, & a_1, & a_2 & \geq 0. \end{array}$$

Let $\mathbf{x}_B = (a_1, a_2, x_5)^T$ and $\mathbf{x}_N = (x_1, x_2, x_3, x_4)^T$ be the initial basic and non-basic vector, respectively. The reduced costs of the non-basic variables then become

$$\mathbf{c}_N^T - \mathbf{c}_B^T \mathbf{B}^{-1} \mathbf{N} = (-2, 1, 1, 2),$$

and hence we choose x_1 as the entering variable. Further, from

$$\begin{aligned} \mathbf{B}^{-1} \mathbf{b} &= (3, 1, 7)^T, \\ \mathbf{B}^{-1} \mathbf{N}_1 &= (1, 1, 2)^T, \end{aligned}$$

$$\arg \min_{i \in \{k \mid (B^{-1}N_1)_k > 0\}} \frac{(B^{-1}b)_i}{(B^{-1}N_1)_i} = \{2\},$$

so we choose a_2 as the leaving variable. The new basic and non-basic vectors are $\mathbf{x}_B = (a_1, x_1, x_5)^T$ and $\mathbf{x}_N = (a_2, x_2, x_3, x_4)^T$, and the reduced costs of the non-basic variables become

$$\mathbf{c}_N^T - \mathbf{c}_B^T B^{-1} \mathbf{N} = (2, -1, 1, -2),$$

so x_4 is the entering variable, and from

$$\begin{aligned} B^{-1}b &= (2, 1, 5)^T, \\ B^{-1}N_4 &= (2, -2, 5)^T, \end{aligned}$$

$$\arg \min_{i \in \{k \mid (B^{-1}N_4)_k > 0\}} \frac{(B^{-1}b)_i}{(B^{-1}N_4)_i} = \{1, 3\},$$

and we choose a_1 to leave the basis. The new basic and non-basic vectors become $\mathbf{x}_B = (x_4, x_1, x_5)^T$ and $\mathbf{x}_N = (a_2, x_2, x_3, a_1)^T$, and the reduced costs of the non-basic variables are

$$\mathbf{c}_N^T - \mathbf{c}_B^T B^{-1} \mathbf{N} = (1, 0, 0, 1),$$

so by choosing the basic variables as $\mathbf{x}_B = (x_4, x_1, x_5)^T$ we get an optimal basic feasible solution of the phase I problem, and $w^* = 0$. This means that by choosing the basic variables as $\mathbf{x}_B = (x_4, x_1, x_5)^T$ we get a basic feasible solution of the phase II problem (9.8).

We return to the problem (9.8). By letting $\mathbf{x}_B = (x_4, x_1, x_5)^T$ and $\mathbf{x}_N = (x_2, x_3)^T$ the reduced costs are

$$\tilde{\mathbf{c}}_N^T = \mathbf{c}_N^T - \mathbf{c}_B^T B^{-1} \mathbf{N} = (0, 2),$$

which means that

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_B \\ \mathbf{x}_N \end{pmatrix} = \begin{pmatrix} x_4 \\ x_1 \\ x_5 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} B^{-1}b \\ \mathbf{0}^2 \end{pmatrix} = \begin{pmatrix} 1 \\ 3 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

is an optimal basic feasible solution to the original problem. (Observe that the BFS found when solving the phase I problem typically is not an optimal solution to the phase II problem!) But since the reduced cost of x_2 is zero there is a possibility that there are alternative optimal solutions. Let x_2 enter the basic vector. From

$$\begin{aligned} B^{-1}b &= (1, 3, 0)^T, \\ B^{-1}N_1 &= (0.5, 0, -0.5)^T, \end{aligned}$$

The simplex method

$$\arg \min_{i \in \{k \mid (\mathbf{B}^{-1}\mathbf{N}_1)_k > 0\}} \frac{(\mathbf{B}^{-1}\mathbf{b})_i}{(\mathbf{B}^{-1}\mathbf{N}_1)_i} = \{1\},$$

so x_4 is the leaving variable. We get $\mathbf{x}_B = (x_2, x_1, x_5)^T$ and $\mathbf{x}_N = (x_4, x_3)^T$, and since the reduced costs become

$$\tilde{\mathbf{c}}_N^T = \mathbf{c}_N^T - \mathbf{c}_B^T \mathbf{B}^{-1} \mathbf{N} = (0, 2),$$

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_B \\ \mathbf{x}_N \end{pmatrix} = \begin{pmatrix} x_2 \\ x_1 \\ x_5 \\ x_4 \\ x_3 \end{pmatrix} = \begin{pmatrix} \mathbf{B}^{-1}\mathbf{b} \\ \mathbf{0}^2 \end{pmatrix} = \begin{pmatrix} 2 \\ 3 \\ 1 \\ 0 \\ 0 \end{pmatrix}$$

is an alternative optimal basic feasible solution. ■

9.1.3 Alternative optimal solutions

As we saw in Example 9.9 there can be alternative optimal solutions to a linear program. However, this can only happen if some of the reduced costs of the non-basic variables of an optimal solution is zero.

Proposition 9.10 (unique optimal solutions in linear programming) *Consider the linear program in standard form*

$$\begin{aligned} \text{minimize} \quad & z = \mathbf{c}^T \mathbf{x}, \\ \text{subject to} \quad & \mathbf{A}\mathbf{x} = \mathbf{b}, \\ & \mathbf{x} \geq \mathbf{0}^n. \end{aligned}$$

Let $\mathbf{x} = (\mathbf{x}_B^T, \mathbf{x}_N^T)^T$ be an optimal basic feasible solution that corresponds to the partition $\mathbf{A} = (\mathbf{B}, \mathbf{N})$. If the reduced costs of the non-basic variables \mathbf{x}_N are all strictly positive, then \mathbf{x} is the unique optimal solution.

Proof. As in the proof of Proposition 9.1 we have that the original linear program is equivalent to

$$\begin{aligned} \text{minimize} \quad & z = \tilde{\mathbf{c}}_N^T \mathbf{x}_N \\ \text{subject to} \quad & \mathbf{x}_B + \mathbf{B}^{-1} \mathbf{N} \mathbf{x}_N = \mathbf{B}^{-1} \mathbf{b}, \\ & \mathbf{x}_B \geq \mathbf{0}^m, \\ & \mathbf{x}_N \geq \mathbf{0}^{n-m}. \end{aligned}$$

Now if the reduced costs of the non-basic variables are all strictly positive, that is, $\tilde{\mathbf{c}}_N > \mathbf{0}^{n-m}$, it follows that a solution for which $(\mathbf{x}_N)_j > 0$ for some $j = 1, \dots, n - m$ cannot be optimal. Hence

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_B \\ \mathbf{x}_N \end{pmatrix} = \begin{pmatrix} \mathbf{B}^{-1}\mathbf{b} \\ \mathbf{0}^{n-m} \end{pmatrix}$$

is the unique optimal solution. ■

9.2 Termination

So far we have not discussed whether the simplex algorithm terminates in a finite number of iterations. Unfortunately, if there exist degenerate BFSs it can happen that the simplex algorithm cycles between degenerate solutions and hence never terminates. However, if all of the BFSs are non-degenerate this kind of cycling never occurs.

Theorem 9.11 (finiteness of the simplex algorithm) *If all of the basic feasible solutions are non-degenerate, then the simplex algorithm terminates after a finite number of iterations.*

Proof. If a basic feasible solution is non-degenerate then it follows that it has exactly m positive components, and hence has a unique associated basis. In this case, in the minimum ratio test,

$$\mu^* = \min_{i \in \{k \mid (\mathbf{B}^{-1}\mathbf{N}_j)_k > 0\}} \frac{(\mathbf{B}^{-1}\mathbf{b})_i}{(\mathbf{B}^{-1}\mathbf{N}_j)_i} > 0.$$

Therefore, at each iteration the objective value decreases, and hence a basic feasible solution that has appeared once can never reappear. Further, from Corollary 3.18 follows that the number of extreme points, hence the number of basic feasible solutions, is finite. We are done. ■

Cycling resulting from degeneracy does not seem to occur often among the numerous degenerate linear programs encountered in practical applications. However, the fact that it can occur is not theoretically satisfactory. Therefore, methods have been developed that avoid cycling. One of them is *Bland's rule*.

Theorem 9.12 (Bland's rule) *Fix an ordering of the variables. (This ordering can be arbitrary, but once it has been selected it cannot be changed.) If at each iteration step the entering and leaving variables are*

chosen as the first variables that are eligible² in the ordering, then the simplex algorithm terminates after a finite number of iteration steps. ■

9.3 Computational complexity

The simplex algorithm is very efficient in practice. Although the total number of basic feasible solutions can be as many as

$$\binom{n}{m} = \frac{n!}{(n-m)!m!}$$

(the number of different ways m objects can be chosen from n objects), which grows exponentially, it is rare that more than $3m$ iterations are needed, and in practice the expected number is in the order of $3m/2$. Since each iteration costs no more than a polynomial ($O(m^3)$ for factorizations and $O(mn)$ for the pricing) the algorithm is polynomial in practice. Its worst-case behaviour is however bad, in fact exponential.

The bad worst-case behaviour of the simplex method led to a huge amount of work being laid down to find polynomial algorithms for solving linear programs. Such a polynomial time competitor to the simplex method nowadays is the class of *interior point algorithms*. Its main feature is that the optimal extreme points are not approached by following the edges, but by moving within the interior of the polyhedron. The famous Karmarkar algorithm is one, which however has been improved much in recent years. An analysis of interior point methods for linear programs is made in Chapter 13, as they are in fact to be seen as instances of the interior penalty algorithm in nonlinear programming.

9.4 Notes and further reading

The simplex method was developed by George Dantzig [Dan51]. The version of the simplex method presented is usually called the *revised simplex method*, and was first described by Dantzig [Dan53] and Orchard-Hays [Orc54]. The first book describing the simplex method was [Dan63].

In the (revised) simplex algorithm several computations are performed using B^{-1} . The major drawback in this approach is that roundoff

²By *eligible entering variables* we mean the variables $(x_N)_j$ for which $(\bar{c}_N)_j < 0$, and when we have chosen the entering variable j , the *eligible leaving variables* are the variables $(x_B)_i$ such that

$$i \in \arg \min_{i \in \{k \mid (B^{-1}N_j)_k > 0\}} \frac{(B^{-1}b)_i}{(B^{-1}N_j)_i}.$$

errors accumulate as the algorithm proceeds. This drawback can however be alleviated by using stable forms of LU decomposition or Cholesky factorization. Most of the software packages for linear programming use LU decomposition. Early references on numerically stable forms of the simplex method are [BaG69, Bar71, Sau72, GiM73]. Books that discuss the subject are [Mur83, NaS96].

The first example of cycling of the simplex algorithm was constructed by Hoffman [Hof53]. Several methods have been developed for avoiding cycling, such as the perturbation method of Charnes [Cha52], the lexicographic method of Dantzig, Orden and Wolfe [DOW55], and Bland's rule [Bla77]. In practice, however, cycling is rarely encountered. Instead, the problem is *stalling*, which means that the value of the objective function does not change (or changes very little) for a very large number of iterations³ before it eventually starts to make substantial progress again. So in practice, we are interested in methods that primarily prevent stalling, and only secondarily avoid cycling (see, e.g., [GMSW89]).

In 1972, Klee and Minty [KlM72] showed that there exist problems of arbitrary size that cause the simplex method to examine every possible basis when the standard (steepest-descent) pricing rule is used, and hence showed that the simplex method is an exponential algorithm in the worst case. It is still an open question, however, whether there exists a rule for choosing entering and leaving basic variables that makes the simplex method polynomial. The first polynomial-time method for linear programming was given by Khachiyan [Kha79, Kha80], by adapting the ellipsoid method for nonlinear programming of Shor [Sho77] and Yudin and Nemirovskii [YuN77]. Karmarkar [Kar84a, Kar84b] showed that interior point methods can be used in order to solve linear programming problems in polynomial time.

General text books that discuss the simplex method are [Dan63, Chv83, Mur83, Sch86, DaT97, Pad99, Van01, DaT03, DaM05].

9.5 Exercises

Exercise 9.1 (checking feasibility: phase I) Consider the system

$$\begin{aligned} 3x_1 + 2x_2 - x_3 &\leq -3, \\ -x_1 - x_2 + 2x_3 &\leq -1, \\ x_1, \quad x_2, \quad x_3 &\geq 0. \end{aligned}$$

Show that this system is infeasible.

³“Very large” normally refers to a number of iterations which is an exponential function of the number of variables of the LP problem.

The simplex method

Exercise 9.2 (the simplex algorithm: phase I & II) Consider the linear program

$$\begin{aligned} \text{minimize } z &= 3x_1 + 2x_2 + x_3, \\ \text{subject to } 2x_1 &+ x_3 \geq 3, \\ 2x_1 + 2x_2 + x_3 &= 5, \\ x_1, x_2, x_3 &\geq 0. \end{aligned}$$

- (a) Solve this problem by using the simplex algorithm with phase I & II.
 (b) Is the optimal solution obtained unique?

Exercise 9.3 (the simplex algorithm) Consider the linear program

$$\begin{aligned} \text{minimize } z &= \mathbf{c}^T \mathbf{x}, \\ \text{subject to } \mathbf{A}\mathbf{x} &= \mathbf{b}, \\ \mathbf{x} &\geq \mathbf{0}^n. \end{aligned}$$

Suppose that at a given step of the simplex algorithm, there is only one possible entering variable, $(\mathbf{x}_N)_j$. Also assume that the current BFS is non-degenerate. Show that $(\mathbf{x}_N)_j > 0$ in any optimal solution.

Exercise 9.4 (cycling of the simplex algorithm) Consider the linear program

$$\begin{aligned} \text{minimize } z &= -\frac{2}{5}x_5 - \frac{2}{5}x_6 + \frac{9}{5}x_7, \\ \text{subject to } x_1 &+ \frac{3}{5}x_5 - \frac{32}{5}x_6 + \frac{24}{5}x_7 = 0, \\ x_2 &+ \frac{1}{5}x_5 - \frac{9}{5}x_6 + \frac{3}{5}x_7 = 0, \\ x_3 &+ \frac{2}{5}x_5 - \frac{8}{5}x_6 + \frac{1}{5}x_7 = 0, \\ x_4 &+ x_6 = 1, \\ x_1, x_2, x_3, x_4, x_5, x_6, x_7 &\geq 0. \end{aligned}$$

In solving this problem by the simplex algorithm starting at the BFS $\mathbf{x}_B := (x_1, x_2, x_3, x_4)^T$, in each iteration step select the entering variable as x_s , where

$$s := \text{minimum } \{j \mid \tilde{\mathbf{c}}_j < 0\},$$

and select the pivot row as the r^{th} row, where

$$r := \text{minimum } \{i \mid \text{row } i \text{ is eligible}\}.$$

Show that cycling occurs.

Linear programming duality and sensitivity analysis



10.1 Introduction

Consider the linear program

$$\begin{aligned} & \text{minimize} && z = \mathbf{c}^T \mathbf{x}, \\ & \text{subject to} && \mathbf{A} \mathbf{x} = \mathbf{b}, \\ & && \mathbf{x} \geq \mathbf{0}^n, \end{aligned} \tag{10.1}$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^m$, and $\mathbf{c} \in \mathbb{R}^n$, and assume that this problem has been solved by the simplex algorithm. Let $\mathbf{x}^* = (\mathbf{x}_B^T, \mathbf{x}_N^T)^T$ be an optimal basic feasible solution corresponding to the partition $\mathbf{A} = (\mathbf{B}, \mathbf{N})$. Introduce the vector $\mathbf{y}^* \in \mathbb{R}^m$ through

$$(\mathbf{y}^*)^T := \mathbf{c}_B^T \mathbf{B}^{-1}.$$

Since \mathbf{x}^* is an optimal solution it follows that the reduced costs of the non-basic variables are greater than or equal to zero, that is,

$$\mathbf{c}_N^T - \mathbf{c}_B^T \mathbf{B}^{-1} \mathbf{N} \geq (\mathbf{0}^{n-m})^T \iff \mathbf{c}_N^T - (\mathbf{y}^*)^T \mathbf{N} \geq (\mathbf{0}^{n-m})^T.$$

Further, $\mathbf{c}_B^T - (\mathbf{y}^*)^T \mathbf{B} = \mathbf{c}_B^T - \mathbf{c}_B^T \mathbf{B}^{-1} \mathbf{B} = (\mathbf{0}^m)^T$, so we have that

$$\mathbf{c}^T - (\mathbf{y}^*)^T \mathbf{A} \geq (\mathbf{0}^n)^T,$$

or equivalently,

$$\mathbf{A}^T \mathbf{y}^* \leq \mathbf{c}.$$

Now, for every $\mathbf{y} \in \mathbb{R}^m$ such that $\mathbf{A}^T \mathbf{y} \leq \mathbf{c}$ and every feasible solution \mathbf{x} to (10.1) it holds that

$$\mathbf{c}^T \mathbf{x} \geq \mathbf{y}^T \mathbf{A} \mathbf{x} = \mathbf{y}^T \mathbf{b} = \mathbf{b}^T \mathbf{y}.$$

But

$$\mathbf{b}^T \mathbf{y}^* = \mathbf{b}^T (\mathbf{B}^{-1})^T \mathbf{c}_B = \mathbf{c}_B^T \mathbf{B}^{-1} \mathbf{b} = \mathbf{c}_B^T \mathbf{x}_B = \mathbf{c}^T \mathbf{x}^*,$$

so in fact we have that \mathbf{y}^* is an optimal solution to the linear program

$$\begin{aligned} & \text{maximize} && \mathbf{b}^T \mathbf{y}, \\ & \text{subject to} && \mathbf{A}^T \mathbf{y} \leq \mathbf{c}, \\ & && \mathbf{y} \text{ free.} \end{aligned} \tag{10.2}$$

Observe that the linear program (10.2) is exactly the Lagrangian dual problem to (10.1) (see Section 6.2.4). Also, note that the linear programs (10.1) and (10.2) have the same optimal objective function values, which is in accordance with the Strong Duality Theorem 6.12 (see also Theorem 10.6 below for an independent proof).

The linear program (10.2) is called the *linear programming dual* to the linear program (10.1) (which is called the *primal linear program*). In this chapter we will study linear programming duality. In Section 10.2 we discuss how to construct the linear programming dual to a general linear program. Section 10.3 presents duality theory, such as that of weak and strong duality and complementary slackness. This theory specializes that of Lagrangian duality presented in Chapter 6. The dual simplex method is developed in Section 10.4. Finally, in Section 10.5 we discuss how the optimal solutions of a linear program change if the right-hand side \mathbf{b} or the objective function coefficients \mathbf{c} are modified.

10.2 The linear programming dual

For every linear program it is possible to construct the Lagrangian dual problem through the Lagrangian relaxation of the affine constraints. We will refer to this problem as the *dual linear program*. It is quite tedious to construct the Lagrangian dual problem for every special case of a linear program, but fortunately the dual of a general linear program can be constructed just by following some simple rules. These rules are presented in this section. (It is, however, a good exercise to show the validity of these rules by constructing the Lagrangian dual in each case.)

10.2.1 Canonical form

When presenting the rules for constructing the linear programming dual we will utilize the notation of *canonical form*. The canonical form is connected with the directions of the inequalities of the problem and with the objective. If the objective is to maximize the objective function, then every inequality of type “ \leq ” is said to be of canonical form. Similarly, if the objective is to minimize the objective function, then every inequality of type “ \geq ” is said to be of canonical form. Further, we consider non-negative variables to be variables in canonical form.

Remark 10.1 (mnemonic rule for canonical form) Consider the LP

$$\begin{array}{ll} \text{minimize} & z = x_1 \\ \text{subject to} & x_1 \leq 1. \end{array}$$

This problem is unbounded from below and hence an optimal solution does not exist. However, if the problem is to

$$\begin{array}{ll} \text{minimize} & z = x_1 \\ \text{subject to} & x_1 \geq 1, \end{array}$$

then an optimal solution exists, namely $x_1 = 1$. Hence it seems natural to consider inequalities of type “ \geq ” as canonical to minimization problems. Similarly, it is natural that inequalities of type “ \leq ” are canonical to maximization problems. ■

10.2.2 Constructing the dual

From the notation of canonical form introduced in Section 10.2.1 we can now construct the dual, (D), to a general linear program, (P), according to the following rules.

Dual variables

To each constraint of (P) a dual variable, y_i , is introduced. If the i^{th} constraint of (P) is an inequality of canonical form, then y_i is a non-negative variable, that is, $y_i \geq 0$. Similarly, if the i^{th} constraint of (P) is an inequality that is not of canonical form, then $y_i \leq 0$. Finally, if the i^{th} constraint of (P) is an equality, then the variable y_i is unrestricted.

Dual objective function

If (P) is a minimization (respectively, a maximization) problem, then (D) is a maximization (respectively, a minimization) problem. The objective

function coefficient for the variable y_i in the dual problem equals the right-hand side constant b_i of the i^{th} constraint of (P).

Constraints of the dual problem

If \mathbf{A} is the constraint matrix of (P), then \mathbf{A}^T is the constraint matrix of (D). The j^{th} right-hand side constant of (D) equals the j^{th} coefficient c_j in the objective function of (P).

If the j^{th} variable of (P) is non-negative, then the j^{th} constraint of (D) is an inequality of canonical form. If the j^{th} variable of (P) is non-positive, then the j^{th} constraint of (D) is an inequality of non-canonical form. Finally, if the j^{th} variable of (P) is unrestricted, then the j^{th} constraint of (D) is an equality.

Summary

The above rules can be summarized as follows:

primal/dual constraint		dual/primal variable
canonical inequality	\Longleftrightarrow	≥ 0
non-canonical inequality	\Longleftrightarrow	≤ 0
equality	\Longleftrightarrow	unrestricted

Consider the following general linear program:

$$\begin{aligned}
 &\text{minimize} && z = \sum_{j=1}^n c_j x_j, \\
 &\text{subject to} && \sum_{j=1}^n a_{ij} x_j \geq b_i, \quad i \in C, \\
 &&& \sum_{j=1}^n a_{ij} x_j \leq b_i, \quad i \in NC, \\
 &&& \sum_{j=1}^n a_{ij} x_j = b_i, \quad i \in E, \\
 &&& x_j \geq 0, \quad j \in P, \\
 &&& x_j \leq 0, \quad j \in N, \\
 &&& x_j \text{ free}, \quad j \in F,
 \end{aligned}$$

where C stands for “canonical”, NC for “non-canonical”, E for “equality”, P for “positive”, N for “negative”, and F for “free”. Note that

$P \cup N \cup F = \{1, \dots, n\}$ and $C \cup NC \cup E = \{1, \dots, m\}$. If we apply the rules above we get the following dual linear program:

$$\begin{aligned}
 &\text{maximize } w = \sum_{i=1}^m b_i y_i, \\
 &\text{subject to } \sum_{i=1}^m a_{ij} y_i \leq c_j, \quad j \in P, \\
 &\quad \sum_{i=1}^m a_{ij} y_i \geq c_j, \quad j \in N, \\
 &\quad \sum_{i=1}^m a_{ij} y_i = c_j, \quad j \in F, \\
 &\quad y_i \geq 0, \quad i \in C, \\
 &\quad y_i \leq 0, \quad i \in NC, \\
 &\quad y_i \text{ free}, \quad i \in E.
 \end{aligned}$$

From this it is easily established that if we construct the dual of the dual linear program, then we return to the original (primal) linear program.

Examples

In order to illustrate how to construct the dual linear program we present two examples. The first example considers a linear program with matrix block structure. This is a usual form of linear programs and it is particularly easy to construct the dual linear program. The other example deals with the transportation problem presented in Section 8.1. The purpose of constructing the dual to this problem is to show how to handle double subscripted variables and indexed constraints.

Example 10.2 (the dual to a linear program of matrix block form) Consider the linear program

$$\begin{aligned}
 &\text{maximize } \mathbf{c}^T \mathbf{x} + \mathbf{d}^T \mathbf{y}, \\
 &\text{subject to } \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} \leq \mathbf{b}, \\
 &\quad \mathbf{D}\mathbf{y} = \mathbf{e}, \\
 &\quad \mathbf{x} \geq \mathbf{0}^{n_1}, \\
 &\quad \mathbf{y} \leq \mathbf{0}^{n_2},
 \end{aligned}$$

LP duality and sensitivity analysis

where $\mathbf{A} \in \mathbb{R}^{m_1 \times n_1}$, $\mathbf{B} \in \mathbb{R}^{m_1 \times n_2}$, $\mathbf{D} \in \mathbb{R}^{m_2 \times n_2}$, $\mathbf{b} \in \mathbb{R}^{m_1}$, $\mathbf{e} \in \mathbb{R}^{m_2}$, $\mathbf{c} \in \mathbb{R}^{n_1}$, and $\mathbf{d} \in \mathbb{R}^{n_2}$. The dual of this linear program is

$$\begin{aligned} & \text{minimize} && \mathbf{b}^T \mathbf{u} + \mathbf{e}^T \mathbf{v}, \\ & \text{subject to} && \mathbf{A}^T \mathbf{u} \geq \mathbf{c}, \\ & && \mathbf{B}^T \mathbf{u} + \mathbf{D}^T \mathbf{v} \leq \mathbf{d}, \\ & && \mathbf{u} \geq \mathbf{0}^{m_1}, \\ & && \mathbf{v} \text{ free.} \end{aligned}$$

Observe that the constraint matrix of the primal problem is

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{0}^{m_2 \times n_1} & \mathbf{D} \end{pmatrix},$$

and if we transpose this matrix we get

$$\begin{pmatrix} \mathbf{A}^T & \mathbf{0}^{n_1 \times m_2} \\ \mathbf{B}^T & \mathbf{D}^T \end{pmatrix}.$$

Also note that the vector of objective function coefficients of the primal problem, $(\mathbf{c}^T, \mathbf{d}^T)^T$, is the right-hand side of the dual problem, and the right-hand side of the primal problem, $(\mathbf{b}^T, \mathbf{e}^T)^T$, is the vector of objective function coefficients of the dual problem. ■

Example 10.3 (the dual of the transportation problem) Consider the transportation problem (see Example 8.1) to

$$\begin{aligned} & \text{minimize} && z = \sum_{i=1}^N \sum_{j=1}^M c_{ij} x_{ij}, \\ & \text{subject to} && \sum_{j=1}^M x_{ij} \leq s_i, \quad i = 1, \dots, N, \\ & && \sum_{i=1}^N x_{ij} \geq d_j, \quad j = 1, \dots, M, \\ & && x_{ij} \geq 0, \quad i = 1, \dots, N, \quad j = 1, \dots, M. \end{aligned}$$

The dual linear program is given by

$$\begin{aligned} & \text{maximize} && w = \sum_{i=1}^N s_i u_i + \sum_{j=1}^M d_j v_j, \\ & \text{subject to} && u_i + v_j \leq c_{ij}, \quad i = 1, \dots, N, \quad j = 1, \dots, M, \\ & && u_i \leq 0, \quad i = 1, \dots, N, \\ & && v_j \geq 0, \quad j = 1, \dots, M. \end{aligned}$$

Observe that there are $N + M$ constraints in the primal problem and hence $N + M$ dual variables. Also, there are NM variables of the primal problem, hence NM constraints in the dual problem. The form of the constraints in the dual problem arises from the fact that x_{ij} appears twice in the column of the constraint matrix corresponding to this variable: once in the constraints over $i = 1, \dots, N$ and once in the constraints over $j = 1, \dots, M$. Also note that all coefficients of the constraint matrix in the primal problem equal $+1$, and since we have one dual constraint for each column, we finally get the dual constraint $u_i + v_j \leq c_{ij}$. ■

10.3 Linear programming duality theory

In this section we present some of the most fundamental duality theorems. Throughout the section we will consider the primal linear program

$$\begin{aligned} \text{minimize} \quad & z = \mathbf{c}^T \mathbf{x}, \\ \text{subject to} \quad & \mathbf{A}\mathbf{x} = \mathbf{b}, \\ & \mathbf{x} \geq \mathbf{0}^n, \end{aligned} \tag{P}$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{b} \in \mathbb{R}^m$, and $\mathbf{c} \in \mathbb{R}^n$, and its dual linear program

$$\begin{aligned} \text{maximize} \quad & w = \mathbf{b}^T \mathbf{y}, \\ \text{subject to} \quad & \mathbf{A}^T \mathbf{y} \leq \mathbf{c}, \\ & \mathbf{y} \text{ free.} \end{aligned} \tag{D}$$

We note that theorems similar to those presented below can be given also for other primal–dual pairs of linear programs.

10.3.1 Weak and strong duality

We begin by proving the Weak Duality Theorem.

Theorem 10.4 (Weak Duality Theorem) *If \mathbf{x} is a feasible solution to (P) and \mathbf{y} a feasible solution to (D), then $\mathbf{c}^T \mathbf{x} \geq \mathbf{b}^T \mathbf{y}$.*

Proof. We have that

$$\begin{aligned} \mathbf{c}^T \mathbf{x} &\geq (\mathbf{A}^T \mathbf{y})^T \mathbf{x} && [\mathbf{c} \geq \mathbf{A}^T \mathbf{y}, \quad \mathbf{x} \geq \mathbf{0}^n] \\ &= \mathbf{y}^T \mathbf{A}\mathbf{x} = \mathbf{y}^T \mathbf{b} && [\mathbf{A}\mathbf{x} = \mathbf{b}] \\ &= \mathbf{b}^T \mathbf{y}, \end{aligned}$$

and we are done. ■

Corollary 10.5 *If \mathbf{x} is a feasible solution to (P), \mathbf{y} is a feasible solution to (D), and $\mathbf{c}^T \mathbf{x} = \mathbf{b}^T \mathbf{y}$, then \mathbf{x} is an optimal solution to (P) and \mathbf{y} is an optimal solution to (D).* ■

Next we show that the duality gap is zero, that is, strong duality holds. Note that this can also be established by the use of the Lagrangian duality theory in Chapter 6.

Theorem 10.6 (Strong Duality Theorem) *If the primal problem (P) and the dual problem (D) have feasible solutions, then there exist optimal solutions to (P) and (D), and their optimal objective function values are equal.*

Proof. Since the dual (D) is feasible it follows from the Weak Duality Theorem 10.4 that the objective function value of (P) is bounded from below. Hence Theorem 8.10 implies that there exists an optimal BFS, $\mathbf{x}^* = (\mathbf{x}_B^T, \mathbf{x}_N^T)^T$, to (P). We construct an optimal solution to (D). (Actually we have already done this in detail in Section 10.1.) Set

$$(\mathbf{y}^*)^T := \mathbf{c}_B^T \mathbf{B}^{-1}. \quad (10.3)$$

Since \mathbf{x}^* is an optimal basic feasible solution the reduced costs of the non-basic variables are non-negative, which gives that (for details see Section 10.1)

$$\mathbf{A}^T \mathbf{y}^* \leq \mathbf{c}.$$

Hence, \mathbf{y}^* is feasible to (D). Further, we have that

$$\mathbf{b}^T \mathbf{y}^* = \mathbf{b}^T (\mathbf{B}^{-1})^T \mathbf{c}_B = \mathbf{c}_B^T \mathbf{B}^{-1} \mathbf{b} = \mathbf{c}_B^T \mathbf{x}_B = \mathbf{c}^T \mathbf{x}^*,$$

so by Corollary 10.5 it follows that \mathbf{y}^* is an optimal solution to (D). ■

See Exercise 10.15 for another formulation of the Strong Duality Theorem.

Remark 10.7 (dual solution from the primal solution) Note that the proof of Theorem 10.6 is constructive. We construct an optimal dual solution from an optimal basic feasible solution through (10.3).

When a linear program is solved by the simplex method we obtain an optimal basic feasible solution (if the LP is not unbounded or infeasible). Hence from (10.3) we then also—without any additional effort—obtain an optimal dual solution from the last pricing step of the simplex algorithm when we conclude that $\tilde{\mathbf{c}}_N \geq \mathbf{0}^{n-m}$. ■

Interpretation of the optimal dual solution

We have from (10.3) that

$$\mathbf{b}^T \mathbf{y}^* = \mathbf{c}_B^T \mathbf{B}^{-1} \mathbf{b},$$

for any optimal basic feasible solution to (P). If $\mathbf{x}_B > \mathbf{0}^m$, then a small change in \mathbf{b} does not change the basis, and so the optimal value of (D) (and (P)), namely

$$v(\mathbf{b}) := \mathbf{b}^T \mathbf{y}^*$$

is linear at, and locally around, the value \mathbf{b} . If, however, some $(\mathbf{x}_B)_i = 0$, then in this degenerate case it could be that the basis changes in a non-differentiable manner with \mathbf{b} . We summarize:

Theorem 10.8 (shadow price) *If, for a given vector $\mathbf{b} \in \mathbb{R}^m$, the optimal solution to (P) corresponds to a non-degenerate basic feasible solution, then its optimal value is differentiable at \mathbf{b} , with*

$$\frac{\partial v(\mathbf{b})}{\partial b_i} = y_i^*, \quad i = 1, \dots, m,$$

that is, $\nabla v(\mathbf{b}) = \mathbf{y}^*$. ■

Remark 10.9 (shadow price) The optimal dual solution is indeed the shadow price for the constraints. If a unit change in one right-hand side b_i does not change the optimal basis, then the above states that the optimal value will change exactly with the amount y_i^* .

It is also clear that non-degeneracy at \mathbf{x}^* in (P) implies that the optimal solution in (D) must be unique. Namely, we can show that the function v is convex on its effective domain (why?) and the non-degeneracy property clearly implies that v is also finite in a neighbourhood of \mathbf{b} . Then, its differentiability at \mathbf{b} is equivalent to the uniqueness of its subgradients at \mathbf{b} ; cf. Proposition 6.17(c). ■

Farkas' Lemma

In Section 3.2 we proved Farkas' Lemma 3.30 by using the Separation Theorem 3.24. Having access to LP duality, Farkas' Lemma can easily be proved by using the Strong Duality Theorem 10.6.

Theorem 10.10 (Farkas' Lemma) *Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$. Then, exactly one of the systems*

$$\begin{aligned} \mathbf{A}\mathbf{x} &= \mathbf{b}, \\ \mathbf{x} &\geq \mathbf{0}^n, \end{aligned} \tag{I}$$

LP duality and sensitivity analysis

and

$$\begin{aligned} \mathbf{A}^T \mathbf{y} &\leq \mathbf{0}^n, \\ \mathbf{b}^T \mathbf{y} &> 0, \end{aligned} \tag{II}$$

has a feasible solution, and the other system is inconsistent.

Proof. If (I) has a solution \mathbf{x} , then

$$\mathbf{b}^T \mathbf{y} = \mathbf{x}^T \mathbf{A}^T \mathbf{y} > 0.$$

But $\mathbf{x} \geq \mathbf{0}^n$, so $\mathbf{A}^T \mathbf{y} \leq \mathbf{0}^n$ cannot hold, which means that (II) is infeasible.

Assume that (II) is infeasible. Consider the linear program

$$\begin{aligned} \text{maximize} \quad & \mathbf{b}^T \mathbf{y}, \\ \text{subject to} \quad & \mathbf{A}^T \mathbf{y} \leq \mathbf{0}^n, \\ & \mathbf{y} \text{ free}, \end{aligned} \tag{10.4}$$

and its dual program

$$\begin{aligned} \text{minimize} \quad & (\mathbf{0}^n)^T \mathbf{x}, \\ \text{subject to} \quad & \mathbf{A} \mathbf{x} = \mathbf{b}, \\ & \mathbf{x} \geq \mathbf{0}^n. \end{aligned} \tag{10.5}$$

Since (II) is infeasible, $\mathbf{y} = \mathbf{0}^m$ is an optimal solution to (10.4). Hence the Strong Duality Theorem 10.6 implies that there exists an optimal solution to (10.5). This solution is feasible in (I).

What we have proved above is the equivalence

$$(I) \iff \neg(II).$$

Logically, this is equivalent to the statement that

$$\neg(I) \iff (II).$$

We have hence established that precisely one of the two systems (I) and (II) has a solution. ■

10.3.2 Complementary slackness

A further relationship between (P) and (D) at an optimal solution is given by the Complementary Slackness Theorem.

Theorem 10.11 (Complementary Slackness Theorem) *Let \mathbf{x} be a feasible solution to (P) and \mathbf{y} a feasible solution to (D). Then*

$$\left. \begin{array}{l} \mathbf{x} \text{ optimal to (P)} \\ \mathbf{y} \text{ optimal to (D)} \end{array} \right\} \iff x_j(c_j - \mathbf{A}_{\cdot j}^T \mathbf{y}) = 0, \quad j = 1, \dots, n, \quad (10.6)$$

where $\mathbf{A}_{\cdot j}$ is the j^{th} column of \mathbf{A} .

Proof. If \mathbf{x} and \mathbf{y} are feasible we get

$$\mathbf{c}^T \mathbf{x} \geq (\mathbf{A}^T \mathbf{y})^T \mathbf{x} = \mathbf{y}^T \mathbf{A} \mathbf{x} = \mathbf{b}^T \mathbf{y}. \quad (10.7)$$

Further, by the Strong Duality Theorem 10.6 and the Weak Duality Theorem 10.4, \mathbf{x} and \mathbf{y} are optimal if and only if $\mathbf{c}^T \mathbf{x} = \mathbf{b}^T \mathbf{y}$, so in fact (10.7) holds with equality, that is,

$$\mathbf{c}^T \mathbf{x} = (\mathbf{A}^T \mathbf{y})^T \mathbf{x} \iff \mathbf{x}^T (\mathbf{c} - \mathbf{A}^T \mathbf{y}) = 0.$$

Since $\mathbf{x} \geq \mathbf{0}^n$ and $\mathbf{A}^T \mathbf{y} \leq \mathbf{c}$, $\mathbf{x}^T (\mathbf{c} - \mathbf{A}^T \mathbf{y}) = 0$ is equivalent to each term in the sum being zero, that is, that (10.6) holds. ■

Often the Complementary Slackness Theorem is stated for the primal–dual pair given by

$$\begin{array}{ll} \text{maximize} & \mathbf{c}^T \mathbf{x} \\ \text{subject to} & \mathbf{A} \mathbf{x} \leq \mathbf{b}, \\ & \mathbf{x} \geq \mathbf{0}^n, \end{array} \quad (10.8)$$

and

$$\begin{array}{ll} \text{minimize} & \mathbf{b}^T \mathbf{y} \\ \text{subject to} & \mathbf{A}^T \mathbf{y} \geq \mathbf{c}, \\ & \mathbf{y} \geq \mathbf{0}^m. \end{array} \quad (10.9)$$

The Complementary Slackness Theorem then becomes as follows. (Its proof is similar to that of Theorem 10.11.)

Theorem 10.12 (Complementary Slackness Theorem) *Let \mathbf{x} be a feasible solution to (10.8) and \mathbf{y} a feasible solution to (10.9). Then \mathbf{x} is optimal to (10.8) and \mathbf{y} optimal to (10.9) if and only if*

$$x_j(c_j - \mathbf{y}^T \mathbf{A}_{\cdot j}) = 0, \quad j = 1, \dots, n, \quad (10.10a)$$

$$y_i(\mathbf{A}_{i \cdot} \mathbf{x} - b_i) = 0, \quad i = 1, \dots, m, \quad (10.10b)$$

where $\mathbf{A}_{\cdot j}$ is the j^{th} column of \mathbf{A} and $\mathbf{A}_{i \cdot}$ the i^{th} row of \mathbf{A} . ■

Remark 10.13 (interpretation of the Complementary Slackness Theorem)

From the Complementary Slackness Theorem 10.12 follows that, for a primal–dual pair of optimal solutions, if there is slack in one constraint, then the respective variable in the other problem is zero. Further, if a variable is positive, then there is no slack in the respective constraint in the other problem. ■

Example 10.14 (the transportation problem) Consider again the transportation problem (cf. Examples 8.1 and 10.3). The following complementarity conditions are particularly simple and intuitive:

$$x_{ij}(u_i + v_j - c_{ij}) = 0, \quad i = 1, \dots, N, \quad j = 1, \dots, M.$$

Hence, given an optimal dual solution $(\mathbf{u}^*, \mathbf{v}^*)$, transportation can occur on link (i, j) , that is, $x_{ij}^* > 0$ may hold, only if $c_{ij} = u_i^* + v_j^*$ holds. The dual variables may be viewed as *node prices* that determine whether the price c_{ij} of transportation is sufficiently low. See Section 10.6 for further reading on duality in linear network flow optimization. ■

The consequence of the Complementary Slackness Theorem is the following characterization of an optimal solution to a linear program. We state it for the primal–dual pair given by (10.8) and (10.9), but it holds as well for each primal–dual pair of linear programs.

Theorem 10.15 (necessary and sufficient conditions for global optimality)

For $\mathbf{x} \in \mathbb{R}^n$ to be an optimal solution to the linear program (10.8), it is both necessary and sufficient that

- (a) \mathbf{x} is a feasible solution to (10.8);
- (b) corresponding to \mathbf{x} there is a dual feasible solution $\mathbf{y} \in \mathbb{R}^m$ to (10.9); and
- (c) the primal–dual pair (\mathbf{x}, \mathbf{y}) satisfies the complementarity conditions (10.10). ■

The simplex method is very well adapted to these conditions. After phase I, (a) holds. Every basic solution (feasible or not) satisfies (c), since if x_j is in the basis, then $\tilde{c}_j = c_j - \mathbf{y}^T \mathbf{A}_{.j} = 0$, and if $\tilde{c}_j \neq 0$, then $x_j = 0$. So, the only condition that the simplex method does not satisfy for every basic feasible solution is (b). The proof of the Strong Duality Theorem 10.6 shows that it is satisfied exactly at an optimal basic feasible solution. The entering criterion is based on trying to better satisfy it. Indeed, by choosing as an entering variable x_j such that

$$j \in \arg \min_{j \in \{1, \dots, n\}} \tilde{c}_j,$$

we actually identify a dual constraint

$$\sum_{i=1}^m a_{ij}y_i \leq c_j,$$

which is among the most violated at the complementary solution $\mathbf{y}^T = \mathbf{c}_B^T \mathbf{B}^{-1}$ given by the current BFS. After the basis change we will have equality in this dual constraint, and hence the basis change corresponds to making a currently most violated dual constraint feasible!

Example 10.16 (illustration of complementary slackness) Consider the primal–dual pair given by

$$\begin{aligned} \text{maximize} \quad & z = 3x_1 + 2x_2, \\ \text{subject to} \quad & x_1 + x_2 \leq 80, \\ & 2x_1 + x_2 \leq 100, \\ & x_1 \leq 40, \\ & x_1, x_2 \geq 0, \end{aligned} \tag{10.11}$$

and

$$\begin{aligned} \text{minimize} \quad & w = 80y_1 + 100y_2 + 40y_3, \\ \text{subject to} \quad & y_1 + 2y_2 + y_3 \geq 3, \\ & y_1 + y_2 \geq 2, \\ & y_1, y_2, y_3 \geq 0. \end{aligned} \tag{10.12}$$

We use Theorem 10.15 to show that $\mathbf{x}^* = (20, 60)^T$ is an optimal solution to (10.11).

- (a) (primal feasibility) Obviously \mathbf{x}^* is a feasible solution to (10.11).
- (c) (complementarity) The complementarity conditions imply that

$$\begin{aligned} y_1^*(x_1^* + x_2^* - 80) &= 0 \\ y_2^*(2x_1^* + x_2^* - 100) &= 0 \\ y_3^*(x_1^* - 40) &= 0 \implies y_3^* = 0 \quad [x_1^* = 20 \neq 40] \\ x_1^*(y_1^* + 2y_2^* + y_3^* - 3) &= 0 \implies y_1^* + 2y_2^* = 3 \quad [x_1^* > 0] \\ x_2^*(y_1^* + y_2^* - 2) &= 0 \implies y_1^* + y_2^* = 2 \quad [x_2^* > 0] \end{aligned}$$

which gives that $y_1^* = 1$, $y_2^* = 1$ and $y_3^* = 0$.

- (b) (dual feasibility) Clearly $\mathbf{y}^* = (1, 1, 0)^T$ is feasible in (10.12).

From Theorem 10.15 it then follows that $\mathbf{x}^* = (20, 60)^T$ is an optimal solution to (10.11) and $\mathbf{y}^* = (1, 1, 0)^T$ an optimal solution to (10.12). ■

10.4 The dual simplex method

The simplex method presented in Chapter 9, which we here refer to as the *primal simplex method*, starts with a basic feasible solution to the primal linear program and then iterates until the primal optimality conditions are fulfilled, that is, until a basic feasible solution is found such that the reduced costs satisfy

$$\tilde{\mathbf{c}}_N^T := \mathbf{c}_N^T - \mathbf{c}_B^T \mathbf{B}^{-1} \mathbf{N} \geq (\mathbf{0}^{n-m})^T.$$

This is equivalent to the dual feasibility condition

$$\mathbf{A}^T \mathbf{y} \leq \mathbf{c},$$

where $\mathbf{y} := (\mathbf{B}^{-1})^T \mathbf{c}_B$. We call a basis such that all of the reduced costs are greater than or equal to zero a *dual feasible basis*; otherwise we call it a *dual infeasible basis*. Hence, the primal simplex method starts with a primal feasible basis and then moves through a sequence of dual infeasible (but primal feasible) bases until a dual feasible basis is found.

The *dual simplex method* is a variant of the primal simplex method that works in a dual manner, in the sense that it starts with a dual feasible basis and then moves through a sequence of primal infeasible (but dual feasible) bases until a primal (and dual) feasible basis is found.

In order to derive the dual simplex algorithm, let \mathbf{x}_B be a dual feasible basis with the corresponding partition (\mathbf{B}, \mathbf{N}) . If

$$\tilde{\mathbf{b}} := \mathbf{B}^{-1} \mathbf{b} \geq \mathbf{0}^m,$$

then \mathbf{x}_B is primal feasible and since it is also dual feasible all of the reduced costs are greater than or equal to zero; hence, \mathbf{x}_B is an optimal BFS. Otherwise some of the components of $\tilde{\mathbf{b}}$ is strictly negative, say \tilde{b}_1 , that is,

$$(\mathbf{x}_B)_1 + \sum_{j=1}^{n-m} (\mathbf{B}^{-1} \mathbf{N})_{1j} (\mathbf{x}_N)_j = \tilde{b}_1 < 0,$$

so $(\mathbf{x}_B)_1 < 0$ in the current basis and will be the leaving variable. If

$$(\mathbf{B}^{-1} \mathbf{N})_{1j} \geq 0, \quad j = 1, \dots, n-m, \quad (10.13)$$

then there exists no primal feasible solution to the problem. (Why?) Hence, if (10.13) is fulfilled, then we say that the *primal infeasibility criterion* is satisfied. Otherwise $(\mathbf{B}^{-1} \mathbf{N})_{1j} < 0$ for some $j = 1, \dots, n-m$. Assume that $(\mathbf{B}^{-1} \mathbf{N})_{1k} < 0$ and choose $(\mathbf{x}_N)_k$ to replace $(\mathbf{x}_B)_1$ in the

basis. (Note that this yields that $(\mathbf{x}_N)_k = \tilde{b}_1/(\mathbf{B}^{-1}\mathbf{N})_{1k} > 0$ in the new basis.) The new reduced costs then become

$$\begin{aligned} (\bar{\mathbf{c}}_B)_1 &:= -\frac{(\tilde{\mathbf{c}}_N)_k}{(\mathbf{B}^{-1}\mathbf{N})_{1k}}, \\ (\bar{\mathbf{c}}_B)_j &:= 0, \quad j = 2, \dots, m, \\ (\bar{\mathbf{c}}_N)_j &:= (\tilde{\mathbf{c}}_N)_j - (\tilde{\mathbf{c}}_N)_k \frac{(\mathbf{B}^{-1}\mathbf{N})_{1j}}{(\mathbf{B}^{-1}\mathbf{N})_{1k}}, \quad j = 1, \dots, n-m. \end{aligned}$$

Since we want the new basis to be dual feasible it must hold that all of the new reduced costs are non-negative, that is,

$$(\bar{\mathbf{c}}_N)_j \geq (\tilde{\mathbf{c}}_N)_k \frac{(\mathbf{B}^{-1}\mathbf{N})_{1j}}{(\mathbf{B}^{-1}\mathbf{N})_{1k}}, \quad j = 1, \dots, n-m,$$

or, equivalently,

$$\frac{(\tilde{\mathbf{c}}_N)_k}{(\mathbf{B}^{-1}\mathbf{N})_{1k}} \geq \frac{(\tilde{\mathbf{c}}_N)_j}{(\mathbf{B}^{-1}\mathbf{N})_{1j}}, \quad \text{for all } j \text{ such that } (\mathbf{B}^{-1}\mathbf{N})_{1j} < 0.$$

Therefore, in order to preserve dual feasibility, as entering variable we must choose $(\mathbf{x}_N)_k$ such that

$$k \in \arg \max_{i \in \{j \mid (\mathbf{B}^{-1}\mathbf{N})_{1j} < 0\}} \frac{(\tilde{\mathbf{c}}_N)_j}{(\mathbf{B}^{-1}\mathbf{N})_{1j}}.$$

We have now derived an infeasibility criterion and criteria for choosing the leaving and the entering variables, and are ready to state the dual simplex algorithm.

The Dual Simplex Algorithm:

Step 0 (initialization: DFS) Assume that $\mathbf{x} = (\mathbf{x}_B^T, \mathbf{x}_N^T)^T$ is a dual feasible basis corresponding to the partition $\mathbf{A} = (\mathbf{B}, \mathbf{N})$.

Step 1 (leaving variable or termination) Calculate

$$\tilde{\mathbf{b}} := \mathbf{B}^{-1}\mathbf{b}.$$

If $\tilde{\mathbf{b}} \geq \mathbf{0}^m$, then stop; the current basis is optimal. Otherwise, choose an s such that $\tilde{b}_s < 0$, and let $(\mathbf{x}_B)_s$ be the leaving variable.

Step 2 (entering variable or termination) If

$$(\mathbf{B}^{-1}\mathbf{N})_{sj} \geq 0, \quad j = 1, \dots, n-m,$$

then stop; the (primal) problem is infeasible. Otherwise, choose a k such that

$$k \in \arg \max_{i \in \{j \mid (B^{-1}N)_{sj} < 0\}} \frac{(\tilde{c}_N)_j}{(B^{-1}N)_{sj}},$$

and let $(x_N)_k$ be the entering variable.

Step 3 (update: change basis) Construct a new partition by swapping $(x_B)_s$ with $(x_N)_k$. Go to Step 1.

Similarly to the primal simplex algorithm it can be shown that the dual simplex algorithm terminates in a finite number of steps if cycling is avoided. Also, there exist rules for choosing the leaving and entering variables (among the eligible ones) such that cycling is avoided.

If a dual feasible solution is not available from the start, it is possible to add a constraint to the original problem that makes it possible to construct a dual feasible basis, and then run the dual simplex algorithm on this modified problem (see Exercise 10.12).

Remark 10.17 (unboundedness of the primal problem) If the dual problem is known to be feasible, the primal problem cannot be unbounded by the Weak Duality Theorem 10.4. Hence the dual simplex algorithm terminates with a basis that satisfies either the optimality criterion or the primal infeasibility criterion. ■

Example 10.18 (illustration of the dual simplex algorithm) Consider the linear program

$$\begin{aligned} \text{minimize} \quad & 3x_1 + 4x_2 + 2x_3 + x_4 + 5x_5, \\ \text{subject to} \quad & x_1 - 2x_2 - x_3 + x_4 + x_5 \leq -3, \\ & -x_1 - x_2 - x_3 + x_4 + x_5 \leq -2, \\ & x_1 + x_2 - 2x_3 + 2x_4 - 3x_5 \leq 4, \\ & x_1, \quad x_2, \quad x_3, \quad x_4, \quad x_5 \geq 0. \end{aligned}$$

By introducing the slack variables x_6, x_7, x_8 , we get the following linear program:

$$\begin{aligned} \text{minimize} \quad & 3x_1 + 4x_2 + 2x_3 + x_4 + 5x_5, \\ \text{subject to} \quad & x_1 - 2x_2 - x_3 + x_4 + x_5 + x_6 = -3, \\ & -x_1 - x_2 - x_3 + x_4 + x_5 + x_7 = -2, \\ & x_1 + x_2 - 2x_3 + 2x_4 - 3x_5 + x_8 = 4, \\ & x_1, \quad x_2, \quad x_3, \quad x_4, \quad x_5, \quad x_6, \quad x_7, \quad x_8 \geq 0. \end{aligned}$$

We see that the basis $\mathbf{x}_B := (x_6, x_7, x_8)^T$ is dual feasible, but primal infeasible. Hence we use the dual simplex algorithm to solve the problem. We have that

$$\tilde{\mathbf{b}} := \mathbf{B}^{-1}\mathbf{b} = (-3, -2, 4)^T,$$

so we choose $(\mathbf{x}_B)_1 = x_6$ to leave the basis. Further we have that the reduced costs of $\mathbf{x}_N := (x_1, x_2, x_3, x_4, x_5)^T$ are

$$\tilde{\mathbf{c}}_N^T = (3, 4, 2, 1, 5),$$

and

$$(\mathbf{B}^{-1}\mathbf{N})_1 = (1, -2, -1, 1, 1), \quad [\text{the 1st row of } \mathbf{B}^{-1}\mathbf{N}]$$

so we choose x_2 as the entering variable. The new basis becomes $\mathbf{x}_B := (x_2, x_7, x_8)^T$, $\mathbf{x}_N := (x_1, x_6, x_3, x_4, x_5)^T$. We get that

$$\tilde{\mathbf{b}} := \mathbf{B}^{-1}\mathbf{b} = (1.5, -0.5, 2.5)^T.$$

Hence, we choose $(\mathbf{x}_B)_2 = x_7$ as the leaving variable. Further,

$$\begin{aligned} \tilde{\mathbf{c}}_N^T &= (5, 2, 0, 3, 7), \\ (\mathbf{B}^{-1}\mathbf{N})_2 &= (1.5, -0.5, -0.5, 0.5, 0.5)^T, \end{aligned}$$

which gives that x_3 is the entering variable. The new basis becomes $\mathbf{x}_B := (x_2, x_3, x_8)^T$. We get that

$$\tilde{\mathbf{b}} := \mathbf{B}^{-1}\mathbf{b} = (1, 1, 5)^T,$$

which means that the optimality criterion (primal feasibility) is satisfied, and an optimal solution to the original problem is given by

$$\mathbf{x}^* = (x_1, x_2, x_3, x_4, x_5)^T = (0, 1, 1, 0, 0)^T.$$

Check that this is indeed true, for example by using Theorem 10.12. ■

10.5 Sensitivity analysis

In this section we study two kinds of perturbations of a linear program in standard form,

$$\begin{aligned} \text{minimize} \quad & z = \mathbf{c}^T \mathbf{x}, \\ \text{subject to} \quad & \mathbf{A}\mathbf{x} = \mathbf{b}, \\ & \mathbf{x} \geq \mathbf{0}^n, \end{aligned} \tag{10.14}$$

namely

1. perturbations in the objective function coefficients c_j ; and
2. perturbations in the right-hand side coefficients b_i .

We assume that $\mathbf{x}^* = (\mathbf{x}_B^T, \mathbf{x}_N^T)^T = ((\mathbf{B}^{-1}\mathbf{b})^T, (\mathbf{0}^{n-m})^T)^T$ is an optimal basic feasible solution to (10.14) with the corresponding partition $\mathbf{A} = (\mathbf{B}, \mathbf{N})$.

10.5.1 Perturbations in the objective function

Assume that the objective function coefficients of the linear program (10.14) are perturbed by the vector $\mathbf{p} \in \mathbb{R}^n$, that is, we consider the perturbed problem to

$$\begin{aligned} \text{minimize} \quad & \tilde{z} = (\mathbf{c} + \mathbf{p})^T \mathbf{x}, \\ \text{subject to} \quad & \mathbf{A}\mathbf{x} = \mathbf{b}, \\ & \mathbf{x} \geq \mathbf{0}^n. \end{aligned} \tag{10.15}$$

The optimal solution \mathbf{x}^* to the unperturbed problem (10.14) is obviously a feasible solution to (10.15), but is it still optimal? To answer this question, we note that a basic feasible solution is optimal if the reduced costs of the non-basic variables are greater than or equal to zero. The reduced costs for the non-basic variables of the perturbed problem (10.15) are given by [let $\mathbf{p} = (\mathbf{p}_B^T, \mathbf{p}_N^T)^T$]

$$\bar{\mathbf{c}}_N^T = (\mathbf{c}_N + \mathbf{p}_N)^T - (\mathbf{c}_B + \mathbf{p}_B)^T \mathbf{B}^{-1} \mathbf{N}.$$

Hence, $\bar{\mathbf{c}}_N \geq \mathbf{0}^{n-m}$ is sufficient for \mathbf{x}^* to be an optimal solution to the perturbed problem (10.15). (Observe, however, that this is not a necessary condition unless \mathbf{x}^* is non-degenerate.)

Perturbations of a non-basic cost coefficient

If only one component of \mathbf{c}_N is perturbed, that is,

$$\mathbf{p} = \begin{pmatrix} \mathbf{p}_B \\ \mathbf{p}_N \end{pmatrix} = \begin{pmatrix} \mathbf{0}^m \\ \varepsilon \mathbf{e}_j \end{pmatrix},$$

for some $\varepsilon \in \mathbb{R}$ and $j \in \{1, \dots, n - m\}$, then we have that \mathbf{x}^* is an optimal solution to the perturbed problem if

$$(\mathbf{c}_N)_j + \varepsilon - \mathbf{c}_B^T \mathbf{B}^{-1} \mathbf{N}_j \geq 0 \quad \Longleftrightarrow \quad \varepsilon + (\bar{\mathbf{c}}_N)_j \geq 0,$$

so in this case we only have to check that the perturbation ε is not less than $-(\bar{\mathbf{c}}_N)_j$ in order to guarantee that \mathbf{x}^* is an optimal solution to the perturbed problem.

Perturbations of a basic cost coefficient

If only one component of \mathbf{c}_B is perturbed, that is,

$$\mathbf{p} = \begin{pmatrix} \mathbf{p}_B \\ \mathbf{p}_N \end{pmatrix} = \begin{pmatrix} \varepsilon \mathbf{e}_j \\ \mathbf{0}^{n-m} \end{pmatrix},$$

for some $\varepsilon \in \mathbb{R}$ and $j \in \{1, \dots, m\}$, then we have that \mathbf{x}^* is an optimal solution to the perturbed problem if

$$\begin{aligned} (\mathbf{c}_N)^T - (\mathbf{c}_B^T + \varepsilon \mathbf{e}_j^T) \mathbf{B}^{-1} \mathbf{N} &\geq (\mathbf{0}^{n-m})^T \\ \iff \\ -\varepsilon \mathbf{e}_j^T \mathbf{B}^{-1} \mathbf{N} + \tilde{\mathbf{c}}_N^T &\geq (\mathbf{0}^{n-m})^T. \end{aligned}$$

In this case all of the reduced costs of the non-basic variables may change, and we must check that the perturbation ε multiplied by the j^{th} row of $-\mathbf{B}^{-1} \mathbf{N}$ plus the original reduced costs $\tilde{\mathbf{c}}_N^T$ is a vector whose components all are greater than or equal to zero.

Perturbations that make \mathbf{x}^* non-optimal

If the perturbation \mathbf{p} is such that some of the reduced costs of the perturbed problem becomes strictly negative for the basis \mathbf{x}_B , then \mathbf{x}^* is perhaps not an optimal solution anymore. If this happens, let some of the variables with strictly negative reduced cost enter the basis and continue the simplex algorithm until an optimal solution is found (or until the unboundedness criterion is satisfied).

10.5.2 Perturbations in the right-hand side coefficients

Now, assume that the right-hand side \mathbf{b} of the linear program (10.14) is perturbed by the vector $\mathbf{p} \in \mathbb{R}^m$, that is, we consider the perturbed problem to

$$\begin{aligned} \text{minimize} \quad & z = \mathbf{c}^T \mathbf{x}, \\ \text{subject to} \quad & \mathbf{A} \mathbf{x} = \mathbf{b} + \mathbf{p}, \\ & \mathbf{x} \geq \mathbf{0}^n. \end{aligned} \tag{10.16}$$

The reduced costs of the unperturbed problem do not change as the right-hand side is perturbed, so the basic feasible solution given by the partition $\mathbf{A} = (\mathbf{B}, \mathbf{N})$ is optimal to the perturbed problem (10.16) if

and only if it is feasible, that is,

$$\begin{pmatrix} \mathbf{x}_B \\ \mathbf{x}_N \end{pmatrix} = \begin{pmatrix} \mathbf{B}^{-1}(\mathbf{b} + \mathbf{p}) \\ \mathbf{0}^{n-m} \end{pmatrix} \geq \mathbf{0}^n,$$

which means that we have to check that $\mathbf{B}^{-1}(\mathbf{b} + \mathbf{p}) \geq \mathbf{0}^m$.

Perturbations of one component of the right-hand side

Suppose that only one of the components of the right-hand side is perturbed, that is,

$$\mathbf{p} = \varepsilon \mathbf{e}_j,$$

for some $\varepsilon \in \mathbb{R}$ and $j \in \{1, \dots, m\}$. The basic feasible solution corresponding to the partition $\mathbf{A} = (\mathbf{B}, \mathbf{N})$ is then feasible if and only if

$$\mathbf{B}^{-1}(\mathbf{b} + \varepsilon \mathbf{e}_j) \geq \mathbf{0}^m \iff \varepsilon \mathbf{B}^{-1} \mathbf{e}_j + \mathbf{B}^{-1} \mathbf{b} \geq \mathbf{0}^m,$$

so it must hold that ε multiplied by the j^{th} column of \mathbf{B}^{-1} plus the vector $\mathbf{B}^{-1} \mathbf{b}$ equals a vector whose components all are greater than or equal to zero.

Perturbations that make \mathbf{x}^* infeasible

If the perturbation \mathbf{p} is such that the basis \mathbf{x}_B becomes infeasible, then some of the components in the updated right-hand side, $\mathbf{B}^{-1}(\mathbf{b} + \mathbf{p})$, is strictly negative. However, the reduced costs are independent of \mathbf{p} , so the basis \mathbf{x}_B is still a *dual* feasible basis. Hence, we can continue with the dual simplex algorithm until an optimal solution is found (or until the primal infeasibility criterion is satisfied).

10.6 Notes and further reading

For an account of the early history of LP duality theory, see [LRS91].

Linear programming duality theory was introduced by John von Neumann [vNe47]. His results build upon his earlier work in game theory. The first published proof of the Strong Duality Theorem is found in Gale, Kuhn, and Tucker [GKT51]. The Complementary Slackness Theorem is due to Dantzig and Orden [DaO53].

Text books that discuss LP duality and sensitivity analysis are [Dan63, Chv83, Mur83, Sch86, DaT97, Pad99, Van01, DaT03, DaM05].

More on the modelling of, and algorithms and duality for, linear network optimization can be found in [AMO93].

10.7 Exercises

Exercise 10.1 (constructing the LP dual) Consider the linear program

$$\begin{aligned}
 &\text{maximize} && z = 6x_1 - 3x_2 - 2x_3 + 5x_4, \\
 &\text{subject to} && 4x_1 + 3x_2 - 8x_3 + 7x_4 = 11, \\
 & && 3x_1 + 2x_2 + 7x_3 + 6x_4 \geq 23, \\
 & && 7x_1 + 4x_2 + 3x_3 + 2x_4 \leq 12, \\
 & && x_1, \quad x_2 \quad \quad \geq 0, \\
 & && \quad \quad \quad x_3 \quad \quad \leq 0, \\
 & && \quad \quad \quad x_4 \quad \text{free}.
 \end{aligned}$$

Construct its linear programming dual.

Exercise 10.2 (constructing the LP dual) Consider the linear program

$$\begin{aligned}
 &\text{minimize} && z = \mathbf{c}^T \mathbf{x}, \\
 &\text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b}, \\
 & && \mathbf{l} \leq \mathbf{x} \leq \mathbf{u}.
 \end{aligned}$$

- (a) Construct its linear programming dual.
- (b) Show that the dual problem is always feasible (independently of \mathbf{A} , \mathbf{b} , \mathbf{l} , and \mathbf{u}).

Exercise 10.3 (application of the Weak and Strong Duality Theorems) Consider the linear program

$$\begin{aligned}
 &\text{minimize} && z = \mathbf{c}^T \mathbf{x}, \\
 &\text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b}, \\
 & && \mathbf{x} \geq \mathbf{0}^n,
 \end{aligned} \tag{P}$$

and the perturbed problem to

$$\begin{aligned}
 &\text{minimize} && z = \mathbf{c}^T \mathbf{x}, \\
 &\text{subject to} && \mathbf{A}\mathbf{x} = \tilde{\mathbf{b}}, \\
 & && \mathbf{x} \geq \mathbf{0}^n.
 \end{aligned} \tag{P'}$$

Show that if (P) has an optimal solution, then the perturbed problem (P') cannot be unbounded (independently of $\tilde{\mathbf{b}}$).

Exercise 10.4 (application of the Weak and Strong Duality Theorems) Consider the linear program

$$\begin{aligned}
 &\text{minimize} && z = \mathbf{c}^T \mathbf{x}, \\
 &\text{subject to} && \mathbf{A}\mathbf{x} \leq \mathbf{b}.
 \end{aligned} \tag{10.17}$$

Assume that the objective function vector \mathbf{c} cannot be written as a linear combination of the rows of \mathbf{A} . Show that (10.17) cannot have an optimal solution.

LP duality and sensitivity analysis

Exercise 10.5 (application of the Weak and Strong Duality Theorems) Consider the linear program

$$\begin{aligned} \text{minimize} \quad & z = \mathbf{c}^T \mathbf{x}, \\ \text{subject to} \quad & \mathbf{A}\mathbf{x} \geq \mathbf{b}, \\ & \mathbf{x} \geq \mathbf{0}^n. \end{aligned} \tag{10.18}$$

Construct a polyhedron that equals the set of optimal solutions to (10.18).

Exercise 10.6 (application of the Weak and Strong Duality Theorems) Consider the linear program

$$\begin{aligned} \text{minimize} \quad & z = \mathbf{c}^T \mathbf{x}, \\ \text{subject to} \quad & \mathbf{A}\mathbf{x} \leq \mathbf{b}, \\ & \mathbf{x} \geq \mathbf{0}^n. \end{aligned} \tag{10.19}$$

Let \mathbf{x}^* be optimal in (10.19) with the optimal value z^* , and let \mathbf{y}^* be optimal in the LP dual of (10.19). Show that

$$z^* = (\mathbf{y}^*)^T \mathbf{A}\mathbf{x}^*.$$

Exercise 10.7 (linear programming primal–dual optimality conditions) Consider the linear program

$$\begin{aligned} \text{maximize} \quad & z = -4x_2 + 3x_3 + 2x_4 - 8x_5, \\ \text{subject to} \quad & 3x_1 + x_2 + 2x_3 + x_4 = 3, \\ & x_1 - x_2 + x_4 - x_5 \geq 2, \\ & x_1, \quad x_2, \quad x_3, \quad x_4, \quad x_5 \geq 0. \end{aligned}$$

Use the LP primal–dual optimality conditions to find an optimal solution.

Exercise 10.8 (linear programming primal–dual optimality conditions) Consider the linear program (the continuous knapsack problem)

$$\begin{aligned} \text{maximize} \quad & z = \mathbf{c}^T \mathbf{x}, \\ \text{subject to} \quad & \mathbf{a}^T \mathbf{x} \leq b, \\ & \mathbf{x} \leq \mathbf{1}^n, \\ & \mathbf{x} \geq \mathbf{0}^n, \end{aligned}$$

where $\mathbf{c} > \mathbf{0}^n$, $\mathbf{a} > \mathbf{0}^n$, $b > 0$, and

$$\frac{c_1}{a_1} \geq \frac{c_2}{a_2} \geq \cdots \geq \frac{c_n}{a_n}.$$

Show that the feasible solution \mathbf{x} given by

$$x_j = 1, \quad j = 1, \dots, r-1, \quad x_r = \frac{b - \sum_{j=1}^{r-1} a_j}{a_r}, \quad x_j = 0, \quad j = r+1, \dots, n,$$

where r is such that $\sum_{j=1}^{r-1} a_j \leq b$ and $\sum_{j=1}^r a_j > b$, is an optimal solution.

Exercise 10.9 Prove Theorem 10.15.

Exercise 10.10 (KKT versus LP primal–dual optimality conditions) Consider the linear program

$$\begin{array}{ll}\text{minimize} & z = \mathbf{c}^T \mathbf{x}, \\ \text{subject to} & \mathbf{A}\mathbf{x} \leq \mathbf{b}.\end{array}$$

Show that the KKT conditions are equivalent to the LP primal–dual optimality conditions.

Exercise 10.11 (Lagrangian primal–dual versus LP primal–dual) Consider the linear program

$$\begin{array}{ll}\text{minimize} & z = \mathbf{c}^T \mathbf{x}, \\ \text{subject to} & \mathbf{A}\mathbf{x} \leq \mathbf{b}.\end{array}$$

Show that the Lagrangian primal–dual optimality conditions are equivalent to the LP primal–dual optimality conditions.

Exercise 10.12 (the dual simplex method) Show that by adding the constraint

$$x_1 + \cdots + x_n \leq M,$$

where M is a positive constant, to a linear program in standard form, it is always possible to construct a dual feasible basis.

Exercise 10.13 (sensitivity analysis: perturbations in the objective function) Consider the linear program

$$\begin{array}{ll}\text{maximize} & z = -x_1 + 18x_2 + c_3x_3 + c_4x_4, \\ \text{subject to} & \begin{array}{l} x_1 + 2x_2 + 3x_3 + 4x_4 \leq 3, \\ -3x_1 + 4x_2 - 5x_3 - 6x_4 \leq 1, \\ x_1, \quad x_2, \quad x_3, \quad x_4 \geq 0. \end{array}\end{array}$$

Find the values of c_3 and c_4 such that the basic solution that corresponds to the partition $\mathbf{x}_B := (x_1, x_2)^T$ is an optimal basic feasible solution to the problem.

Exercise 10.14 (sensitivity analysis: perturbations in the right-hand side) Consider the linear program

$$\begin{array}{ll}\text{minimize} & z = -x_1 + 2x_2 + x_3, \\ \text{subject to} & \begin{array}{l} 2x_1 + x_2 - x_3 \leq 7, \\ -x_1 + 2x_2 + 3x_3 \geq 3 + \delta, \\ x_1, \quad x_2, \quad x_3 \geq 0. \end{array}\end{array}$$

LP duality and sensitivity analysis

(a) Let $\delta = 0$. Show that the basic solution that corresponds to the partition $\mathbf{x}_B := (x_1, x_3)^T$ is an optimal solution to the problem.

(b) Find the values of the perturbation $\delta \in \mathbb{R}$ such that the above BFS is optimal.

(c) Find an optimal solution when $\delta = -7$.

Exercise 10.15 (a version of the Strong Duality Theorem) Consider the linear program

$$\begin{aligned} \text{minimize} \quad & z = \mathbf{c}^T \mathbf{x}, \\ \text{subject to} \quad & \mathbf{A}\mathbf{x} = \mathbf{b}, \\ & \mathbf{x} \geq \mathbf{0}^n, \end{aligned} \tag{P}$$

and its dual linear program

$$\begin{aligned} \text{maximize} \quad & w = \mathbf{b}^T \mathbf{y}, \\ \text{subject to} \quad & \mathbf{A}^T \mathbf{y} \leq \mathbf{c}, \\ & \mathbf{y} \text{ free.} \end{aligned} \tag{D}$$

Show that if one of the problems (P) and (D) has a finite optimal solution, then so does its dual, and their optimal objective function values are equal.

Exercise 10.16 (an LP duality paradox) For a standard primal–dual pair of LPs, consider the following string of inequalities:

$$\begin{aligned} \text{maximum} \{ \mathbf{c}^T \mathbf{x} \mid \mathbf{A}\mathbf{x} \leq \mathbf{b}; \mathbf{x} \geq \mathbf{0}^n \} &\leq \text{minimum} \{ \mathbf{b}^T \mathbf{y} \mid \mathbf{A}^T \mathbf{y} \geq \mathbf{c}; \mathbf{y} \geq \mathbf{0}^m \} \\ &\leq \text{maximum} \{ \mathbf{b}^T \mathbf{y} \mid \mathbf{A}^T \mathbf{y} \geq \mathbf{c}; \mathbf{y} \geq \mathbf{0}^m \} \\ &\leq \text{minimum} \{ \mathbf{c}^T \mathbf{x} \mid \mathbf{A}\mathbf{x} \geq \mathbf{b}; \mathbf{x} \leq \mathbf{0}^n \} \\ &\leq \text{maximum} \{ \mathbf{c}^T \mathbf{x} \mid \mathbf{A}\mathbf{x} \geq \mathbf{b}; \mathbf{x} \leq \mathbf{0}^n \} \\ &\leq \text{minimum} \{ \mathbf{b}^T \mathbf{y} \mid \mathbf{A}^T \mathbf{y} \leq \mathbf{c}; \mathbf{y} \leq \mathbf{0}^m \} \\ &\leq \text{maximum} \{ \mathbf{b}^T \mathbf{y} \mid \mathbf{A}^T \mathbf{y} \leq \mathbf{c}; \mathbf{y} \leq \mathbf{0}^m \} \\ &\leq \text{minimum} \{ \mathbf{c}^T \mathbf{x} \mid \mathbf{A}\mathbf{x} \leq \mathbf{b}; \mathbf{x} \geq \mathbf{0}^n \} \\ &\leq \text{maximum} \{ \mathbf{c}^T \mathbf{x} \mid \mathbf{A}\mathbf{x} \leq \mathbf{b}; \mathbf{x} \geq \mathbf{0}^n \}. \end{aligned}$$

Since equality must hold throughout, the range of $\mathbf{c}^T \mathbf{x}$ is a constant over the primal polyhedron, and $\mathbf{b}^T \mathbf{y}$ is constant over the dual polyhedron, yet \mathbf{c} , \mathbf{A} , and \mathbf{b} are arbitrary. What is wrong in the above line of arguments?

[Note: This and other paradoxes in optimization are found on Harvey Greenberg's page <http://www.cudenver.edu/~hgreenbe/myths/myths.html>.]

Part V

Algorithms

Unconstrained optimization

XI

11.1 Introduction

We consider the unconstrained optimization problem to

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \ f(\mathbf{x}), \quad (11.1)$$

where $f \in C^0$ on \mathbb{R}^n (f is continuous). Mostly, we will assume that $f \in C^1$ holds (f is continuously differentiable), in some cases even $f \in C^2$.

The method of choice for this problem depends on many factors:

- What is the size of the problem (that is, n)?
- Are $\nabla f(\mathbf{x})$ and $\nabla^2 f(\mathbf{x})$ available, and if so at what cost?
- What is the solution requirement? (Do we need a global minimum or a local minimum or simply a stationary point?)
- What are the convexity properties of f ?
- Do we have a good estimate of the location of a stationary point \mathbf{x}^* ? (Can we use locally-only convergent methods?)

We will discuss some basic approaches to the problem (11.1) and refer to questions such as the ones just mentioned during the development.

Example 11.1 (non-linear least squares data fitting) Suppose that we have m data points (t_i, b_i) which we believe are related through an algebraic expression of the form

$$x_1 + x_2 \exp(x_3 t_i) + x_4 \exp(x_5 t_i) = b_i, \quad i = 1, \dots, m,$$

where however the parameters x_1, \dots, x_5 are unknown. (Here, $\exp(x) = e^x$.) In order to best describe the above model, we minimize the total “residual error” given by the norm of the residual

$$f_i(\mathbf{x}) := b_i - [x_1 + x_2 \exp(x_3 t_i) + x_4 \exp(x_5 t_i)], \quad i = 1, \dots, m.$$

Unconstrained optimization

A minimization will then yield the best fit with respect to the data points available. The following then is the resulting optimization problem to be solved:

$$\underset{\mathbf{x} \in \mathbb{R}^5}{\text{minimize}} \ f(\mathbf{x}) := \sum_{i=1}^m |f_i(\mathbf{x})|^2 = \sum_{i=1}^m [f_i(\mathbf{x})]^2.$$

This type of problem is very often solved within numerical analysis and mathematical statistics. Note that the 2-norm is not the only measure of the residual used; sometimes the maximum norm is used. ■

What is the typical form of an algorithm in unconstrained optimization (in fact, for almost every problem class)? Take a look at Figure 11.1 of the level curves¹ of a convex, quadratic function, and the algorithm description below.

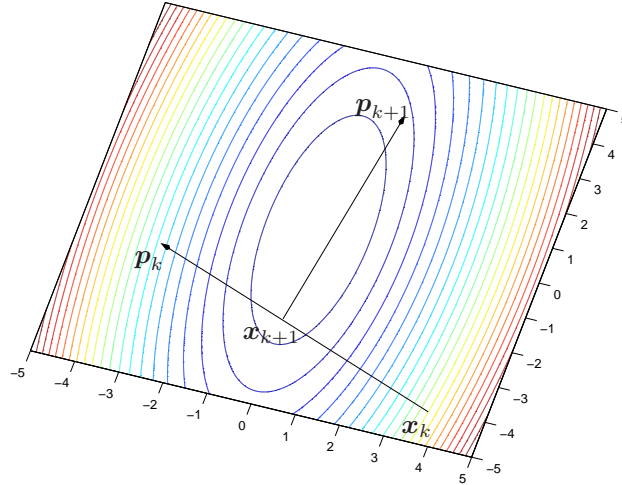


Figure 11.1: At \mathbf{x}_k , the descent direction \mathbf{p}_k is generated. A step α_k is taken in this direction, producing \mathbf{x}_{k+1} . At this point, a new descent direction \mathbf{p}_{k+1} is generated, and so on.

Descent algorithm:

Step 0 (initialization). Determine a *starting point* $\mathbf{x}_0 \in \mathbb{R}^n$. Set $k := 0$.

¹A *level curve* (or, *iso-curve*, or *iso-cost line*) is a set of the form $\{\mathbf{x} \in \mathbb{R}^n \mid f(\mathbf{x}) = k\}$ for a fixed value of $k \in \mathbb{R}$.

- Step 1** (descent direction). Determine a *descent direction* $\mathbf{p}_k \in \mathbb{R}^n$.
- Step 2** (line search). Determine a *step length* $\alpha_k > 0$ such that $f(\mathbf{x}_k + \alpha_k \mathbf{p}_k) < f(\mathbf{x}_k)$ holds.
- Step 3** (update). Let $\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{p}_k$.
- Step 4** (termination check). If a *termination criterion* is fulfilled, then stop! Otherwise, let $k := k + 1$ and go to step 1.

This type of algorithm is inherently local, since we cannot in general use more than the information that can be calculated at the current point \mathbf{x}_k , that is, $f(\mathbf{x}_k)$, $\nabla f(\mathbf{x}_k)$, and $\nabla^2 f(\mathbf{x}_k)$. As far as our local “sight” is concerned, we sometimes call this type of method (for *maximization* problems) the “near-sighted mountain climber,” reflecting the situation in which the mountain climber is in a deep fog and can only check her barometer for the height and feel the steepness of the slope under her feet. Notice then that Figure 11.1 was plotted using several thousands of function evaluations; in reality—and definitely in higher dimension than two—we *never* have this type of orienteering map.

We begin by analyzing Step 1, the most important step of the above-described algorithm. Based on the result in Proposition 4.16 it makes good sense to generate \mathbf{p}_k such that it is a direction of descent.

11.2 Descent directions

11.2.1 Introduction

Recall Definition 4.15 of a direction of descent at a given point. Usually, we have many possible such choices; see for example Proposition 4.16 for a sufficient criterion for a continuously differentiable function. In this section we discuss some details on how descent directions can be generated, depending on a particular situation.

Example 11.2 (example descent directions) (a) Let $f \in C^1(N)$ in some neighborhood N of $\mathbf{x}_k \in \mathbb{R}^n$. If $\nabla f(\mathbf{x}_k) \neq \mathbf{0}^n$, then $\mathbf{p} = -\nabla f(\mathbf{x}_k)$ is a descent direction for f at \mathbf{x}_k (this follows directly from Proposition 4.16). This is the search direction used in the steepest descent method, and it naturally bears the name of *steepest descent direction* because it solves the minimization problem to²

$$\underset{\mathbf{p} \in \mathbb{R}^n: \|\mathbf{p}\|=1}{\text{minimize}} \quad \nabla f(\mathbf{x}_k)^T \mathbf{p}. \quad (11.2)$$

²We have that $\nabla f(\mathbf{x})^T \mathbf{p} = \|\nabla f(\mathbf{x})\| \cdot \|\mathbf{p}\| \cos \theta$, where θ is the angle between the vectors $\nabla f(\mathbf{x})$ and \mathbf{p} ; this expression is clearly minimized by making $\cos \theta = -1$, that is, by letting \mathbf{p} have the angle 180° with $\nabla f(\mathbf{x})$; in other words, $\mathbf{p} = -\nabla f(\mathbf{x})/\|\nabla f(\mathbf{x})\|$.

(b) Let $f \in C^2(N)$ in some neighborhood N of \mathbf{x}_k . If $\nabla f(\mathbf{x}_k) = \mathbf{0}^n$ we cannot use the steepest descent direction anymore. However, we can work with second order information provided by the Hessian to find a descent direction in this case also, provided that f is non-convex at \mathbf{x}_k . Assume that $\nabla^2 f(\mathbf{x}_k)$ is not positive semidefinite (otherwise, \mathbf{x}_k is likely to be a local minimum; see Theorem 4.17). If $\nabla^2 f(\mathbf{x}_k)$ is indefinite we call the stationary point \mathbf{x}_k a *saddle point* of f . Let \mathbf{p} be an eigenvector corresponding to a negative eigenvalue λ of $\nabla^2 f(\mathbf{x}_k)$. Then, we call \mathbf{p} a *direction of negative curvature* for f at \mathbf{x}_k , and \mathbf{p} is a descent direction since for all $\alpha > 0$ small enough, $f(\mathbf{x}_k + \alpha\mathbf{p}) - f(\mathbf{x}_k) = \alpha \nabla f(\mathbf{x}_k)^T \mathbf{p} + \frac{\alpha^2}{2} \mathbf{p}^T \nabla^2 f(\mathbf{x}_k) \mathbf{p} + o(\alpha^2) = \frac{\alpha^2}{2} \lambda \|\mathbf{p}\|^2 + o(\alpha^2) < 0$.

(c) Assume the conditions of (a), and let $\mathbf{Q} \in \mathbb{R}^{n \times n}$ be an arbitrary symmetric, positive definite matrix. Then $\mathbf{p} = -\mathbf{Q} \nabla f(\mathbf{x}_k)$ is a descent direction for f at \mathbf{x}_k : $\nabla f(\mathbf{x}_k)^T \mathbf{p} = -\nabla f(\mathbf{x}_k)^T \mathbf{Q} \nabla f(\mathbf{x}_k) < 0$, due to the positive definiteness of \mathbf{Q} . (This is of course true only if \mathbf{x}_k is non-stationary, as assumed.)

Pre-multiplying by \mathbf{Q} may be interpreted as a scaling of ∇f if we choose a diagonal matrix \mathbf{Q} ; the use of more general matrices is of course possible and leads to exceptionally good computational results for clever choices of \mathbf{Q} . Newton and quasi-Newton methods are based on constructing directions in this way. Note that setting $\mathbf{Q} = \mathbf{I}^n$ (the identity matrix in $\mathbb{R}^{n \times n}$), we obtain the steepest descent direction. ■

To find some arbitrary direction of descent is not a very difficult task as demonstrated by Example 11.2 [in fact, the situation when $\nabla f(\mathbf{x}_k) = \mathbf{0}^n$ appearing in (b) is quite an exotic one already, so typically one can always use directions constructed in (a), or, more generally (c), as descent directions]. However, in order to secure the convergence of numerical algorithms we must provide descent directions that “behave well” numerically. Typical requirements, additional to the basic requirement of being a direction of descent, are:

$$|\nabla f(\mathbf{x}_k)^T \mathbf{p}_k| \geq s_1 \|\nabla f(\mathbf{x}_k)\|^2, \quad \text{and} \quad \|\mathbf{p}_k\| \leq s_2 \|\nabla f(\mathbf{x}_k)\|, \quad (11.3)$$

or

$$-\frac{\nabla f(\mathbf{x}_k)^T \mathbf{p}_k}{\|\nabla f(\mathbf{x}_k)\| \cdot \|\mathbf{p}_k\|} \geq s_1, \quad \text{and} \quad \|\mathbf{p}_k\| \geq s_2 \|\nabla f(\mathbf{x}_k)\|, \quad (11.4)$$

where $s_1, s_2 > 0$, and \mathbf{x}_k and \mathbf{p}_k are, respectively, iterates and search directions of some iterative algorithm.

The purpose of these condition is to prevent the descent directions to deteriorate in quality, in terms of providing good enough descent.

For example, the first condition in (11.3) states that if the directional derivative of f tends to zero then it must be that the gradient of f also tends to zero, while the second condition makes sure that a bad direction in terms of the directional derivative is not compensated by the search direction becoming extremely long in norm. The first condition in (11.4) is equivalent to the requirement that the cosine of the angle between $-\nabla f(\mathbf{x}_k)$ and \mathbf{p}_k is positive and bounded away from zero by the value of s_1 , that is, the angle must be acute and not too close to $\pi/2$; this is another way of saying that the direction \mathbf{p}_k must be steep enough. The purpose of the second condition in (11.4) then is to ensure that if the search direction vanishes then so does the gradient. Methods satisfying (11.3), (11.4) are sometimes referred to as *gradient related*, since they cannot be based on search directions that are very far from those of the steepest descent method.

The choice $\mathbf{p}_k = -\nabla f(\mathbf{x}_k)$ fulfills (11.3), (11.4) with $s_1 = s_2 = 1$.

Another example is as follows: set $\mathbf{p}_k = -\mathbf{Q}_k \nabla f(\mathbf{x}_k)$, where $\mathbf{Q}_k \in \mathbb{R}^{n \times n}$ is a symmetric and *positive definite* matrix such that $m\|\mathbf{s}\|^2 \leq \mathbf{s}^T \mathbf{Q}_k \mathbf{s} \leq M\|\mathbf{s}\|^2$, for all $\mathbf{s} \in \mathbb{R}^n$, holds. [All eigenvalues of \mathbf{Q}_k lie in the interval $[m, M] \subset (0, \infty)$.] Then, the requirement (11.3) is verified with $s_1 = m$, $s_2 = M$, and (11.4) holds with $s_1 = m/M$, $s_2 = m$.

11.2.2 Newton's method and extensions

What should a good descent direction accomplish? Roughly speaking, it should provide as large descent as possible, that is, minimize $f(\mathbf{x} + \mathbf{p}) - f(\mathbf{x})$ over some large enough region of \mathbf{p} around the origin. In principle, this is the idea behind the optimization problem (11.2), because, according to (2.1), $f(\mathbf{x} + \mathbf{p}) - f(\mathbf{x}) \approx \nabla f(\mathbf{x})^T \mathbf{p}$.

Therefore, more insights into how the scaling matrices \mathbf{Q} appearing in Example 11.2(c) should be constructed and, in particular, reasons why the steepest descent direction is not a very wise choice, can be gained if we consider more general approximations than the ones given by (2.1). Namely, assume that $f \in C^1$ near \mathbf{x} , and that for some positive definite matrix \mathbf{Q} it holds that

$$f(\mathbf{x} + \mathbf{p}) - f(\mathbf{x}) \approx \varphi_{\mathbf{x}}(\mathbf{p}) := \nabla f(\mathbf{x})^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \mathbf{Q}^{-1} \mathbf{p}. \quad (11.5)$$

For example, if $f \in C^2$, $\nabla^2 f(\mathbf{x}) \succ \mathbf{0}^{n \times n}$, and assuming that $o(\|\mathbf{p}\|^2) \approx 0$ [cf. (2.3)] we may use $\mathbf{Q}^{-1} = \nabla^2 f(\mathbf{x})$.

Using the optimality conditions, we can easily check that the search direction defined in Example 11.2(c) is a solution to the following optimization problem:

$$\underset{\mathbf{p} \in \mathbb{R}^n}{\text{minimize}} \quad \varphi_{\mathbf{x}}(\mathbf{p}), \quad (11.6)$$

where $\varphi_x(\mathbf{p})$ is defined by (11.5). The closer $\varphi_x(\mathbf{p})$ approximates $f(\mathbf{x} + \mathbf{p}) - f(\mathbf{x})$, the better we can expect the quality of the search directions generated by the method described in Example 11.2(c) to be.

As was already mentioned, setting $\mathbf{Q} = \mathbf{I}^n$, which absolutely fails to take into account any information about f (that is, it is a “one-size-fits-all” approximation), gives us the steepest descent direction. (Cases can easily be constructed such that the algorithm converges extremely slowly; convergence can actually be so bad that the authors of the book [BGLS03] decree that the steepest descent method should be *forbidden*!) On the other hand, the “best” second-order approximation is given by the Taylor expansion (2.3), and therefore we would like to set $\mathbf{Q} = [\nabla^2 f(\mathbf{x})]^{-1}$; this is exactly the choice made in Newton’s method.

Remark 11.3 (a motivation for the descent property in Newton’s method)

The search direction in Newton’s method is based on the solution of the following linear system of equations: find $\mathbf{p} \in \mathbb{R}^n$ such that

$$\nabla_{\mathbf{p}} \varphi_x(\mathbf{p}) := \nabla f(\mathbf{x}) + \nabla^2 f(\mathbf{x}) \mathbf{p} = \mathbf{0}^n.$$

Consider the case of $n = 1$. We should then solve

$$f'(x) + f''(x)p = 0. \quad (11.7)$$

It is obvious that unless $f'(x) = 0$ (whence we are at a stationary point and $p = 0$ solves the equation) we cannot solve (11.7) unless $f''(x) \neq 0$. Then, the solution $\bar{p} := -f'(x)/f''(x)$ to (11.7) is well-defined. We distinguish between two cases:

(a) $f''(x) > 0$. The derivative of the second-order approximation $p \mapsto f'(x)p + \frac{1}{2}f''(x)p^2$ then has a positive slope. Hence, if $f'(x) > 0$ then $\bar{p} < 0$, and if $f'(x) < 0$ then $\bar{p} > 0$ holds. In both cases, therefore, the directional derivative $f'(x)\bar{p} < 0$, that is, \bar{p} is a descent direction.

(b) $f''(x) < 0$. The derivative of the second-order approximation $p \mapsto f'(x)p + \frac{1}{2}f''(x)p^2$ then has a negative slope. Hence, if $f'(x) > 0$ then $\bar{p} > 0$, and if $f'(x) < 0$ then $\bar{p} < 0$ holds. In both cases, therefore, the directional derivative $f'(x)\bar{p} > 0$, that is, \bar{p} is an ascent direction.

From the above it is clear that Newton’s method³ provides the same search direction regardless of whether the optimization problem is a minimization or a maximization problem; the reason is that the search direction is based on the *stationarity* of the second-order approximation and not its minimization/maximization. We also see that the Newton direction \bar{p} is a descent direction if the function f is of the strictly convex

³For $n = 1$ it is often referred to as the *Newton–Raphson method*; cf. Section 4.6.4.2.

type around x [that is, if $f''(x) > 0$], and an ascent direction if it is of the strictly concave type around x [that is, if $f''(x) < 0$]. In other words, if the objective function is (strictly) convex or concave, the Newton equation will give us the right direction, if it gives us a direction at all. In the case when $n > 1$, Newton's method acts as a descent method if the Hessian matrix $\nabla^2 f(\mathbf{x})$ is positive definite, and as an ascent method if it is negative definite, which is appropriate. ■

An essential problem arises if the above-described is *not* what we want; for example, we may be interested in maximizing a function which is neither convex nor concave, and around a current point the function is of strictly convex type (that is, the Hessian is positive definite). In this case the Newton direction will not point in an ascent direction, but instead the opposite. How to solve a problem with a Newton-type method in a non-convex world is the main topic of what follows. As always, we consider minimization to be the direction of interest for f .

So, why might one want to choose a matrix \mathbf{Q} different from the “best” choice $[\nabla^2 f(\mathbf{x})]^{-1}$? There are several reasons:

Lack of positive definiteness The matrix $\nabla^2 f(\mathbf{x})$ may not be positive definite. As a result, the problem (11.6) may even lack solutions and $-\nabla^2 f(\mathbf{x})^{-1} \nabla f(\mathbf{x})$ may in any case not be a descent direction.

This problem can be cured by adding to $\nabla^2 f(\mathbf{x})$ a diagonal matrix \mathbf{E} , so that $\nabla^2 f(\mathbf{x}) + \mathbf{E}$ is positive definite. For example, $\mathbf{E} = \gamma \mathbf{I}^n$, for $-\gamma$ smaller than all the non-positive eigenvalues of $\nabla^2 f(\mathbf{x})$, may be used because such a modification “shifts” the original eigenvalues of $\nabla^2 f(\mathbf{x})$ by $\gamma > 0$. The value of γ needed will automatically be found when solving the “Newton equation” $\nabla^2 f(\mathbf{x}) \mathbf{p} = -\nabla f(\mathbf{x})$, since eigenvalues of $\nabla^2 f(\mathbf{x})$ are pivot elements in Gaussian-elimination procedures. This modification bears the name *Levenberg–Marquardt*.

[Note: as γ becomes large, \mathbf{p} resembles more and more the steepest descent direction.]

Lack of enough differentiability The function f might not be in C^2 , or the matrix of second derivatives might be too costly to compute.

Either being the case, in *quasi-Newton methods* one approximates the Newton equation by replacing $\nabla^2 f(\mathbf{x}_k)$ with a matrix \mathbf{B}_k that is cheaper to compute, typically by only using values of ∇f at the current and some previous points.

Using a first-order Taylor expansion (2.1) for $\nabla f(\mathbf{x}_k)$ we know that

$$\nabla^2 f(\mathbf{x}_k)(\mathbf{x}_k - \mathbf{x}_{k-1}) \approx \nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}_{k-1}),$$

so the matrix \mathbf{B}_k is taken to satisfy the similar system

$$\mathbf{B}_k(\mathbf{x}_k - \mathbf{x}_{k-1}) = \nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}_{k-1}).$$

Notice that for $n = 1$, this corresponds to the *secant method*, in which at iteration k we approximate the second derivative as

$$f''(x_k) \approx \frac{f'(x_k) - f'(x_{k-1})}{x_k - x_{k-1}}.$$

The matrix \mathbf{B}_k has n^2 elements and is hence under-determined by these n equations; additional requirements, such as ones that make sure that \mathbf{B}_k is symmetric and positive definite, result in particular quasi-Newton methods. Typically, starting from $\mathbf{B}_0 = \mathbf{I}^n$, \mathbf{B}_{k+1} is calculated from \mathbf{B}_k using a rank-one or rank-two matrix update; in particular, this allows us to update the factorization of \mathbf{B}_k to efficiently obtain the factorization of \mathbf{B}_{k+1} using standard algorithms in linear algebra.

There are infinitely many choices that may be used, and the following (called the Broyden–Fletcher–Goldfarb–Shanno, or BFGS, method after the original publications [Bro70, Fle70, Gol70, Sha70]) is considered to be the most effective one:

$$\mathbf{B}_{k+1} = \mathbf{B}_k - \frac{(\mathbf{B}_k \mathbf{s}_k)(\mathbf{B}_k \mathbf{s}_k)^T}{\mathbf{s}_k^T \mathbf{B}_k \mathbf{s}_k} + \frac{\mathbf{y}_k \mathbf{y}_k^T}{\mathbf{y}_k^T \mathbf{s}_k},$$

where $\mathbf{s}_k = \mathbf{x}_{k+1} - \mathbf{x}_k$, and $\mathbf{y}_k = \nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k)$. Interestingly enough, should f be quadratic, \mathbf{B}_k will be identical to the Hessian of f after a finite number of steps (namely, n).

Quasi-Newton methods with various updating rules for \mathbf{B}_k are very popular for unconstrained optimization; see Section 11.9.

Computational burden The solution of a linear system $\mathbf{B}_k \mathbf{p}_k = -\nabla f(\mathbf{x}_k)$, or (which is the same if we identify $\mathbf{Q}^{-1} = \mathbf{B}_k$) finding the optimum of (11.6), may be too costly. This is exactly the situation when one would like to use the steepest descent method, which avoids any such calculations.

Other possibilities are: (a) In a quasi-Newton method, keep the matrix \mathbf{B}_k (and, obviously, its factorization) fixed for $k_0 > 1$ subsequent steps; this way, we need only to perform matrix factorization (the most computationally consuming part) every k_0 steps.

(b) Solve the optimization problem (11.6) only approximately, based on the following arguments. Assume that \mathbf{x}_k violates the second-order necessary optimality conditions for f , and consider the problem (11.6), where we replace the matrix \mathbf{Q}^{-1} with an iteration-dependent, perhaps

only positive semidefinite matrix \mathbf{B}_k . As a first example, suppose we consider the Newton method, whence we choose $\mathbf{B}_k = \nabla^2 f(\mathbf{x}_k)$. Then, by the assumption that the second order necessary optimality conditions are violated, $\mathbf{p} = \mathbf{0}^n$ is not a minimum of $\varphi_{\mathbf{x}_k}(\mathbf{p})$ in the problem (11.6). Let $\tilde{\mathbf{p}} \neq \mathbf{0}^n$ be any vector with $\varphi_{\mathbf{x}_k}(\tilde{\mathbf{p}}) < \varphi_{\mathbf{x}_k}(\mathbf{0}^n) = 0$. Then,

$$\varphi_{\mathbf{x}_k}(\tilde{\mathbf{p}}) = \nabla f(\mathbf{x}_k)^T \tilde{\mathbf{p}} + \underbrace{\frac{1}{2} \tilde{\mathbf{p}}^T \mathbf{B}_k \tilde{\mathbf{p}}}_{\geq 0} < \varphi_{\mathbf{x}_k}(\mathbf{0}^n) = 0,$$

which implies that $\nabla f(\mathbf{x}_k)^T \tilde{\mathbf{p}} < 0$. This means that if the Newton equations are solved inexactly, a descent direction is still obtained. This can of course be generalized for quasi-Newton methods as well, since we only assumed that the matrix \mathbf{B}_k is positive semidefinite.

We summarize the above development of search directions in Table 11.1. The iterate is \mathbf{x}_k ; for each algorithm, we describe the linear system solved in order to generate the search direction \mathbf{p}_k . In the table, $\gamma_k \geq 0$ and $\mathbf{B}_k \in \mathbb{R}^{n \times n}$ is a symmetric and positive definite matrix.

Table 11.1: Search directions.

<i>Algorithm</i>	<i>Linear system</i>
Steepest descent	$\mathbf{p}_k = -\nabla f(\mathbf{x}_k)$
Newton's method	$\nabla^2 f(\mathbf{x}_k) \mathbf{p}_k = -\nabla f(\mathbf{x}_k)$
Levenberg–Marquardt	$[\nabla^2 f(\mathbf{x}_k) + \gamma_k \mathbf{I}^n] \mathbf{p}_k = -\nabla f(\mathbf{x}_k)$
Quasi-Newton	$\mathbf{B}_k \mathbf{p}_k = -\nabla f(\mathbf{x}_k)$

11.3 The line search problem

11.3.1 A characterization of the line search problem

Executing Step 2 in the iterative algorithm is naturally done by finding an approximate solution to the one-dimensional problem to

$$\underset{\alpha \geq 0}{\text{minimize}} \quad \varphi(\alpha) := f(\mathbf{x}_k + \alpha \mathbf{p}_k). \quad (11.8)$$

Its optimality conditions are that⁴

$$\varphi'(\alpha^*) \geq 0; \quad \alpha^* \cdot \varphi'(\alpha^*) = 0; \quad \alpha^* \geq 0, \quad (11.10)$$

⁴These conditions are the same as those in Proposition 4.23(b). To establish this fact, let's suppose first that we satisfy (4.10) which here becomes the statement that

$$\varphi'(\alpha^*)(\alpha - \alpha^*) \geq 0, \quad \alpha \geq 0. \quad (11.9)$$

that is,

$$\nabla f(\mathbf{x}_k + \alpha^* \mathbf{p}_k)^T \mathbf{p}_k \geq 0; \quad \alpha^* \cdot \nabla f(\mathbf{x}_k + \alpha^* \mathbf{p}_k)^T \mathbf{p}_k = 0; \quad \alpha^* \geq 0,$$

holds. So, if $\alpha^* > 0$, then $\varphi'(\alpha^*) = 0$ must hold, which therefore means that $\nabla f(\mathbf{x}_k + \alpha^* \mathbf{p}_k)^T \mathbf{p}_k = 0$; that is, the search direction \mathbf{p}_k is orthogonal to the gradient of f at the point $\mathbf{x}_k + \alpha^* \mathbf{p}_k$.

Figure 11.2 shows an example of the one-dimensional function φ along a descent direction with a well-defined minimum.

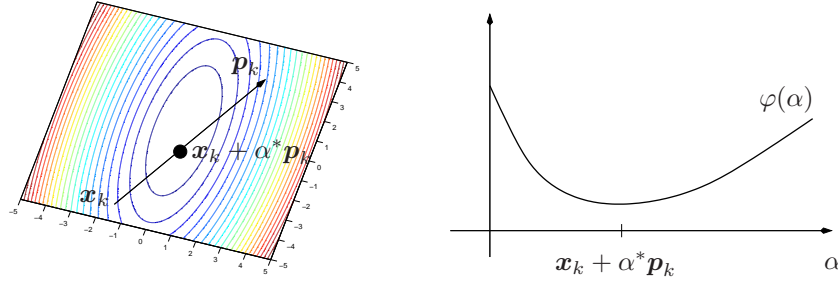


Figure 11.2: A line search in a descent direction.

In the quest for a stationary point it is of relatively minor importance to perform a line search accurately—the stationary point is most probably not situated somewhere along that half-line anyway. Therefore, most line search strategies used in practice are approximate. It should also be noted that if the function f is non-convex then so is probably the case with φ as well, and globally minimizing a non-convex function is difficult even in one variable.

11.3.2 Approximate line search strategies

First, we consider the case where f is quadratic; this is the only general case where an accurate line search is practical.

Let $f(\mathbf{x}) := (1/2)\mathbf{x}^T \mathbf{Q} \mathbf{x} - \mathbf{q}^T \mathbf{x} + a$, where $\mathbf{Q} \in \mathbb{R}^{n \times n}$ is symmetric, $\mathbf{q} \in \mathbb{R}^n$, and $a \in \mathbb{R}$. Suppose we wish to minimize the function φ for this

Setting first $\alpha = 0$ in (11.9), then $\alpha^* \cdot \varphi'(\alpha^*) \leq 0$ follows. On the other hand, setting $\alpha = 2\alpha^*$ in (11.9), then $\alpha^* \cdot \varphi'(\alpha^*) \geq 0$ follows. So, $\alpha^* \cdot \varphi'(\alpha^*) = 0$ must hold. Also, setting $\alpha = \alpha^* + 1$ in (11.9), we obtain that $\varphi'(\alpha^*) \geq 0$. This establishes that (11.10) follows from (4.10). To establish the reverse conclusion and therefore prove that the two conditions are the same, we note that if we satisfy (11.10), then it follows that for every $\alpha \geq 0$, $\varphi'(\alpha^*)(\alpha - \alpha^*) = \alpha \varphi'(\alpha^*) \geq 0$, and we are done.

special case. Then, we can solve the equation $\varphi'(\alpha) = 0$ analytically:

$$\begin{aligned}\varphi'(\alpha) &= \nabla f(\mathbf{x} + \alpha \mathbf{p})^T \mathbf{p} = [\mathbf{Q}(\mathbf{x} + \alpha \mathbf{p}) - \mathbf{q}]^T \mathbf{p} = \alpha \mathbf{p}^T \mathbf{Q} \mathbf{p} - (\mathbf{q} - \mathbf{Q} \mathbf{x})^T \mathbf{p} = 0 \\ &\Leftrightarrow \\ \alpha &= (\mathbf{q} - \mathbf{Q} \mathbf{x})^T \mathbf{p} / \mathbf{p}^T \mathbf{Q} \mathbf{p}.\end{aligned}$$

Let's check the validity and meaning of this solution. We suppose naturally that \mathbf{p} is a descent direction, whence $\varphi'(0) = -(\mathbf{q} - \mathbf{Q} \mathbf{x})^T \mathbf{p} < 0$ holds. Therefore, if \mathbf{Q} is positive definite, we are guaranteed that the value of α will be positive.

Among the classic approximate line searches we mention very briefly the following:

Interpolation Take $f(\mathbf{x}_k), \nabla f(\mathbf{x}_k), \nabla f(\mathbf{x}_k)^T \mathbf{p}_k$ to model a quadratic function approximating f along \mathbf{p}_k . Minimize it by using the analytic formula above.

Newton's method Repeat the improvements gained from a quadratic approximation: $\alpha := \alpha - \varphi'(\alpha) / \varphi''(\alpha)$.

Golden Section The golden section method is a derivative-free method for minimizing *unimodal* functions.⁵ The method reduces an interval wherein the reduction is based only on evaluating φ . The portion left of the length of the previous interval after reduction is $\frac{\sqrt{5}-1}{2} \approx 0.618$.

An approximate line search methodology often used is known as the *Armijo* step length rule. The idea is to quickly generate a step length α which provides a “sufficient” decrease in the value of f . Note that $f(\mathbf{x}_k + \alpha \mathbf{p}_k) \approx f(\mathbf{x}_k) + \alpha \cdot \nabla f(\mathbf{x}_k)^T \mathbf{p}_k$ for very small values of $\alpha > 0$. The requirement of the step length rule is that we get a decrease in the left-hand side of the above approximate relation which is at least a fraction of that predicted in the right-hand side.

Let $\mu \in (0, 1)$ be the fraction of decrease required. Then, the step lengths accepted by the Armijo step length rule are the positive values α which satisfy the inequality

$$\varphi(\alpha) - \varphi(0) \leq \mu \alpha \varphi'(0), \quad (11.11a)$$

that is,

$$f(\mathbf{x}_k + \alpha \mathbf{p}_k) - f(\mathbf{x}_k) \leq \mu \alpha \nabla f(\mathbf{x}_k)^T \mathbf{p}_k. \quad (11.11b)$$

Figure 11.3 illustrates the Armijo step length rule.

⁵ φ is *unimodal* in an interval $[a, b]$ of \mathbb{R} if it has a unique global minimum in $[a, b]$, and is strictly increasing to the left as well as to the right of the minimum. This notion is equivalent to that of φ having a minimum over $[a, b]$ and being strictly quasi-convex there.

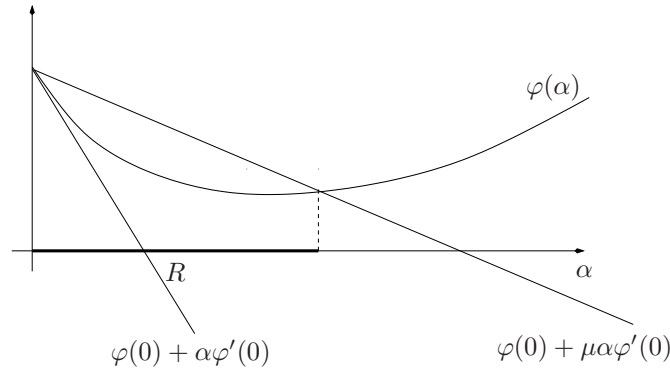


Figure 11.3: The interval R accepted by the Armijo step length rule.

The typical choices are the following: the value of μ is small [$\mu \in (0.001, 0.01)$]; and take $\alpha := 1$. If $\alpha = 1$ does not satisfy the inequality (11.11), then take $\alpha := \alpha/2$, and check the inequality (11.11) again, and so on. The choice of initial trial step $\alpha = 1$ is especially of interest in Newton-type methods, where, locally around a stationary point \mathbf{x}^* where $\nabla^2 f(\mathbf{x}^*)$ is positive definite, local convergence with step length one is guaranteed. (See also Section 4.6.4.2.)

We can select any starting guess $\alpha > 0$ and any fraction $\beta \in (0, 1)$ in place of the choice $\beta = \frac{1}{2}$ made above.

The Armijo condition is satisfied for any sufficiently small step length, provided that $\nabla f(\mathbf{x}_k)^T \mathbf{p}_k < 0$ holds (see Exercise 11.1). In itself it therefore does not guarantee that the next iterate is much better in terms of the objective value than the current one. Often, therefore, it is combined with a condition such that

$$|\varphi'(\alpha_k)| \leq \eta |\varphi'(0)|,$$

that is,

$$|\nabla f(x_k + \alpha p_k)^T p_k| \leq \eta |\nabla f(x_k)^T p_k|,$$

holds for some $\eta \in [0, 1)$. This is called the *Wolfe condition*. A relaxed condition, the *weak Wolfe condition*, of the form

$$\varphi'(\alpha_k) \geq \eta \varphi'(0)$$

is often preferred, since the latter takes less computations to fulfill. The choices $0 < \mu < \eta < 1$ lead to interesting descent algorithms when the Armijo and weak Wolfe conditions are combined, and it is possible (why?) to find positive step lengths that satisfy these two conditions

provided only that f is bounded from below and \mathbf{p}_k is a direction of descent.

11.4 Convergent algorithms

This section presents two basic convergence results for descent methods under different step length rules.

Theorem 11.4 (convergence of a gradient related algorithm) *Suppose that $f \in C^1$, and that for the initial point \mathbf{x}_0 the level set $\text{lev}_f(f(\mathbf{x}_0)) := \{\mathbf{x} \in \mathbb{R}^n \mid f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$ is bounded. Consider the iterative algorithm defined by the description in Section 11.1. In this algorithm, we make the following choices, valid for each iteration k :*

- \mathbf{p}_k satisfies the sufficient descent condition (11.4);
- $\|\mathbf{p}_k\| \leq M$, where M is some positive constant; and
- the Armijo step length rule (11.11) is used.

Then, the sequence $\{\mathbf{x}_k\}$ is bounded, the sequence $\{f(\mathbf{x}_k)\}$ is descending and lower bounded and therefore has a limit, and every limit point of $\{\mathbf{x}_k\}$ is stationary.

Proof. That the sequence $\{\mathbf{x}_k\}$ is bounded follows since the algorithm, as stated, is a descent method, and we assumed that the level set of f at the starting point is bounded; therefore, the sequence of iterates must remain in that set and is therefore bounded.

The rest of the proof is by contradiction. Suppose that $\bar{\mathbf{x}}$ is a limit point of $\{\mathbf{x}_k\}$ but that $\nabla f(\bar{\mathbf{x}}) \neq \mathbf{0}^n$. It is clear that by the continuity of f , $f(\mathbf{x}_k) \rightarrow f(\bar{\mathbf{x}})$. Hence, $f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \rightarrow 0$ must hold. According to the Armijo rule, then, $\alpha_k \nabla f(\mathbf{x}_k)^T \mathbf{p}_k \rightarrow 0$. Here, there are two possibilities. Suppose that $\alpha_k \rightarrow 0$. Then, there must be some iteration \bar{k} after which the initial step length is not accepted by the inequality (11.11), and therefore,

$$f(\mathbf{x}_k + (\alpha_k/\beta)\mathbf{p}_k) - f(\mathbf{x}_k) > \mu(\alpha_k/\beta)\nabla f(\mathbf{x}_k)^T \mathbf{p}_k, \quad k \geq \bar{k}.$$

Dividing both sides by α_k/β we obtain in the limit that

$$(1 - \mu)\nabla f(\bar{\mathbf{x}})^T \mathbf{p}^\infty \geq 0,$$

for any limit point \mathbf{p}^∞ of the bounded sequence $\{\mathbf{p}_k\}$. But in the limit of the inequalities in (11.4) we then clearly reach a contradiction to our claim. So, in fact, we must have that $\alpha_k \not\rightarrow 0$. In this case, then, by the

Unconstrained optimization

above we must have that $\nabla f(\mathbf{x}_k)^T \mathbf{p}_k \rightarrow 0$ holds, so by letting k tend to infinity we obtain that

$$\nabla f(\bar{\mathbf{x}})^T \mathbf{p}^\infty = 0,$$

which again produces a contradiction to the initial claim because of (11.4). We conclude that $\nabla f(\bar{\mathbf{x}}) = \mathbf{0}^n$ must therefore hold. ■

The above proof can be repeated almost in verbatim to establish that *any* step length rule that provides reduction in the value of f that is at least as good as that guaranteed by the Armijo rule will inherit its convergence properties. The main argument is based on the inequality

$$f(\mathbf{x}_{k+1} - \mathbf{x}_k) - f(\mathbf{x}_k) \leq f(\bar{\mathbf{x}}_{k+1}) - f(\mathbf{x}_k) \leq \mu \bar{\alpha}_k \nabla f(\mathbf{x}_k)^T \mathbf{p}_k,$$

where $\bar{\mathbf{x}}_{k+1}$ and $\bar{\alpha}_k$ are the next iterate and step length resulting from the use of the Armijo rule, respectively. If we repeat the arguments in the above proof, replacing α_k with $\bar{\alpha}_k$, we obtain the same contradictions to the condition (11.4). For example, this argument can be used to establish the convergence of gradient related algorithms using exact line searches.

We further note that there is no guarantee that the limit point $\bar{\mathbf{x}}$ is a local minimum; it may also be a *saddle point*, that is, a stationary point where $\nabla^2 f(\bar{\mathbf{x}})$ is indefinite, if it exists.

Another result is cited below from [BeT00]. It allows the Armijo step length rule to be replaced by a much simpler type of step length rule which is also used to minimize a class of non-differentiable functions (cf. Section 6.4). The proof requires the addition of a technical assumption:

Definition 11.5 (Lipschitz continuity) *A C^1 function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is said to have a Lipschitz continuous gradient mapping on \mathbb{R}^n if there exists a scalar $L \geq 0$ such that*

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L \|\mathbf{x} - \mathbf{y}\| \quad (11.12)$$

holds for every $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$. ■

Check that the gradient of a C^2 function f is Lipschitz continuous whenever its Hessian matrix is bounded over \mathbb{R}^n .

Theorem 11.6 (on the convergence of gradient related methods) *Let $f \in C^1$. Consider the sequence $\{\mathbf{x}_k\}$ generated by the formula $\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{p}_k$. Suppose that:*

- ∇f is Lipschitz continuous on \mathbb{R}^n ;
- $c_1 \|\nabla f(\mathbf{x}_k)\|^2 \leq -\nabla f(\mathbf{x}_k)^T \mathbf{p}_k$, $c_1 > 0$;

- $\|p_k\| \leq c_2 \|\nabla f(x_k)\|$, $c_2 > 0$;
- $\alpha_k > 0$ satisfies that $\alpha_k \rightarrow 0$ and $\lim_{k \rightarrow \infty} \sum_{s=1}^k \alpha_s = \infty$.

Then, either $f(x_k) \rightarrow -\infty$ holds, or $f(x_k) \rightarrow \bar{f}$ and $\nabla f(x_k) \rightarrow \mathbf{0}^n$. ■

In Theorem 11.4 convergence is only established in terms of that of subsequences, and the requirements include a level set boundedness condition that can be difficult to check. A strong convergence result is available for the case of convex functions f and the steepest descent method whenever we know that there exists at least one optimal solution; it follows immediately from Theorem 12.4 on the gradient projection method for differentiable optimization over convex sets. In fact, we have already seen such a result in Theorem 6.25 for possibly even non-differentiable convex functions.

Theorem 11.7 (convergence of the steepest descent method under convexity) *Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex and in C^1 on \mathbb{R}^n . Suppose further that the problem (11.1) has at least one optimal solution. Consider the steepest descent algorithm, where the step lengths α_k are determined by the Armijo step length rule. Then, the sequence $\{x_k\}$ converges to some optimal solution to (11.1).* ■

We have so far neglected Step 4 in the algorithm description in Section 11.1 in that we assume in the above results that the sequence $\{x_k\}$ is infinite. A termination criterion must obviously be applied if we are to obtain a result in a finite amount of time. This is the subject of the next section.

11.5 Finite termination criteria

As noted above, convergence to a stationary point is only asymptotic. How does one know when to terminate? A criterion based only on a small size of $\|\nabla f(x_k)\|$ is no good. Why? Because we compare with 0!

The recommendation is the combination of the following:

1. $\|\nabla f(x_k)\| \leq \varepsilon_1(1 + |f(x_k)|)$, $\varepsilon_1 > 0$ small;
2. $f(x_{k-1}) - f(x_k) \leq \varepsilon_2(1 + |f(x_k)|)$, $\varepsilon_2 > 0$ small; and
3. $\|x_{k-1} - x_k\| \leq \varepsilon_3(1 + \|x_k\|)$, $\varepsilon_3 > 0$ small.

The right-hand sides are constructed in order to eliminate some of the possible influences of bad scaling of the variable values, of the objective function, and of the gradient, and also of the possibility that some values are zero at the limit point.

Notice that using the criterion 2. only might mean that we terminate too soon if f is very flat; similarly, using only 3., we terminate prematurely if f is steep around the stationary point we are approaching. The presence of the constant 1 is to remove the dependency of the criterion on the absolute values of f and \mathbf{x}_k , particularly if they are near zero.

We also note that using the $\|\cdot\|_2$ norm may not be good when n is very large: suppose that $\nabla f(\bar{\mathbf{x}}) = (\gamma, \gamma, \dots, \gamma)^T = \gamma(1, 1, \dots, 1)^T$. Then, $\|\nabla f(\bar{\mathbf{x}})\|_2 = \sqrt{n} \cdot \gamma$, which illustrates that the dimension of the problem may enter the norm. Better then is to use the ∞ -norm: $\|\nabla f(\bar{\mathbf{x}})\|_\infty := \max_{1 \leq j \leq n} |\frac{\partial f(\bar{\mathbf{x}})}{\partial x_j}| = |\gamma|$, which does not depend on n .

Norms may have other bad effects. From

$$\begin{aligned}\mathbf{x}_{k-1} &= (1.44453, 0.00093, 0.0000079)^T, \\ \mathbf{x}_k &= (1.44441, 0.00012, 0.0000011)^T;\end{aligned}$$

$$\|\mathbf{x}_{k-1} - \mathbf{x}_k\|_\infty = \|(0.00012, 0.00081, 0.0000068)^T\|_\infty = 0.00081$$

follows. Here, the termination test would possibly pass, although the number of significant digits is very small (the first significant digit is still changing in two components of \mathbf{x} !) Norms emphasize larger elements, so small ones may have bad relative accuracy. This is a case where *scaling* is needed.

Suppose we know that $\mathbf{x}^* = (1, 10^{-4}, 10^{-6})^T$. If, by transforming the space, we obtain the optimal solution $\hat{\mathbf{x}}^* = (1, 1, 1)^T$, then the same relative accuracy would be possible to achieve for all variables. Let then

$$\hat{\mathbf{x}} = \mathbf{D}\mathbf{x}, \quad \text{where} \quad \mathbf{D} := \begin{pmatrix} 1 & 0 & 0 \\ 0 & 10^4 & 0 \\ 0 & 0 & 10^6 \end{pmatrix}.$$

Let $f(\mathbf{x}) := \frac{1}{2}\mathbf{x}^T \mathbf{Q}\mathbf{x} - \mathbf{q}^T \mathbf{x}$, where

$$\mathbf{Q} := \begin{pmatrix} 8 & 3 \cdot 10^4 & 0 \\ 3 \cdot 10^4 & 4 \cdot 10^8 & 10^{10} \\ 0 & 10^{10} & 6 \cdot 10^{12} \end{pmatrix} \quad \text{and} \quad \mathbf{q} := \begin{pmatrix} 11 \\ 8 \cdot 10^4 \\ 7 \cdot 10^6 \end{pmatrix}.$$

Hence, $\mathbf{x}^* = \mathbf{Q}^{-1}\mathbf{q} = (1, 10^{-4}, 10^{-6})^T$.

With $\hat{\mathbf{x}} = \mathbf{D}\mathbf{x}$, we get the transformed problem to minimize $\hat{f}(\hat{\mathbf{x}}) := \frac{1}{2}\hat{\mathbf{x}}^T (\mathbf{D}^{-1}\mathbf{Q}\mathbf{D}^{-1})\hat{\mathbf{x}} - (\mathbf{D}^{-1}\mathbf{q})^T \hat{\mathbf{x}}$, with

$$\mathbf{D}^{-1}\mathbf{Q}\mathbf{D}^{-1} = \begin{pmatrix} 8 & 3 & 0 \\ 3 & 4 & 1 \\ 0 & 1 & 6 \end{pmatrix}; \quad \mathbf{D}^{-1}\mathbf{q} = \begin{pmatrix} 11 \\ 8 \\ 7 \end{pmatrix},$$

and $\hat{\mathbf{x}}^* = (1, 1, 1)^T$. Notice the change in the condition number of the matrix!

The steepest descent algorithm takes only $\nabla f(\mathbf{x})$ into account, not $\nabla^2 f(\mathbf{x})$. Therefore, if the problem is badly scaled, it will suffer from a poor convergence behaviour. Introducing elements of $\nabla^2 f(\mathbf{x})$ into the search direction helps in this respect. This is the precisely the effect of using second-order (Newton-type) algorithms.

11.6 A comment on non-differentiability

The subject of non-differentiable optimization will not be taken up in generality here; it has been analyzed more fully for Lagrangian dual problems in Chapter 6. The purpose of this discussion is to explain, by means of an example, that things can go terribly wrong if we apply methods for the minimization of differentiable functions when the function is non-differentiable.

Rademacher's Theorem states that a Lipschitz continuous function [cf. (11.12) for a statement of the Lipschitz condition for a vector-valued function] automatically is differentiable almost everywhere. It seems to imply that we need not worry about differentiability, because it is very unlikely that a non-differentiable point will be “hit” by mistake. This is certainly true if the subject is simply to pick points at random, but the subject of optimization deals with searching for a particular, *extremal* point in the sense of the objective function, and such points tend to be non-differentiable with a higher probability than zero! Suppose for example that we consider the convex (why?) function

$$f(\mathbf{x}) := \max_{i \in \{1, \dots, m\}} \{\mathbf{c}_i^T \mathbf{x} + b_i\}, \quad \mathbf{x} \in \mathbb{R}^n,$$

that is, a function defined by the point-wise maximum of affine functions. It has the appearance shown in Figure 11.4.

Clearly, the minimum of this function is located at a point where it is non-differentiable.

We next look at a specific problem to which we will apply the method of steepest descent. Suppose that we are given the following convex (why?) objective function:⁶

$$f(x_1, x_2) := \begin{cases} 5(9x_1^2 + 16x_2^2)^{1/2}, & \text{if } x_1 > |x_2|, \\ 9x_1 + 16|x_2|, & \text{if } x_1 \leq |x_2|. \end{cases}$$

For $x_1 > |x_2|$, f is actually continuously differentiable.

⁶This example is due to Wolfe [Wol75].

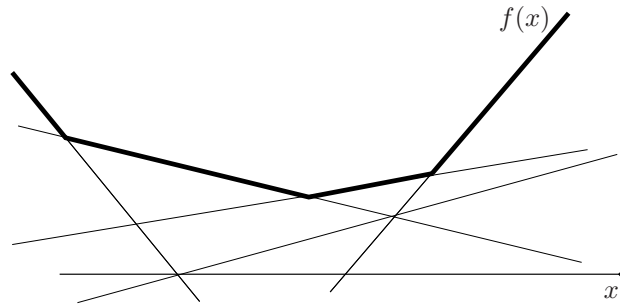


Figure 11.4: A piece-wise affine convex function.

If we start at \mathbf{x}_0 anywhere in the region $x_1 > |x_2| > (9/16)^2|x_1|$ then we obtain a sequence generated by steepest descent with exact line searches that defines a polygonal path with successive orthogonal segments, converging to $\bar{\mathbf{x}} = (0, 0)^T$.

But $\bar{\mathbf{x}}$ is not a stationary point! The reason why it went wrong is that the gradients calculated say very little about the behaviour of f at the limit point $(0, 0)^T$. In fact, f is non-differentiable there. In this example, it in fact holds that $\lim_{x_1 \rightarrow -\infty} f(x_1, 0) = -\infty$, so steepest descent has failed miserably.

In order to resolve this problem, we need to take some necessary measures; the ones below applies to convex functions:

- a) At a non-differentiable point, $\nabla f(x)$ must be replaced by a well-defined extension. Usually, we would replace it with a *subgradient*, that is, one of the vectors that define a supporting hyperplane to the graph of f . At $\bar{\mathbf{x}}$ it is the set defined by the convex hull of the two vectors $(9, 16)^T$ and $(9, -16)^T$.
- b) The step lengths must be chosen differently; exact line searches are clearly forbidden, as we have just seen.

From such considerations, we may develop algorithms that find optima to non-differentiable problems. They are referred to as *subgradient algorithms*, and are analyzed in Section 6.4.

11.7 Trust region methods

Trust region methods use quadratic “models” like Newton-type methods do, but avoiding a line search by instead bounding the length of the search direction, thereby also influencing its direction.

Let $\varphi_k(\mathbf{p}) := f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \nabla^2 f(\mathbf{x}_k) \mathbf{p}$. We say that the model φ_k is *trusted* only in a neighbourhood of $\mathbf{x}_k : \|\mathbf{p}\| \leq \Delta_k$. The use of this bound is apparent when $\nabla^2 f(\mathbf{x}_k)$ is not positive semidefinite. The problem to minimize $\varphi_k(\mathbf{p})$ subject to $\|\mathbf{p}\| \leq \Delta_k$ can be solved (approximately) quite efficiently. The idea is that when $\nabla^2 f(\mathbf{x}_k)$ is badly conditioned, the value of Δ_k should be kept low—thus turning the algorithm more into a steepest descent-like method [recall (11.2)]—while if $\nabla^2 f(\mathbf{x}_k)$ is well conditioned, Δ_k should become large and allow unit steps to be taken. (Prove that the direction of \mathbf{p}_k tends to that of the steepest descent method when $\Delta_k \rightarrow 0$!)

The vector \mathbf{p}_k that solves the trust region problem satisfies $[\nabla^2 f(\mathbf{x}_k) + \gamma_k \mathbf{I}^n] \mathbf{p}_k = -\nabla f(\mathbf{x}_k)$ for some $\gamma_k \geq 0$ such that $\nabla^2 f(\mathbf{x}_k) + \gamma \mathbf{I}^n$ is positive semidefinite. The bounding enforced hence has a similar effect to that of the Levenberg–Marquardt strategy discussed in Section 11.2.2. Provided that the value of Δ_k is small enough, $f(\mathbf{x}_k + \mathbf{p}_k) < f(\mathbf{x}_k)$ holds. Even if $\nabla f(\mathbf{x}_k) = \mathbf{0}^n$ holds, $f(\mathbf{x}_k + \mathbf{p}_k) < f(\mathbf{x}_k)$ if $\nabla^2 f(\mathbf{x}_k)$ is not positive semidefinite. So, progress is made also from stationary points if they are saddle points or local maxima. The robustness and strong convergence characteristics have made trust region methods quite popular.

The update of the trust region size is based on the following measure of similarity between the model φ_k and f : Let

$$\rho_k = \frac{f(\mathbf{x}_k) - f(\mathbf{x}_k + \mathbf{p}_k)}{f(\mathbf{x}_k) - \varphi_k(\mathbf{p}_k)} = \frac{\text{actual reduction}}{\text{predicted reduction}}.$$

If $\rho_k \leq \mu$ let $\mathbf{x}_{k+1} = \mathbf{x}_k$ (unsuccessful step), else
 $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{p}_k$ (successful step).

The value of Δ_k is updated in the following manner, depending on the value of ρ_k :

$$\begin{aligned} \rho_k \leq \mu &\implies \Delta_{k+1} = \frac{1}{2} \Delta_k, \\ \mu < \rho_k < \eta &\implies \Delta_{k+1} = \Delta_k, \\ \rho_k \geq \eta &\implies \Delta_{k+1} = 2 \Delta_k. \end{aligned}$$

Here, $0 < \mu < \eta < 1$, with typical choices being $\mu = \frac{1}{4}$ and $\eta = \frac{3}{4}$; μ is a bound used for deciding when the model can or cannot be trusted even within the region given, while η is used for deciding when the model is good enough to be used in a larger neighbourhood.

Figure 11.5 illustrates the trust region subproblem.

11.8 Conjugate gradient methods

When applied to nonlinear unconstrained optimization problems conjugate direction methods are methods intermediate between the steepest

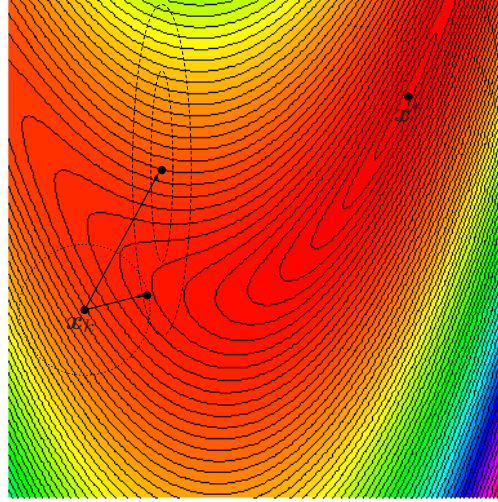


Figure 11.5: Trust region and line search step. The dashed ellipses are two level curves of the quadratic model constructed at \mathbf{x}_k , while the dotted circle is the boundary of the trust region. A step to the minimum of the quadratic model is here clearly inferior to the step taken within the trust region.

descent and Newton methods. The motivation behind them is similar to that for quasi-Newton methods: accelerating the steepest descent method but avoid the evaluation, storage and inversion of the Hessian matrix. They are analyzed for quadratic problems only; extensions to non-quadratic problems utilize that close to an optimal solution every problem is nearly quadratic. Even for non-quadratic problems, the last few decades of developments have resulted in conjugate direction methods being among the most efficient general methodologies available.

11.8.1 Conjugate directions

Definition 11.8 (conjugate direction) *Let $\mathbf{Q} \in \mathbb{R}^{n \times n}$ be symmetric. Two vectors \mathbf{p}_1 and \mathbf{p}_2 in \mathbb{R}^n are \mathbf{Q} -orthogonal, or, conjugate with respect to \mathbf{Q} , if $\mathbf{p}_1^T \mathbf{Q} \mathbf{p}_2 = 0$. ■*

Note that if \mathbf{Q} is the zero matrix then every pair of vectors in \mathbb{R}^n are conjugate; when \mathbf{Q} is the unit matrix, conjugacy reduces to orthogonality. The following result is easy to prove (see Exercise 11.14).

Proposition 11.9 (conjugate vectors are linearly independent) *If $\mathbf{Q} \in \mathbb{R}^{n \times n}$ is symmetric and positive definite and the collection $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_k$ are mutually conjugate with respect to \mathbf{Q} , then they are also linearly independent.* ■

The usefulness of conjugate directions for the quadratic problem to

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{x}) := \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} - \mathbf{q}^T \mathbf{x}, \quad (11.13)$$

where from now on \mathbf{Q} is assumed to be symmetric and positive definite, is clear from the following identification: if the vectors $\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_{n-1}$ are \mathbf{Q} -orthogonal, then Proposition 11.9 implies that there exists a vector $\mathbf{w} \in \mathbb{R}^n$ with

$$\mathbf{x}^* = \sum_{i=0}^{n-1} w_i \mathbf{p}_i; \quad (11.14)$$

multiplying the equation by \mathbf{Q} and scalar multiplying the result by \mathbf{p}_i yields

$$w_i = \frac{\mathbf{p}_i^T \mathbf{Q} \mathbf{x}^*}{\mathbf{p}_i^T \mathbf{Q} \mathbf{p}_i} = \frac{\mathbf{p}_i^T \mathbf{q}}{\mathbf{p}_i^T \mathbf{Q} \mathbf{p}_i}, \quad (11.15)$$

so that

$$\mathbf{x}^* = \sum_{i=0}^{n-1} \frac{\mathbf{p}_i^T \mathbf{q}}{\mathbf{p}_i^T \mathbf{Q} \mathbf{p}_i} \mathbf{p}_i. \quad (11.16)$$

Two ideas are embedded in (11.16): by selecting a proper set of orthogonal vectors \mathbf{p}_i , and by taking the appropriate scalar product all terms but i in (11.14) disappear. This could be accomplished by using any n orthogonal vectors, but (11.15) shows that by making them \mathbf{Q} -orthogonal we can express w_i without knowing \mathbf{x}^* .

11.8.2 Conjugate direction methods

The corresponding conjugate direction method for (11.13) is given by

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k, \quad k = 0, \dots, n-1,$$

where $\mathbf{x}_0 \in \mathbb{R}^n$ is arbitrary and α_k is obtained from an exact line search with respect to f in the direction of \mathbf{p}_k ; cf. (11.8). The principal result about conjugate direction methods is that successive iterates minimize f over a progressively expanding linear manifold that after at most n iterations includes the minimizer of f over \mathbb{R}^n . In other words, defining

$$M_k := \{ \mathbf{x} \in \mathbb{R}^n \mid \mathbf{x} = \mathbf{x}_0 + \text{subspace spanned by } \{\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_k\} \},$$

$$\{\mathbf{x}_{k+1}\} = \arg \min_{\mathbf{x} \in M_k} f(\mathbf{x}) \quad (11.17)$$

holds. To show this, note that by the exact line search rule, for all i ,

$$\left. \frac{\partial f(\mathbf{x}_i + \alpha \mathbf{p}_i)}{\partial \alpha} \right|_{\alpha=\alpha_i} = \nabla f(\mathbf{x}_{i+1})^T \mathbf{p}_i = 0.$$

and for $i = 0, 1, \dots, k-1$,

$$\begin{aligned} \nabla f(\mathbf{x}_{k+1})^T \mathbf{p}_i &= (\mathbf{Q}\mathbf{x}_{k+1} - \mathbf{q})^T \mathbf{p}_i \\ &= \left(\mathbf{x}_{i+1} + \sum_{j=i+1}^k \alpha_j \mathbf{p}_j \right)^T \mathbf{Q}\mathbf{p}_i - \mathbf{q}^T \mathbf{p}_i \\ &= \mathbf{x}_{i+1}^T \mathbf{Q}\mathbf{p}_i - \mathbf{q}^T \mathbf{p}_i \\ &= \nabla f(\mathbf{x}_{i+1})^T \mathbf{p}_i, \end{aligned}$$

where we used the conjugacy of \mathbf{p}_i and \mathbf{p}_j , $j = i+1, \dots, k$. Hence, $\nabla f(\mathbf{x}_{k+1})^T \mathbf{p}_i = 0$ for every $i = 0, 1, \dots, k$, which verifies (11.17).

It is easy to get a picture of what is going on if we look at the case where $\mathbf{Q} = \mathbf{I}^n$ and $\mathbf{q} = \mathbf{0}^n$; since the level curves are circles, minimizing over the n coordinates one by one gives us \mathbf{x}^* in n steps; in each iteration we also identify the optimal value of one of the variables. Conjugate directions in effect does this, although in a transformed space.⁷

The discussion so far has been based on an arbitrary selection of conjugate directions. There are many ways in which conjugate directions could be generated. For example, we could let the vectors \mathbf{p}_i , $i = 0, \dots, n-1$ be defined by the eigenvectors of \mathbf{Q} , as they are mutually orthogonal as well as conjugate with respect to \mathbf{Q} . (Why?) Such a procedure would however be too costly in large-scale applications. The remarkable feature of the conjugate gradient method to be presented below is that the new vector \mathbf{p}_k can be generated directly from the vector \mathbf{p}_{k-1} ; there is no need to remember any of the vectors $\mathbf{p}_0, \dots, \mathbf{p}_{k-2}$, and yet \mathbf{p}_k will be conjugate to them all.

11.8.3 Generating conjugate directions

Given a set of linearly independent vectors $\mathbf{d}_0, \mathbf{d}_1, \dots, \mathbf{d}_k$ we can generate a set of mutually \mathbf{Q} -orthogonal vectors $\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_k$ such that they

⁷Compare this to Newton's method as applied to the problem (11.13); its convergence in one step corresponds to the convergence in one step of the steepest descent method when we first have performed a coordinate transformation such that the level curves become circular.

span the same subspace, by using the Gram–Schmidt procedure. We start the recursion with $\mathbf{p}_0 = \mathbf{d}_0$. Suppose that for some $i < k$ we have $\mathbf{d}_0, \mathbf{d}_1, \dots, \mathbf{d}_i$ such that they span the same subspace as $\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_i$. Then, let \mathbf{p}_{i+1} take the following form:

$$\mathbf{p}_{i+1} = \mathbf{d}_{i+1} + \sum_{m=0}^i c_m^{i+1} \mathbf{p}_m,$$

choosing c_m^{i+1} so that \mathbf{p}_{i+1} is \mathbf{Q} -orthogonal to $\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_i$. This will be true if, for each $j = 0, 1, \dots, i$,

$$\mathbf{p}_{i+1}^T \mathbf{Q} \mathbf{p}_j = \mathbf{d}_{i+1}^T \mathbf{Q} \mathbf{p}_j + \left(\sum_{m=0}^i c_m^{i+1} \mathbf{p}_m \right)^T \mathbf{Q} \mathbf{p}_j = 0.$$

Since $\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_i$ are \mathbf{Q} -orthogonal we have $\mathbf{p}_m^T \mathbf{Q} \mathbf{p}_j = 0$ if $m \neq j$, so

$$c_j^{i+1} = -\frac{\mathbf{d}_{i+1}^T \mathbf{Q} \mathbf{p}_j}{\mathbf{p}_j^T \mathbf{Q} \mathbf{p}_j}, \quad j = 0, 1, \dots, i.$$

Some notes are in order regarding the above development:

1. $\mathbf{p}_j^T \mathbf{Q} \mathbf{p}_j \neq 0$.
2. $\mathbf{p}_{i+1} \neq \mathbf{0}^n$; otherwise it would contradict the linear independence of $\mathbf{d}_0, \mathbf{d}_1, \dots, \mathbf{d}_k$.
3. \mathbf{d}_{i+1} lies in the subspace spanned by $\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_{i+1}$, while \mathbf{p}_{i+1} lies in the subspace spanned by $\mathbf{d}_0, \mathbf{d}_1, \dots, \mathbf{d}_{i+1}$, since these vectors span the same space. Therefore, the subspace identification above is true for $i+1$, and we have shown that the Gram–Schmidt procedure has the property asked for.

11.8.4 Conjugate gradient methods

The conjugate gradient method applies the above Gram–Schmidt procedure to the vectors

$$\mathbf{d}_0 = -\nabla f(\mathbf{x}_0), \quad \mathbf{d}_1 = -\nabla f(\mathbf{x}_1), \quad \dots, \quad \mathbf{d}_{n-1} = -\nabla f(\mathbf{x}_{n-1}).$$

Thus, the conjugate gradient method is to take $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$, where α_k is determined through an exact line search and \mathbf{p}_k is obtained through step k of the Gram–Schmidt procedure to the vector $\mathbf{d}_k = -\nabla f(\mathbf{x}_k)$ and the previous vectors $\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_{k-1}$. In particular,

$$\mathbf{p}_k = -\nabla f(\mathbf{x}_k) + \sum_{j=0}^{k-1} \frac{\nabla f(\mathbf{x}_k)^T \mathbf{Q} \mathbf{p}_j}{\mathbf{p}_j^T \mathbf{Q} \mathbf{p}_j} \mathbf{p}_j. \quad (11.18)$$

Unconstrained optimization

It holds that $\mathbf{p}_0 = -\nabla f(\mathbf{x}_0)$, and termination occurs at step k if $\nabla f(\mathbf{x}_k) = \mathbf{0}^n$; the latter happens exactly when $\mathbf{p}_k = \mathbf{0}^n$. (Why?)

[Note: the search directions are based on negative gradients of f , $-\nabla f(\mathbf{x}_k) = \mathbf{q} - \mathbf{Q}\mathbf{x}_k$, which are identical to the residual in the linear system $\mathbf{Q}\mathbf{x} = \mathbf{q}$ that identifies the optimal solution to (11.13).]

The formula (11.18) can in fact be simplified. The reason is that, because of the successive optimization over subspaces, $\nabla f(\mathbf{x}_k)$ is orthogonal to the subspace spanned by $\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_{k-1}$.

Theorem 11.10 (the conjugate gradient method) *The directions of the conjugate gradient method are generated by*

$$\mathbf{p}_0 = -\nabla f(\mathbf{x}_0); \quad (11.19a)$$

$$\mathbf{p}_k = -\nabla f(\mathbf{x}_k) + \beta_k \mathbf{p}_{k-1}, \quad k = 1, 2, \dots, n-1, \quad (11.19b)$$

where

$$\beta_k = \frac{\nabla f(\mathbf{x}_k)^T \nabla f(\mathbf{x}_k)}{\nabla f(\mathbf{x}_{k-1})^T \nabla f(\mathbf{x}_{k-1})}. \quad (11.19c)$$

Moreover, the method terminates after at most n steps.

Proof. We first use induction to show that the gradients $\nabla f(\mathbf{x}_k)$ are linearly independent. It is clearly true for $k = 0$. Suppose that the method has not terminated after k steps, and that $\nabla f(\mathbf{x}_0), \nabla f(\mathbf{x}_1), \dots, \nabla f(\mathbf{x}_{k-1})$ are linearly independent. Being a conjugate gradient method we know that the subspace spanned by these vectors is the same as that spanned by the vectors $\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_{k-1}$:

$$\text{span}(\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_{k-1}) = \text{span}[\nabla f(\mathbf{x}_0), \nabla f(\mathbf{x}_1), \dots, \nabla f(\mathbf{x}_{k-1})]. \quad (11.20)$$

Now, either $\nabla f(\mathbf{x}_k) = \mathbf{0}^n$, whence the algorithm terminates at the optimal solution, or $\nabla f(\mathbf{x}_k) \neq \mathbf{0}^n$, in which case (by the expanding manifold property) it is orthogonal to $\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_{k-1}$. By (11.20) $\nabla f(\mathbf{x}_k)$ is linearly independent of $\nabla f(\mathbf{x}_0), \nabla f(\mathbf{x}_1), \dots, \nabla f(\mathbf{x}_{k-1})$, completing the induction. Since we have at most n linearly independent vectors in \mathbb{R}^n the algorithm must stop after at most n steps.

The proof is completed by showing that the simplification in (11.19c) is possible. For all j with $\nabla f(\mathbf{x}_j) \neq \mathbf{0}^n$ we have that

$$\nabla f(\mathbf{x}_{j+1}) - \nabla f(\mathbf{x}_j) = \mathbf{Q}(\mathbf{x}_{j+1} - \mathbf{x}_j) = \alpha_j \mathbf{Q}\mathbf{p}_j,$$

and, since $\alpha_j \neq 0$,

$$\begin{aligned}\nabla f(\mathbf{x}_i)^T \mathbf{Q} \mathbf{p}_j &= \frac{1}{\alpha_j} \nabla f(\mathbf{x}_i)^T [\nabla f(\mathbf{x}_{j+1}) - \nabla f(\mathbf{x}_j)] \\ &= \begin{cases} 0, & \text{if } j = 0, 1, \dots, i-2, \\ \frac{1}{\alpha_j} \nabla f(\mathbf{x}_i)^T \nabla f(\mathbf{x}_i), & \text{if } j = i-1, \end{cases}\end{aligned}$$

and also that

$$\mathbf{p}_j^T \mathbf{Q} \mathbf{p}_j = \frac{1}{\alpha_j} \mathbf{p}_j^T [\nabla f(\mathbf{x}_{j+1}) - \nabla f(\mathbf{x}_j)].$$

Substituting these two relations into the Gram–Schmidt formula, we obtain that (11.19b) holds, with

$$\beta_k = \frac{\nabla f(\mathbf{x}_k)^T \nabla f(\mathbf{x}_k)}{\mathbf{p}_{k-1}^T [\nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}_{k-1})]}.$$

From (11.19b) follows that $\mathbf{p}_{k-1} = -\nabla f(\mathbf{x}_{k-1}) + \beta_{k-1} \mathbf{p}_{k-2}$. Using this equation and the orthogonality of $\nabla f(\mathbf{x}_k)$ and $\nabla f(\mathbf{x}_{k-1})$ we can write the denominator in the expression for β_k as desired. We are done. ■

We can deduce further properties of the algorithm. If the matrix \mathbf{Q} has the eigenvalues $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ then we have the following estimate of the distance to the optimal solution after iteration $k+1$:

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\|_Q^2 \leq \left(\frac{\lambda_{n-k} - \lambda_1}{\lambda_{n-k} + \lambda_1} \right)^2 \|\mathbf{x}_0 - \mathbf{x}^*\|_Q^2,$$

where $\|z\|_Q^2 = z^T \mathbf{Q} z$, $z \in \mathbb{R}^n$. What does this estimate tell us about the behaviour of the conjugate gradient algorithm? Suppose that we have a situation where the matrix \mathbf{Q} has m large eigenvalues, and the remaining $n-m$ eigenvalues all are approximately equal to 1. Then the above tells us that after $m+1$ steps of the conjugate gradient algorithm,

$$\|\mathbf{x}_{m+1} - \mathbf{x}^*\|_Q \approx (\lambda_{n-m} - \lambda_1) \|\mathbf{x}_0 - \mathbf{x}^*\|_Q.$$

For a small value of $\lambda_{n-m} - \lambda_1$ this implies that the algorithm gives a good estimate of \mathbf{x}^* already after $m+1$ steps. The conjugate gradient algorithm hence eliminates the effect of the largest eigenvalues first, as the convergence rate after the first $m+1$ steps does not depend on the $m+1$ largest eigenvalues.

The exercises offer additional insight into this convergence theory.

This is in sharp contrast with the convergence rate of the steepest descent algorithm, which is known to be

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\|_Q^2 \leq \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \right)^2 \|\mathbf{x}_k - \mathbf{x}^*\|_Q^2;$$

in other words, the rate of convergence worsens as the condition number of the matrix \mathbf{Q} , $\kappa(\mathbf{Q}) := \lambda_n/\lambda_1$, increases.⁸

Nevertheless, the conjugate gradient method often comes with a *pre-conditioning*, which means that the system solved is not $\mathbf{Q}\mathbf{x} = \mathbf{q}$ but $\mathbf{M}\mathbf{Q}\mathbf{x} = \mathbf{M}\mathbf{q}$ for some invertible square matrix \mathbf{M} , constructed such that the eigenvalues of $\mathbf{M}\mathbf{Q}$ are better clustered than \mathbf{Q} itself. (In other words, the condition number is reduced.)

11.8.5 Extension to non-quadratic problems

Due to the orthogonality of $\nabla f(\mathbf{x}_k)$ and $\nabla f(\mathbf{x}_{k-1})$, we could rewrite (11.19c) as

$$\beta_k = \frac{\nabla f(\mathbf{x}_k)^T [\nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}_{k-1})]}{\nabla f(\mathbf{x}_{k-1})^T \nabla f(\mathbf{x}_{k-1})}. \quad (11.21)$$

The formula (11.19c) is often referred to as the *Fletcher–Reeves formula* (after the paper [FLR64]), while the formula (11.21) is referred to as the *Polak–Ribière formula* (after the paper [PoR69]).

For the quadratic programming problem, the two formulas are identical. However, they would not produce the same sequence of iterates if f were non-quadratic, and the conjugate gradient method has been extended also to such cases. The normal procedure is then to utilize the above algorithm for $k < n$ steps, after which a restart is made at the current iterate using the steepest descent direction; that is, we use the conjugate gradient algorithm several times in succession, in order to not lose conjugacy. The algorithm is not any more guaranteed to terminate after n steps, of course, but the algorithm has been observed to be quite efficient when the objective function and gradient values are cheap to evaluate; especially, this is true when comparing the algorithm class to that of quasi-Newton. (See [Lue84, Ber99] for further discussions on

⁸This type of bound on the convergence rate of the steepest descent method can also be extended to non-quadratic problems: suppose \mathbf{x}^* is the unique optimal solution to the problem of minimizing the C^2 function f and that $\nabla^2 f(\mathbf{x}^*)$ is positive definite. Then, with $0 < \lambda_1 \leq \dots \leq \lambda_n$ being the eigenvalues of $\nabla^2 f(\mathbf{x}^*)$ we have that for all k ,

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \right)^2 [f(\mathbf{x}_k) - f(\mathbf{x}^*)].$$

such computational issues.) It is also remarked in several sources that the Polak–Ribière formula (11.21) is preferable in the non-quadratic case.

11.9 A quasi-Newton method: DFP

As we have already touched upon in Section 11.2.2, most quasi-Newton methods are based on the idea to try to construct the (inverse) Hessian, or an approximation of it, through the use of information gathered in the process of solving the problem; the algorithm then works as a deflected gradient method where the matrix scaling of the negative of the gradient vector is the current approximation of the inverse Hessian matrix.

The BFGS updating formula that was given in Section 11.2.2 is a rank-two update of the Hessian matrix. There are several other versions of the quasi-Newton method, the most popular being based on rank-two updates but of the inverse of the Hessian rather than the Hessian matrix itself. We present one such method below.

The Davidon–Fletcher–Powell algorithm is given in the two papers [Dav59, FLP63]. The algorithm is of interest to us here especially because we can show that through a special choice of matrix update, the quasi-Newton method implemented with an exact line search works exactly like a conjugate gradient method! Moreover, since quasi-Newton methods do not rely on exact line searches for convergence, we learn that quasi-Newton methods are, in this sense, more general than conjugate gradient methods.

The algorithm can be explained like this, as applied to the problem (11.13), the matrix \mathbf{Q} being symmetric and positive definite: start with a symmetric and positive definite matrix $\mathbf{H}_0 \in \mathbb{R}^{n \times n}$, a point $\mathbf{x}_0 \in \mathbb{R}^n$, and with $k = 0$; then set

$$\mathbf{d}_k = -\mathbf{H}_k \nabla f(\mathbf{x}_k); \quad (11.22a)$$

$$\{\alpha_k\} = \arg \min_{\alpha \geq 0} f(\mathbf{x}_k + \alpha \mathbf{d}_k); \quad (11.22b)$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k; \quad (11.22c)$$

$$\mathbf{p}_k = \alpha_k \mathbf{d}_k; \quad (11.22d)$$

$$\mathbf{q}_k = \nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k); \quad (11.22e)$$

$$\mathbf{H}_{k+1} = \mathbf{H}_k + \frac{\mathbf{p}_k \mathbf{p}_k^T}{\mathbf{p}_k^T \mathbf{q}_k} - \frac{(\mathbf{H}_k \mathbf{q}_k)(\mathbf{q}_k^T \mathbf{H}_k)}{\mathbf{q}_k^T \mathbf{H}_k \mathbf{q}_k}; \quad (11.22f)$$

and repeat with $k := k + 1$.

We note that the matrix update in (11.22f) is a rank two update, since the two matrices added to \mathbf{H}_k both are defined by the outer product of a given vector with itself.

Unconstrained optimization

We first demonstrate that the matrices \mathbf{H}_k are positive definite. For any $\mathbf{x} \in \mathbb{R}^n$ we have

$$\mathbf{x}^\top \mathbf{H}_{k+1} \mathbf{x} = \mathbf{x}^\top \mathbf{H}_k \mathbf{x} + \frac{(\mathbf{x}^\top \mathbf{p}_k)^2}{\mathbf{p}_k^\top \mathbf{q}_k} - \frac{(\mathbf{x}^\top \mathbf{H}_k \mathbf{q}_k)^2}{\mathbf{q}_k^\top \mathbf{H}_k \mathbf{q}_k}.$$

Defining $\mathbf{a} = \mathbf{H}_k^{1/2} \mathbf{x}$ and $\mathbf{b} = \mathbf{H}_k^{1/2} \mathbf{q}_k$ we can write this as

$$\mathbf{x}^\top \mathbf{H}_{k+1} \mathbf{x} = \frac{(\mathbf{a}^\top \mathbf{a})(\mathbf{b}^\top \mathbf{b}) - (\mathbf{a}^\top \mathbf{b})^2}{\mathbf{b}^\top \mathbf{b}} + \frac{(\mathbf{x}^\top \mathbf{p}_k)^2}{\mathbf{p}_k^\top \mathbf{q}_k}.$$

We also have that

$$\mathbf{p}_k^\top \mathbf{q}_k = \mathbf{p}_k^\top \nabla f(\mathbf{x}_{k+1}) - \mathbf{p}_k^\top \nabla f(\mathbf{x}_k) = -\mathbf{p}_k^\top \nabla f(\mathbf{x}_k),$$

since

$$\mathbf{p}_k^\top \nabla f(\mathbf{x}_{k+1}) = 0 \quad (11.23)$$

due to the line search being exact. Therefore, by the definition of \mathbf{p}_k ,

$$\mathbf{p}_k^\top \mathbf{q}_k = \alpha_k \nabla f(\mathbf{x}_k)^\top \mathbf{H}_k \nabla f(\mathbf{x}_k),$$

and hence

$$\mathbf{x}^\top \mathbf{H}_{k+1} \mathbf{x} = \frac{(\mathbf{a}^\top \mathbf{a})(\mathbf{b}^\top \mathbf{b}) - (\mathbf{a}^\top \mathbf{b})^2}{\mathbf{b}^\top \mathbf{b}} + \frac{(\mathbf{x}^\top \mathbf{p}_k)^2}{\alpha_k \nabla f(\mathbf{x}_k)^\top \mathbf{H}_k \nabla f(\mathbf{x}_k)}.$$

Both terms in the right-hand side are non-negative, the first because of the Cauchy–Bunyakowski–Schwarz inequality. We must finally show that not both can be zero at the same time. The first term disappears precisely when \mathbf{a} and \mathbf{b} are parallel. This in turn implies that \mathbf{x} and \mathbf{q}_k are parallel, say, $\mathbf{x} = \beta \mathbf{q}_k$ for some $\beta \in \mathbb{R}$. But this would mean that

$$\mathbf{p}_k^\top \mathbf{x} = \beta \mathbf{p}_k^\top \mathbf{q}_k = \beta \alpha_k \nabla f(\mathbf{x}_k)^\top \mathbf{H}_k \nabla f(\mathbf{x}_k) \neq 0,$$

whence $\mathbf{x}^\top \mathbf{H}_{k+1} \mathbf{x} > 0$ holds.

Notice that the fact that the line search is exact is not actually used; it is enough that the α_k chosen yields that $\mathbf{p}_k^\top \mathbf{q}_k > 0$.

The following proposition shows that the Davidon–Fletcher–Powell (DFP) algorithm (11.22) is a conjugate gradient algorithm which provides an optimal solution to (11.13) in at most n steps, when $\mathbf{H}_n = \mathbf{Q}^{-1}$.

Theorem 11.11 (finite convergence of the DFP algorithm) *Consider the algorithm (11.22) for the problem (11.13). Then,*

$$\mathbf{p}_i^\top \mathbf{Q} \mathbf{p}_j = 0, \quad 0 \leq i < j \leq k, \quad (11.24a)$$

$$\mathbf{H}_{k+1} \mathbf{Q} \mathbf{p}_i = \mathbf{p}_i, \quad 0 \leq i \leq k \quad (11.24b)$$

holds.

Proof. We have that

$$\mathbf{q}_k = \nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k) = \mathbf{Q}\mathbf{x}_{k+1} - \mathbf{Q}\mathbf{x}_k = \mathbf{Q}\mathbf{p}_k, \quad (11.25)$$

and

$$\mathbf{H}_{k+1}\mathbf{Q}\mathbf{p}_k = \mathbf{H}_{k+1}\mathbf{q}_k = \mathbf{p}_k, \quad (11.26)$$

the latter from (11.22f).

Proving (11.24) by induction, we see from the above equation that it is true for $k = 0$. Assume (11.24) is true for $k - 1$. We have that

$$\nabla f(\mathbf{x}_k) = \nabla f(\mathbf{x}_{i+1}) + \mathbf{Q}(\mathbf{p}_{i+1} + \cdots + \mathbf{p}_{k-1}).$$

Therefore, from (11.24a) and (11.23),

$$\mathbf{p}_i^T \nabla f(\mathbf{x}_k) = \mathbf{p}_i^T \nabla f(\mathbf{x}_{i+1}) = 0, \quad 0 \leq i < k.$$

Hence from (11.24b)

$$\mathbf{p}_i^T \mathbf{Q}\mathbf{H}_k \nabla f(\mathbf{x}_k) = 0.$$

Thus, since $\mathbf{p}_k = -\alpha_k \mathbf{H}_k \nabla f(\mathbf{x}_k)$ and since $\alpha_k \neq 0$, we obtain

$$\mathbf{p}_i^T \mathbf{Q}\mathbf{p}_k = 0, \quad i < k, \quad (11.27)$$

which proves (11.24a) for k .

Now, since from (11.24b) for $k - 1$, (11.25), and (11.27)

$$\mathbf{q}_k^T \mathbf{H}_k \mathbf{Q}\mathbf{p}_i = \mathbf{q}_k^T \mathbf{p}_i = \mathbf{p}_k^T \mathbf{Q}\mathbf{p}_i = 0, \quad 0 \leq i < k,$$

we have that

$$\mathbf{H}_{k+1}\mathbf{Q}\mathbf{p}_i = \mathbf{H}_k\mathbf{Q}\mathbf{p}_i = \mathbf{p}_i, \quad 0 \leq i < k.$$

This together with (11.26) proves (11.24b) for k . ■

Since the \mathbf{p}_k -vectors are \mathbf{Q} -orthogonal and since we minimize f successively over these directions, the DFP algorithm is a conjugate direction method. Especially, if the initial matrix \mathbf{H}_0 is taken to be the unit matrix, it becomes the conjugate gradient method. In any case, however, convergence is obtained after at most n steps.

Finally, we note that (11.24b) shows that the vectors $\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_k$ are eigenvectors corresponding to unity eigenvalues of the matrix $\mathbf{H}_{k+1}\mathbf{Q}$. These eigenvectors are linearly independent, since they are \mathbf{Q} -orthogonal, and therefore we have that $\mathbf{H}_n = \mathbf{Q}^{-1}$. In other words, with any choice of initial matrix \mathbf{H}_0 (as long as it is symmetric and positive definite) n steps of the 2-rank updates in (11.22f) result in the final matrix being identical to the inverse of the Hessian.

11.10 Convergence rates

The *local convergence rate* is a statement about the speed in which one iteration takes the guess closer to the solution.

Definition 11.12 (local convergence rate) Suppose that $\{\mathbf{x}_k\} \subset \mathbb{R}^n$ and that $\mathbf{x}_k \rightarrow \mathbf{x}^*$. Consider for large k the quotients

$$q_k := \frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|}.$$

(a) [linear convergence rate] The speed of convergence is linear if

$$\limsup_{k \rightarrow \infty} q_k < 1.$$

A linear convergence rate is roughly equivalent to the statement that we get one new correct digit per iteration.

(b) [superlinear convergence rate] The speed of convergence is superlinear if

$$\lim_{k \rightarrow \infty} q_k = 0.$$

(c) [quadratic convergence rate] The speed of convergence is quadratic if

$$\limsup_{k \rightarrow \infty} \frac{q_k}{\|\mathbf{x}_k - \mathbf{x}^*\|} \leq c, \quad c \geq 0.$$

A quadratic convergence rate is roughly equivalent to the statement that the number of correct digits is doubled in every iteration. ■

The steepest descent method has, at most, a linear rate of convergence, moreover often with a constant q_k near unity. Newton-like algorithms have, however, superlinear convergence if $\nabla^2 f(\mathbf{x}^*)$ is positive definite, and even quadratic local convergence can be achieved for Newton's method if $\nabla^2 f$ is Lipschitz continuous in a neighbourhood of \mathbf{x}^* .

11.11 Implicit functions

Suppose that the value of $f(\mathbf{x})$ is given through a simulation procedure, according to Figure 11.6.

If the response $\mathbf{y} = \mathbf{y}(\mathbf{x})$ from the input \mathbf{x} is unknown explicitly, then we cannot differentiate $\mathbf{x} \mapsto f(\mathbf{x}, \mathbf{y}(\mathbf{x}))$ with respect to \mathbf{x} . If, however, we believe that $\mathbf{y}(\cdot)$ is *differentiable*, which means that \mathbf{y} is very stable with respect to changes in \mathbf{x} , then $\nabla_{\mathbf{x}} \mathbf{y}(\mathbf{x})$, and hence $\nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}(\mathbf{x}))$ can be calculated *numerically*. The use of the Taylor expansion technique that follows is only practical if $\mathbf{y}(\mathbf{x})$ is cheap to calculate.

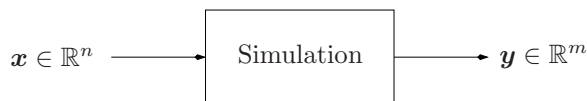


Figure 11.6: A simulation procedure.

Let $\mathbf{e}_i = (0, 0, \dots, 0, 1, 0, \dots, 0)^T$ be the unit vector in \mathbb{R}^n (the only non-zero entry is in position i). Then,

$$\begin{aligned} f(\mathbf{x} + h\mathbf{e}_i) &= f(\mathbf{x}) + h\mathbf{e}_i^T \nabla f(\mathbf{x}) + (h^2/2)\mathbf{e}_i^T \nabla^2 f(\mathbf{x})\mathbf{e}_i + \dots \\ &= f(\mathbf{x}) + h\partial f(\mathbf{x})/\partial x_i + (h^2/2)\partial^2 f(\mathbf{x})/\partial x_i^2 + \dots \end{aligned}$$

So, for small $h > 0$,

$$\begin{aligned} \frac{\partial f(\mathbf{x})}{\partial x_i} &\approx \frac{f(\mathbf{x} + h\mathbf{e}_i) - f(\mathbf{x})}{h}, & \text{(forward difference)} \\ \frac{\partial f(\mathbf{x})}{\partial x_i} &\approx \frac{f(\mathbf{x} + h\mathbf{e}_i) - f(\mathbf{x} - h\mathbf{e}_i)}{2h}. & \text{(central difference)} \end{aligned}$$

The value of h is typically set to a function of the machine precision; if chosen too large, we get a bad approximation of the partial derivative, while a too small value might result in numerical cancellation.

The *automatic differentiation* technique exploits the inherent structure of most practical functions, that they almost always are evaluated through a sequence of elementary operations. Automatic differentiation represents this structure in the form of a computational graph; when forming partial derivatives this graph is utilized in the design of chain rules for the automatic derivative calculation. In applications to simulation models, this means that differentiation is performed within the simulation package, thereby avoiding some of the computational cost and the potential instability inherent in difference formulas.

11.12 Notes and further reading

The material of this chapter is classic; text books covering similar material in more depth include [OrR70, DeS83, Lue84, Fle87, BSS93, BGLS03]. Line search methods were first developed by Newton [New1687], and the steepest descent method is due to Cauchy [Cau1847]. The Armijo rule is due to Armijo [Arm66], and the Wolfe condition is due to Wolfe [Wol69]. The classic book by Brent [Bre73] analyzes algorithms that do not use derivatives, especially line search methods.

Rademacher's Theorem [Rad19] states that a Lipschitz continuous function is differentiable everywhere except on sets of Lebesgue measure zero. The Lipschitz condition is due to Lipschitz [Lip1877]. Algorithms for the minimization of non-differentiable convex functions are given in [Sho85, HiL93, Ber99, BGLS03].

Trust region methods are given a thorough treatment in the book [CGT00]. The material on the conjugate gradient and BFGS methods was collected from [Lue84, Ber99]; another good source is [NoW99].

A popular class of algorithms for problems with an implicit objective function is the class of *pattern search methods*. With such algorithms the search for a good gradient-like direction is replaced by calculations of the objective function along directions specified by a pattern of possible points. For an introduction to the field, see [KLT03].

Automatic differentiation is covered in the monograph [Gri00].

11.13 Exercises

Exercise 11.1 (well-posedness of the Armijo rule) Through an argument by contradiction, establish the following: If $f \in C^1$, $\mathbf{x}_k \in \mathbb{R}^n$ and $\mathbf{p}_k \in \mathbb{R}^n$ satisfies $\nabla f(\mathbf{x}_k)^T \mathbf{p}_k < 0$, then for every choice of $\mu \in (0, 1)$ there exists $\bar{\alpha} > 0$ such that every $\alpha \in (0, \bar{\alpha}]$ satisfies (11.11). In other words, which ever positive first trial step length α we choose, we will find a step length that satisfies (11.11) in a finite number of trials.

Exercise 11.2 (descent direction) Investigate whether the direction of $\mathbf{p} = (2, -1)^T$ is a direction of descent with respect to the function $f(\mathbf{x}) := x_1^2 + x_1x_2 - 4x_2^2 + 10$ at $\mathbf{x} := (1, 1)^T$.

Exercise 11.3 (Newton's method) Suppose that you wish to solve the unconstrained problem to minimize $f(\mathbf{x})$ over $\mathbf{x} \in \mathbb{R}^n$, where $f \in C^2(\mathbb{R}^n)$. You are naturally interested in using Newton's method (with line searches).

(a) At some iteration you get the error message, "Step length is zero." Which reason(s) can there be for such a message?

(b) At some iteration you get the error message, "Search direction does not exist." Which reason(s) can there be for such a message?

(c) Describe at least one means to modify Newton's method such that neither of the above two error messages will ever appear.

Exercise 11.4 (steepest descent) Is it possible to reach the (unique) optimal solution to the problem to

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \ f(\mathbf{x}) := (x_1 - 2)^2 + 5(x_2 + 6)^2$$

by the use of the steepest descent algorithm, if we first perform a variable substitution? If so, perform it and thus find the optimal solution.

Exercise 11.5 (steepest descent) Consider the problem to

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{x}) := (2x_1^2 - x_2)^2 + 3x_1^2 - x_2.$$

- (a) Perform one iteration of the steepest descent method using an exact line search, starting at $\mathbf{x}_0 := (1/2, 5/4)^T$.
- (b) Is the function convex around \mathbf{x}_1 ?
- (c) Will it converge to a global optimum? Why/why not?

Exercise 11.6 (Newton's method with exact line search) Consider the problem to

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{x}) := (x_1 + 2x_2 - 3)^2 + (x_1 - 2)^2.$$

- (a) Start from $\mathbf{x}_0 := (0, 0)^T$, and perform one iteration of Newton's method with an exact line search.
- (b) Are there any descent directions from \mathbf{x}_1 ?
- (c) Is \mathbf{x}_1 optimal? Why/why not?

Exercise 11.7 (Newton's method with Armijo line search) Consider the problem to

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{x}) := \frac{1}{2}(x_1 - 2x_2)^2 + x_1^4.$$

- (a) Start from $\mathbf{x}_0 := (2, 1)^T$, and perform one iteration of Newton's method with the Armijo rule, using the fraction requirement $\mu = 0.1$.
- (b) Determine the values of $\mu \in (0, 1)$ such that the step length $\alpha = 1$ will be accepted.

Exercise 11.8 (Newton's method for nonlinear equations) Suppose the function $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is continuously differentiable and consider the following system of nonlinear equations:

$$\mathbf{f}(\mathbf{x}) = \mathbf{0}^n.$$

Newton's method for the solution of unconstrained optimization problems has its correspondence for the above problem.

Given an iterate \mathbf{x}_k we construct the following linear approximation of the nonlinear function:

$$\mathbf{f}(\mathbf{x}_k) + \nabla \mathbf{f}(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k) = \mathbf{0}^n,$$

or, equivalently,

$$\nabla \mathbf{f}(\mathbf{x}_k)\mathbf{x} = \nabla \mathbf{f}(\mathbf{x}_k)\mathbf{x}_k - \mathbf{f}(\mathbf{x}_k),$$

where

$$\nabla \mathbf{f}(\mathbf{x}) = \begin{pmatrix} \nabla f_1(\mathbf{x})^T \\ \nabla f_2(\mathbf{x})^T \\ \vdots \\ \nabla f_n(\mathbf{x})^T \end{pmatrix}$$

Unconstrained optimization

is the *Jacobian* of \mathbf{f} at \mathbf{x} . Assuming that $\nabla \mathbf{f}(\mathbf{x})$ is non-singular, this linear system has a unique solution which defines the new iterate, \mathbf{x}_{k+1} , that is,

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \nabla \mathbf{f}(\mathbf{x}_k)^{-1} \mathbf{f}(\mathbf{x}_k).$$

(One can show that if \mathbf{f} satisfies some additional requirements, this sequence of iterates will converge to a solution to the original nonlinear system, either from any starting point—global convergence—or from a point sufficiently close to a solution—local convergence.)

(a) Consider the nonlinear system

$$\mathbf{f}(x_1, x_2) = \begin{pmatrix} f_1(x_1, x_2) \\ f_2(x_1, x_2) \end{pmatrix} = \begin{pmatrix} 2(x_1 - 2)^3 + x_1 - 2x_2 \\ 4x_2 - 2x_1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}.$$

Perform one iteration of the above algorithm, starting from $\mathbf{x}_0 = (1, 0)^T$. Calculate the value of

$$\|\mathbf{f}(x_1, x_2)\| = \sqrt{f_1(x_1, x_2)^2 + f_2(x_1, x_2)^2}$$

both at \mathbf{x}_0 and \mathbf{x}_1 . (Observe that $\|\mathbf{f}(\mathbf{x})\| = 0$ if and only if $\mathbf{f}(\mathbf{x}) = \mathbf{0}^n$, whence the values of $\|\mathbf{f}(\mathbf{x}_k)\|$, $k = 1, 2, \dots$, can be used as a measure of convergence of the iterates.)

(b) Explain why the above method generalizes Newton's method for unconstrained optimization to a larger class of problems.

Exercise 11.9 (over-determined linear equations) Consider the problem to

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2,$$

where $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $\mathbf{b} \in \mathbb{R}^m$. Assume that $m \geq n$ and that $\text{rank } \mathbf{A} = n$.

(a) Write down the necessary optimality conditions for this problem. Are they also sufficient for global optimality? Why/why not?

(b) Write down the globally optimal solution in closed form.

Exercise 11.10 (sufficient descent conditions) Consider the first sufficient descent condition in (11.4). Why does it have that form, and why is the alternative form

$$-\nabla f(\mathbf{x}_k)^T \mathbf{p}_k \geq s_1$$

not acceptable?

Exercise 11.11 (Newton's method under affine transformations) Suppose we make the following change of variables: $\mathbf{y} := \mathbf{A}\mathbf{x} + \mathbf{b}$, where $\mathbf{A} \in \mathbb{R}^{n \times n}$ is invertible. Show that Newton's method is invariant to such changes of variables.

Exercise 11.12 (Levenberg–Marquardt, exam 990308) Consider the unconstrained optimization problem to

$$\text{minimize } f(\mathbf{x}) := \mathbf{q}^T \mathbf{x} + \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x}, \quad (11.28a)$$

$$\text{subject to } \mathbf{x} \in \mathbb{R}^n, \quad (11.28b)$$

where $\mathbf{Q} \in \mathbb{R}^{n \times n}$ is symmetric and positive semidefinite but not positive definite. We attack the problem through a Levenberg–Marquardt strategy, that is, we utilize a Newton-type method where a multiple $\gamma > 0$ of the unit matrix is added to the Hessian of f (that is, to the matrix \mathbf{Q}) in order to guarantee that the (modified) Newton equation is uniquely solvable. (See Section 11.2.2.) This implies that, given an iteration point \mathbf{x}_k , the search direction \mathbf{p}_k is determined by solving the linear system

$$\begin{aligned} [\nabla^2 f(\mathbf{x}_k) + \gamma \mathbf{I}^n] \mathbf{p} &= -\nabla f(\mathbf{x}_k) && \Longleftrightarrow \\ [\mathbf{Q} + \gamma \mathbf{I}^n] \mathbf{p} &= -(\mathbf{Q} \mathbf{x}_k + \mathbf{q}). \end{aligned} \quad (11.29)$$

(a) Consider the formula

$$\mathbf{x}_{k+1} := \mathbf{x}_k + \mathbf{p}_k, \quad k = 0, 1, \dots, \quad (11.30)$$

that is, the algorithm that is obtained by utilizing the Newton-like search direction \mathbf{p}_k from (11.29) and the step length 1 in every iteration. Show that this iterative step is the same as that to let \mathbf{x}_{k+1} be given by the solution to the problem to

$$\text{minimize } f(\mathbf{y}) + \frac{\gamma}{2} \|\mathbf{y} - \mathbf{x}_k\|^2, \quad (11.31a)$$

$$\text{subject to } \mathbf{y} \in \mathbb{R}^n. \quad (11.31b)$$

(b) Suppose that an optimal solution to (11.28) exists. Suppose also that the sequence $\{\mathbf{x}_k\}$ generated by the algorithm (11.30) converges to a point \mathbf{x}^∞ . (This can actually be shown to hold.) Show that \mathbf{x}^∞ is optimal in (11.28).

[Note: This algorithm is in fact a special case of the *proximal point algorithm*. Suppose that f is a convex function on \mathbb{R}^n and the variables are constrained to a non-empty, closed and convex set $S \subseteq \mathbb{R}^n$.

We extend the iteration formula (11.31) to the following:

$$\text{minimize } f(\mathbf{y}) + \frac{\gamma_k}{2} \|\mathbf{y} - \mathbf{x}_k\|^2, \quad (11.32a)$$

$$\text{subject to } \mathbf{y} \in S, \quad (11.32b)$$

where $\{\gamma_k\} \subset (0, 2)$ is a sequence of positive numbers that is bounded away from zero, and where \mathbf{x}_{k+1} is taken as the unique vector \mathbf{y} solving (11.32). If an optimal solution exists, it is possible to show that the sequence given by (11.32) converges to a solution. See [Pat98, Ber99] for overviews of this class of methods. (It is called “proximal point” because of the above interpretation: that the next iterate is close, proximal, to the previous one.)]

Unconstrained optimization

Exercise 11.13 (unconstrained optimization algorithms, exam 980819) Consider the unconstrained optimization problem to minimize $f(\mathbf{x})$ over $\mathbf{x} \in \mathbb{R}^n$, where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is in C^1 . Let $\{\mathbf{x}_k\}$ be a sequence of iteration points generated by some algorithm for solving this problem, and suppose that it holds that $\nabla f(\mathbf{x}_k) \rightarrow \mathbf{0}^n$, that is, the gradient value tends to zero (which of course is a favourable behaviour of the algorithm). The question is what this means in terms of the convergence of the more important sequence $\{\mathbf{x}_k\}$.

Consider therefore the sequence $\{\mathbf{x}_k\}$, and also the sequence $\{f(\mathbf{x}_k)\}$ of function values. Given the assumption that $\nabla f(\mathbf{x}_k) \rightarrow \mathbf{0}^n$, is it true that $\{\mathbf{x}_k\}$ and/or $\{f(\mathbf{x}_k)\}$ converges or are even bounded? Provide every possible case in terms of the convergence of these two sequences, and give examples, preferably simple ones for $n = 1$.

Exercise 11.14 (conjugate directions) Prove Proposition 11.9.

Exercise 11.15 (conjugate gradient method) Apply the conjugate gradient method to the system $\mathbf{Q}\mathbf{x} = \mathbf{q}$, where

$$\mathbf{Q} = \begin{pmatrix} 2 & 1 & 0 \\ 1 & 2 & 1 \\ 0 & 1 & 2 \end{pmatrix} \quad \text{and} \quad \mathbf{q} = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}.$$

Exercise 11.16 (convergence of the conjugate gradient method, I) In the conjugate gradient method, prove that the vector \mathbf{p}_i can be written as a linear combination of the set of vectors $\{\mathbf{q}, \mathbf{Q}\mathbf{q}, \mathbf{Q}^2\mathbf{q}, \dots, \mathbf{Q}^i\mathbf{q}\}$. Also prove that \mathbf{x}_{i+1} minimizes the quadratic function $\mathbb{R}^n \ni \mathbf{x} \mapsto f(\mathbf{x}) := \frac{1}{2}\mathbf{x}^T \mathbf{Q}\mathbf{x} - \mathbf{q}^T \mathbf{x}$ over all the linear combinations of these vectors.

Exercise 11.17 (convergence of the conjugate gradient method, II) Use the result of the previous problem to establish that the conjugate gradient method converges in a number of iterations equal to the number of distinct eigenvalues of the matrix \mathbf{Q} .

Optimization over convex sets

XII

12.1 Feasible direction methods

Consider the problem to

$$\text{minimize } f(\mathbf{x}), \quad (12.1a)$$

$$\text{subject to } \mathbf{x} \in X, \quad (12.1b)$$

where $X \subseteq \mathbb{R}^n$ is a nonempty, closed and convex set and $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is a C^1 function on X .

In almost all cases, algorithms for this problem will not be based on staying feasible ($\mathbf{x}_k \in X$) but rather to reach feasibility and optimality at the same time. Why? If X is defined by (convex) inequalities of the form “ $g_i(\mathbf{x}) \leq b_i$,” where g_i is nonlinear, then checking, for example, whether \mathbf{p} is a feasible direction at \mathbf{x} , or what the maximum feasible step from \mathbf{x} in the direction of \mathbf{p} is, is very difficult. For example, in the latter case, for which step length $\alpha > 0$ does it happen that $g_i(\mathbf{x} + \alpha\mathbf{p}) = b_i$? This is a nonlinear equation!

The notable exception is when X is a polyhedral set, which we from here on in this chapter will assume is the case. How to characterize a feasible direction in the polyhedral case has already been analyzed in Example 4.22.

A general framework of feasible-direction methods for the problem (12.1) can be extended from the unconstrained world as follows (notice the difference to the description given in Section 11.1):

Feasible descent algorithm:

Step 0 (initialization). Determine a *starting point* $\mathbf{x}_0 \in \mathbb{R}^n$ such that $\mathbf{x}_0 \in X$. Set $k := 0$.

- Step 1** (feasible descent direction). Determine a *search direction* $\mathbf{p}_k \in \mathbb{R}^n$ such that \mathbf{p}_k is a feasible direction.
- Step 2** (line search). Determine a *step length* $\alpha_k > 0$ such that $f(\mathbf{x}_k + \alpha_k \mathbf{p}_k) < f(\mathbf{x}_k)$ and $\mathbf{x}_k + \alpha_k \mathbf{p}_k \in X$.
- Step 3** (update). Let $\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{p}_k$.
- Step 4** (termination check). If a *termination criterion* is fulfilled, then stop! Otherwise, let $k := k + 1$ and go to step 1.

This type of algorithm is local, just as in the unconstrained case, and there are more difficulties compared to the unconstrained case, associated with generating search directions \mathbf{p}_k : they must simultaneously be feasible and provide descent. The linearity of the constraints however ensures that it is possible to solve approximations of the original problem, such as approximations based on a Taylor expansion of f around the current iterate \mathbf{x}_k , or to determine the active constraint set at \mathbf{x}_k and generating a feasible search direction in a neighbourhood determined by them. Moreover, we also need to determine a maximum step length in the line search, and the termination criteria need to be different from those in the unconstrained case because the gradient of f need not be zero at a stationary point.

In the following, we will analyze three natural algorithms for the solution of the problem (12.1). The first two, the Frank–Wolfe and simplicial decomposition algorithms, are based on generating search directions by solving *linear* problems, while the third one, the gradient projection algorithm, corresponds to solving a more difficult, yet related, convex *quadratic* problem. In each of these cases, the algorithms are derived from the necessary optimality conditions associated with the problem (12.1), which can be found in Section 4.4.

We will establish convergence for a general case, where all that can be guaranteed is that any limit point of the sequence $\{\mathbf{x}_k\}$ is stationary. We will however also establish what can be achieved in addition when the problem is *convex*. It is then not only the (obvious) case that every limit point then is globally optimal, but we can on occasion prove something stronger. In the Frank–Wolfe method, we can utilize simpler step length rules than the Armijo rule; in the simplicial decomposition method, we can establish finite convergence even when previous information is discarded; and in the gradient projection method, we can prove convergence to an optimal solution. The latter has interesting consequences for iterative methods, like the steepest descent algorithm, in the unconstrained case, which it generalizes to the constrained case.

12.2 The Frank–Wolfe algorithm

Consider the problem (12.1). We suppose that $X \subset \mathbb{R}^n$ is a bounded polyhedron, for the simplicity of the presentation. (See, however, Exercise 12.1.) We suppose further that $f \in C^1(\mathbb{R}^n)$.

The *Frank–Wolfe algorithm* works as follows:

Frank–Wolfe algorithm:

Step 0 (initialization). Generate the starting point $\mathbf{x}_0 \in X$ (for example by letting it be any extreme point in X). Set $k := 0$.

Step 1 (feasible descent direction). Solve the problem to

$$\underset{\mathbf{y} \in X}{\text{minimize}} \quad z_k(\mathbf{y}) := \nabla f(\mathbf{x}_k)^T(\mathbf{y} - \mathbf{x}_k). \quad (12.2)$$

Let \mathbf{y}_k be a solution to this LP problem, and $\mathbf{p}_k := \mathbf{y}_k - \mathbf{x}_k$ be the search direction.

Step 2 (line search). Approximately solve the one-dimensional problem to minimize $f(\mathbf{x}_k + \alpha \mathbf{p}_k)$ over $\alpha \in [0, 1]$. Let α_k be the resulting step length.

Step 3 (update). Let $\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{p}_k$.

Step 4 (termination check). If, for example, $z_k(\mathbf{y}_k)$ or α_k is close to zero, then terminate! Otherwise, let $k := k + 1$ and go to Step 1.

It should be clear that \mathbf{p}_k is a feasible direction: note that we can write the new point \mathbf{x}_{k+1} as $\alpha \mathbf{y}_k + (1 - \alpha) \mathbf{x}_k$ for some $\alpha \in [0, 1]$, that is, we construct an optimal convex combination of two feasible points.

In Step 2, we can utilize the Armijo step length rule or a more accurate line search such as one based on a quadratic interpolation.

The LP problem (12.2) was introduced already in Section 4.4, cf., for example, (4.14). It was also there shown that $z_k(\mathbf{y}_k) < 0$ holds if and only if \mathbf{x}_k is not stationary. Hence, we have shown that the Frank–Wolfe algorithm is a descent algorithm.

We have also shown in Section 4.4 that a lower bound on the optimal value f^* is available whenever f is convex: the first termination criterion in Step 4 then states that this lower bound is close enough to f^* .

Theorem 12.1 (convergence of the Frank–Wolfe algorithm) *Suppose that $X \subseteq \mathbb{R}^n$ is a nonempty, bounded polyhedron, and that the function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is in C^1 on X . Suppose that in Step 2 of the Frank–Wolfe algorithm we use the Armijo step length rule. Then, the sequence $\{\mathbf{x}_k\}$ is bounded, and every limit point (at least one exists) is stationary.*

If f is convex on X , then every limit point is globally optimal.

Proof. The boundedness of X ensures that the sequences $\{\mathbf{x}_k\}$ and $\{\mathbf{y}_k\}$ are bounded. Let \mathbf{x}^∞ and \mathbf{y}^∞ be limit points, corresponding to some subsequence \mathcal{K} . Also the sequence $\{\mathbf{p}_k\}$ is bounded; let further $\mathbf{p}^\infty := \mathbf{y}^\infty - \mathbf{x}^\infty$. We first show that

$$\{\nabla f(\mathbf{x}_k)^\top \mathbf{p}_k\}_{k \in \mathcal{K}} \rightarrow 0 \quad (12.3)$$

holds. As $\{f(\mathbf{x}_k)\}$ is descending, it must hold that its limit in \mathcal{K} is $f(\mathbf{x}^\infty)$. Hence, $\{f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k)\}_{k \in \mathcal{K}} \rightarrow 0$. Further,

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) \leq \mu \alpha_k \nabla f(\mathbf{x}_k)^\top \mathbf{p}_k < 0, \quad k = 0, 1, \dots, \quad (12.4)$$

which means that either $\{\nabla f(\mathbf{x}_k)^\top \mathbf{p}_k\}_{k \in \mathcal{K}} \rightarrow 0$ (whence we are done) or $\{\alpha_k\}_{k \in \mathcal{K}} \rightarrow 0$ holds. In the latter case, there must be an index κ such that for every $k \geq \kappa$ in \mathcal{K} the initial step length is not accepted by the Armijo rule, that is,

$$f(\mathbf{x}_k + (\alpha_k/\beta)\mathbf{p}_k) - f(\mathbf{x}_k) > \mu(\alpha_k/\beta)\nabla f(\mathbf{x}_k)^\top \mathbf{p}_k, \quad k \geq \kappa, \quad k \in \mathcal{K}.$$

Dividing both sides of this inequality by α_k/β , in the limit in \mathcal{K} of the resulting inequality we reach the conclusion that $\nabla f(\mathbf{x}^\infty)^\top \mathbf{p}^\infty \geq 0$, while in the limit in \mathcal{K} of the descent inequality $\nabla f(\mathbf{x}_k)^\top \mathbf{p}_k < 0$, the reverse inequality is obtained. We conclude that (12.3) follows.

From the above follows that $\nabla f(\mathbf{x}^\infty)^\top (\mathbf{y}^\infty - \mathbf{x}^\infty) = 0$, whence stationarity follows, since in the limit in \mathcal{K} of the characterization

$$\nabla f(\mathbf{x}_k)^\top (\mathbf{y}_k - \mathbf{x}_k) \leq \nabla f(\mathbf{x}_k)^\top (\mathbf{y} - \mathbf{x}_k), \quad \mathbf{y} \in X,$$

we obtain that

$$\nabla f(\mathbf{x}^\infty)^\top (\mathbf{y} - \mathbf{x}^\infty) \geq \nabla f(\mathbf{x}^\infty)^\top (\mathbf{y}^\infty - \mathbf{x}^\infty) = 0, \quad \mathbf{y} \in X.$$

Since the limit point was arbitrarily chosen, the first result follows.

The second part of the theorem follows from Theorem 4.24. ■

Figure 12.1 illustrates the LP problem in Step 1 at a non-stationary point \mathbf{x}_k , the resulting extreme point \mathbf{y}_k , and search direction \mathbf{p}_k .

We have above established the result for the Armijo rule. By applying the same technique as that discussed after Theorem 11.4 for gradient related methods in unconstrained optimization, we can also establish convergence to stationary points under the use of exact line searches.

Under additional technical assumptions we can establish that the sequence $\nabla f(\mathbf{x}_k)^\top \mathbf{p}_k \rightarrow 0$; see Exercise 12.2.

Under the assumption that f is convex, several additional techniques for choosing the step lengths are available; see the notes for references. We refer to one such choice below.

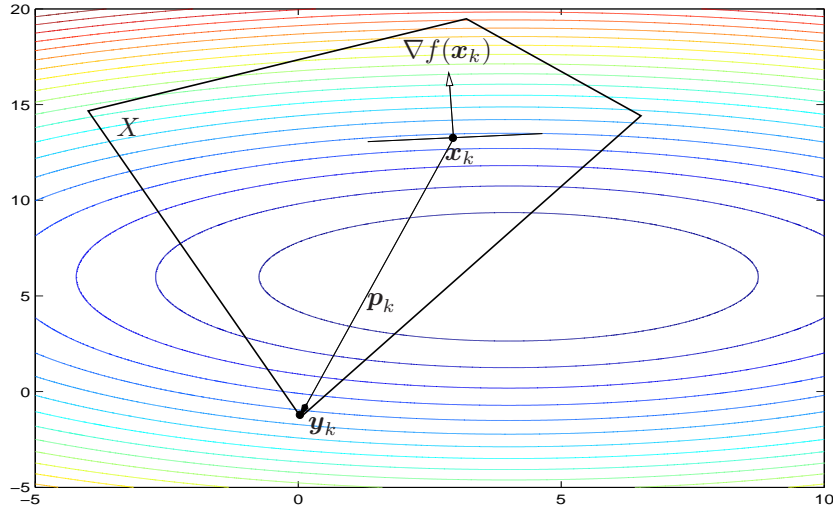


Figure 12.1: Step 1 of the Frank–Wolfe algorithm.

Theorem 12.2 (convergence of the Frank–Wolfe algorithm) *Suppose that $X \subset \mathbb{R}^n$ is nonempty, convex, closed, and bounded and that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex and in C^1 on X .*

(a) *In the Frank–Wolfe algorithm, suppose the step lengths $\alpha_k \in (0, 1]$ satisfy that, for some $C > 0$ and every k large enough,¹*

$$\alpha_k \leq C/k, \quad (12.5a)$$

$$1 - \alpha_{k+1} = \frac{\alpha_{k+1}}{\alpha_k}. \quad (12.5b)$$

If the sequence $\{x_k\}$ is finite, then the last iterate solves (12.1). Otherwise, $f(x_k) \rightarrow f^$, and the sequence $\{x_k\}$ converges to the set of solutions to (12.1): $\text{dist}_{X^*}(x_k) \rightarrow 0$. In particular, any limit point of $\{x_k\}$ solves (12.1).*

(b) *Suppose that ∇f is Lipschitz continuous on X . In the Frank–Wolfe algorithm, suppose the step lengths $\alpha_k \in (0, 1]$ are chosen according to the quadratically convergent divergent step length rule (6.41), (6.42). Then, the conclusions in (a) hold. ■*

¹According to this step length rule, $\alpha_k \approx 1/k$ for large k .

12.3 The simplicial decomposition algorithm

Consider the problem (12.1) under the same conditions as stated in Section 12.2. The *simplicial decomposition algorithm* builds on the Representation Theorem 3.22.

In the below description we let \mathcal{P} denote the set of extreme points of X . We also denote by \mathcal{P}_k a subset of the extreme points which have been generated prior to iteration k and which are kept in memory; an element of this set is denoted by \mathbf{y}^i in order to not mix these extreme points with the vectors \mathbf{y}_k solving the LP problem.

The simplicial decomposition algorithm works as follows:

Simplicial decomposition algorithm:

Step 0 (initialization). Generate the starting point $\mathbf{x}_0 \in X$ (for example by letting it be any extreme point in X). Set $k := 0$. Let $\widehat{\mathcal{P}}_0 = \mathcal{P}_0 := \emptyset$. Let $\bar{\mathbf{x}}_0 = \mathbf{x}_0$.

Step 1 (feasible descent direction). Let \mathbf{y}_k be a solution to the LP problem (12.2).

Let $\mathcal{P}_{k+1} := \widehat{\mathcal{P}}_k \cup \{k\}$, for some subset $\widehat{\mathcal{P}}_k$ of \mathcal{P}_k .

Step 2 (multidimensional line search). Let $\boldsymbol{\nu}_{k+1}$ be an approximate solution to the *restricted master problem* to

$$\underset{\boldsymbol{\nu}}{\text{minimize}} \quad f \left(\bar{\mathbf{x}}_k + \sum_{i \in \mathcal{P}_{k+1}} \nu_i (\mathbf{y}^i - \bar{\mathbf{x}}_k) \right), \quad (12.6a)$$

$$\text{subject to} \quad \sum_{i \in \mathcal{P}_{k+1}} \nu_i \leq 1, \quad (12.6b)$$

$$\nu_i \geq 0, \quad i \in \mathcal{P}_{k+1}, \quad (12.6c)$$

where $\bar{\mathbf{x}}_k \in X_k := \text{conv}(\{\bar{\mathbf{x}}_{k-1}\} \cup \{\mathbf{y}^i \mid i \in \mathcal{P}_k\})$.

Step 3 (update). Let $\mathbf{x}_{k+1} := \bar{\mathbf{x}}_k + \sum_{i \in \mathcal{P}_{k+1}} (\boldsymbol{\nu}_{k+1})_i (\mathbf{y}^i - \bar{\mathbf{x}}_k)$.

Step 4 (termination check). If, for example, $z_k(\mathbf{y}_k)$ is close to zero, or if $\mathcal{P}_{k+1} = \mathcal{P}_k$, then terminate! Otherwise, let $k := k + 1$ and go to Step 1.

The description of the simplicial decomposition algorithm does not completely specify the sets \mathcal{P}_k or the points $\bar{\mathbf{x}}_k$. We say that we use *column dropping* in Step 1 of the algorithm when $\widehat{\mathcal{P}}_k \subset \mathcal{P}_k$. Some classic column dropping rules are given in Table 12.1.

To begin with, suppose that we use the principle in (a). According to this principle, we run the algorithm by adding one new extreme point to the previous set of extreme points known so far, solve the problem to

Table 12.1: Some column dropping principles for Step 1

-
- (a) [no column dropping]: For all k , $\widehat{\mathcal{P}}_k := \mathcal{P}_k$, and $\bar{\mathbf{x}}_k := \mathbf{x}_k$.
 (b) [Zero weight column dropping]: For $k \geq 1$,

$$\widehat{\mathcal{P}}_k := \{ i \in \mathcal{P}_k \mid (\boldsymbol{\nu}_k)_i > 0 \}.$$

For all k , $\bar{\mathbf{x}}_k := \mathbf{x}_0$.

- (c) [bounded size of \mathcal{P}_k]: Let r be a positive integer. For $k \geq 1$, let

$$\widehat{\mathcal{P}}_k := \{ i \in \mathcal{P}_k \mid (\boldsymbol{\nu}_k)_i > 0 \}.$$

If $|\widehat{\mathcal{P}}_k| < r$, then let $\bar{\mathbf{x}}_k := \bar{\mathbf{x}}_{k-1}$. If $|\widehat{\mathcal{P}}_k| = r$, then let

$$\widehat{\mathcal{P}}_k := \widehat{\mathcal{P}}_k \setminus i_k^*, \quad i_k^* \in \arg \min_{i \in \widehat{\mathcal{P}}_k} \{ (\boldsymbol{\nu}_{k-1})_i \},$$

where ties are broken arbitrarily, and $\bar{\mathbf{x}}_k := \mathbf{x}_k$.

- (d) [Frank–Wolfe]: For all k , $\widehat{\mathcal{P}}_k := \emptyset$ and $\bar{\mathbf{x}}_k := \mathbf{x}_k$.
-

minimize f over the convex hull of them, and repeat until we either get close enough to a stationary point or if the last LP did not give us a new extreme point. (In the latter case we are at a stationary point! Why?)

Suppose instead that we drop every extreme point that got a zero weight in the last restricted master problem, that is, we work according to the principle in (b). We then remove all the extreme points that we believe will not be useful in order to describe the optimal solution as a convex combination of them.

The algorithm corresponding to the principle in (c) is normally called the *restricted simplicial decomposition* algorithm; it allows us to drop extreme points in order to keep the memory requirements below a certain threshold. In order to do so, we may need to also throw away an extreme point that had a positive weight at the optimal solution to the previous restricted master problem, and we implement this by removing one with the least weight.

The most extreme case of the principle in (c) is to throw away every point that was previously generated, and keep only the most recent one. (It corresponds to letting $r = 1$.) Then, according to the principle in (d), we are back at the Frank–Wolfe algorithm!

The restricted master problem (12.6) does not contain the slack variable associated with the convexity weight for the vector $\bar{\mathbf{x}}_k$. Introducing it as $\mu \geq 0$, we obtain an equivalent statement of the problem:

$$\begin{aligned} \underset{(\mu, \nu)}{\text{minimize}} \quad & f \left(\mu \bar{\mathbf{x}}_k + \sum_{i \in \mathcal{P}_{k+1}} \nu_i \mathbf{y}^i \right), \end{aligned} \quad (12.7a)$$

$$\text{subject to} \quad \mu + \sum_{i \in \mathcal{P}_{k+1}} \nu_i = 1, \quad (12.7b)$$

$$\mu, \nu_i \geq 0, \quad i \in \mathcal{P}_{k+1}. \quad (12.7c)$$

We then recognize the feasible set of the restricted master problem as the particular explicit convex hull that it actually is.

The advantage of using an inner representation of X and to algorithmically improve it is that it is *much* simpler to deal with in an optimization algorithm than the linear constraints that define X above.

The disadvantage of the representation is that the set \mathcal{P} of extreme points is both very large for a large-scale problem, and it is not known; compare with the case of the simplex method, where we cannot simply enumerate the extreme points in order to then pick the best one. In this nonlinear case, we may also need a large number of them in order to span an optimal solution. The trick that makes the algorithm work well is the column dropping in Step 1, which keeps down the size of the problems to be solved.

The simplicial decomposition algorithm is quite similar to the Frank–Wolfe algorithm (notice this from (12.7) when $|\mathcal{P}_{k+1}| = 1$). The main, and very important, difference, is that the Frank–Wolfe algorithm drops all the previous extreme points visited, and only optimizes over line segments. In simplicial decomposition the information is kept, and thanks to this extra information the algorithm can make much better progress in each iteration. The method therefore becomes much more efficient than the Frank–Wolfe algorithm in practice.

As far as convergence is concerned, the algorithm clearly does at least as well as the Frank–Wolfe algorithm does. In fact, if we use the principle (a) in Table 12.1 then we always add an extreme point in each iteration, and hence we must arrive at a stationary point in a finite number of iterations since \mathcal{P} is finite. If we use the principle in (b) finite convergence is still ensured under mild additional conditions on f , while for finite convergence under the principle in (c) the value of r is crucial: it must be at least as large as the number of extreme points needed to describe the optimal solution. In other words, we must have memory enough to be able to span \mathbf{x}^* . If the value is too small, then the algorithm cannot converge finitely, and then the behaviour can be as

bad as in the Frank–Wolfe algorithm, which is rather bad indeed. (The convergence rate cannot even be linear.)

We finally note that in both of the Frank–Wolfe and simplicial decomposition algorithms, the sequence $\{f(\mathbf{x}_k)\}$ is decreasing provided that the line search is exact enough; on the other hand, the sequence $\{z_k(\mathbf{y}_k)\}$ is non-monotone in general, although the limit (almost) always is zero. We refer to Theorem 12.1 for the corresponding result for the Frank–Wolfe algorithm. The result for simplicial decomposition is the same, except for the cases when convergence to a stationary point is finite, as discussed above.²

12.4 The gradient projection algorithm

As was observed in Exercise 4.5 the result of the operation

$$\mathbf{y} := \text{Proj}_X[\mathbf{x} - \nabla f(\mathbf{x})]$$

at $\mathbf{x} \in X$ is that $\mathbf{y} = \mathbf{x}$ if and only if \mathbf{x} is a stationary point and if it is not then $\mathbf{p} := \mathbf{y} - \mathbf{x}$ defines a descent direction with respect to f at \mathbf{x} . This characterization is true also if we introduce an arbitrary scalar $\alpha > 0$ as follows:

$$\mathbf{y} := \text{Proj}_X[\mathbf{x} - \alpha \nabla f(\mathbf{x})].$$

(Why?) In other words, supposing that \mathbf{x}_k is not stationary, we generate the next iteration point as follows:

$$\mathbf{x}_{k+1} := \text{Proj}_X[\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)], \quad k = 1, \dots, \quad (12.8)$$

where $\{\alpha_k\} \subset \mathbb{R}_{++}$. As the vector $\mathbf{p}_k := \text{Proj}_X[\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)] - \mathbf{x}_k$ defines a feasible descent direction with respect to f at \mathbf{x}_k it follows that $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$ provided that the value of α_k is sufficiently small. The *gradient projection algorithm* is based on this observation. How do we choose α_k and what does it mean to perform a line search in α_k ?

We propose here to utilize the Armijo rule, which was introduced for unconstrained optimization in (11.11). It is however modified such that the trial points are not $\mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k)$ as in the steepest descent method, but the projected, feasible, points $\text{Proj}_X[\mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k)]$. Starting with

²The line segment representing the feasible set in the line search problem in Step 2 of the Frank–Wolfe algorithm satisfies $[x_k, y_k] \subset X_{k+1}$, that is, the restricted master problem in the simplicial decomposition algorithm is always defined over a set that is at least as large. As a consequence, the latter algorithm always will be able to achieve an improvement in the value of the objective function that is at least as great as that of the former. From this observation it is relatively easy to establish a basic convergence result along the lines of Theorem 12.1.

a trial step $\bar{\alpha} > 0$, we check the Armijo criterion in (11.11) for the feasible point $\text{Proj}_X[\mathbf{x}_k - \bar{\alpha}\nabla f(\mathbf{x}_k)]$, and then replace $\bar{\alpha}$ by $\bar{\alpha}\beta$ for some $\beta \in (0, 1)$ if it is not satisfied, and so on, until the Armijo criterion is satisfied. Eventually, then, we will satisfy the following inequality:

$$f(\text{Proj}_X[\mathbf{x}_k + \alpha_k \mathbf{p}_k]) - f(\mathbf{x}_k) \leq \mu \alpha_k \nabla f(\mathbf{x}_k)^\top \mathbf{p}_k.$$

The resulting step length then is $\alpha_k = \bar{\alpha}\beta^i$ for some integer $i \geq 0$ (zero if the initial step is accepted, otherwise positive), and the new iteration point is the last point projected, $\mathbf{x}_{k+1} := \text{Proj}_X[\mathbf{x}_k - (\bar{\alpha}\beta^i)\nabla f(\mathbf{x}_k)]$.

Consider Figure 12.2. It illustrates a case where we imagine that the initial step $\bar{\alpha}$ has to be reduced twice (here, $\beta = \frac{1}{2}$) before the step is accepted. As we can see from the figure the “line search” is not really a line search, since the feasible points checked rather follow a piece-wise linear curve than a line; in this example we trace the boundary of X , and we sometimes refer to this type of line search as a *boundary search* or a search *along the projection arc*.

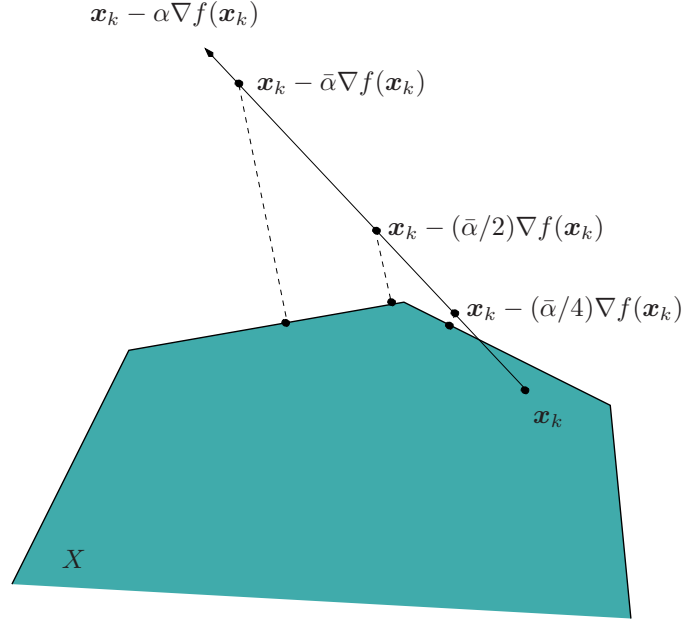


Figure 12.2: Trial steps in the boundary search.

Although the technique looks more complex than the use of the Armijo rule for the steepest descent method, their convergence behaviour

are the same. Theorem 11.4 on the convergence of gradient related methods in unconstrained optimization can be extended to the case of the gradient projection method, to state the following (see Exercise 12.6):

Theorem 12.3 (convergence of a gradient projection algorithm) *Suppose that $X \subseteq \mathbb{R}^n$ is nonempty, closed, and convex, and that the function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is in C^1 on X . Suppose further that for the starting point \mathbf{x}_0 the level set $\text{lev}_f(f(\mathbf{x}_0))$ intersected with X is bounded. Consider the iterative algorithm defined by the iteration (12.8), where the step length α_k is determined by the Armijo step length rule along the projection arc. Then, the sequence $\{\mathbf{x}_k\}$ is bounded, the sequence $\{f(\mathbf{x}_k)\}$ is descending, lower bounded and therefore has a limit, and every limit point of $\{\mathbf{x}_k\}$ is stationary. ■*

The following theorem shows that the gradient projection method has a much stronger convergence property in the convex case.

Theorem 12.4 (convergence of a gradient projection algorithm) *Suppose that $X \subseteq \mathbb{R}^n$ is nonempty, closed, and convex, and that the function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is in C^1 on X . Suppose further that f is convex and that the problem (12.1) has at least one optimal solution. Consider the iterative algorithm defined by the iteration (12.8), where the step length α_k is determined by the Armijo step length rule along the projection arc. Then, $\mathbf{x}_k \rightarrow \mathbf{x}^*$ for some optimal solution \mathbf{x}^* to (12.1).*

Proof. If $\{\mathbf{x}_k\}$ is finite, then the stopping criterion implies that the last iterate is optimal. Suppose therefore that the sequence is infinite.

Let \mathbf{x}^* be an arbitrary optimal solution to (12.1). We have that

$$\|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2 + \|\mathbf{x}_k - \mathbf{x}^*\|^2 - \|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 = 2(\mathbf{x}_{k+1} - \mathbf{x}_k)^T(\mathbf{x}^* - \mathbf{x}_k). \quad (12.9)$$

Further, from the variational inequality characterization of the projection resulting in \mathbf{x}_{k+1} follows that

$$\begin{aligned} 0 &\leq [\mathbf{x}_{k+1} - \mathbf{x}_k + \alpha_k \nabla f(\mathbf{x}_k)]^T(\mathbf{x}^* - \mathbf{x}_{k+1}) \\ &= [\mathbf{x}_{k+1} - \mathbf{x}_k + \alpha_k \nabla f(\mathbf{x}_k)]^T(\mathbf{x}^* - \mathbf{x}_k) \\ &\quad + [\mathbf{x}_{k+1} - \mathbf{x}_k + \alpha_k \nabla f(\mathbf{x}_k)]^T(\mathbf{x}_k - \mathbf{x}_{k+1}), \end{aligned}$$

which yields

$$\begin{aligned}
 (\mathbf{x}_{k+1} - \mathbf{x}_k)^\top (\mathbf{x}^* - \mathbf{x}_k) &\geq \alpha_k \nabla f(\mathbf{x}_k)^\top (\mathbf{x}_k - \mathbf{x}^*) \\
 &\quad - [\mathbf{x}_{k+1} - \mathbf{x}_k + \alpha_k \nabla f(\mathbf{x}_k)]^\top (\mathbf{x}_k - \mathbf{x}_{k+1}) \\
 &\geq \alpha_k [f(\mathbf{x}_k) - f(\mathbf{x}^*)] \\
 &\quad + [\mathbf{x}_{k+1} - \mathbf{x}_k + \alpha_k \nabla f(\mathbf{x}_k)]^\top (\mathbf{x}_{k+1} - \mathbf{x}_k) \\
 &\geq \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2 + \alpha_k \nabla f(\mathbf{x}_k)^\top (\mathbf{x}_{k+1} - \mathbf{x}_k),
 \end{aligned} \tag{12.10}$$

where we used the convexity characterization in Theorem 3.40(a) in the second inequality, and the optimality of \mathbf{x}^* in the third. Combining (12.9) and (12.10) now yields that

$$\begin{aligned}
 \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2 + \|\mathbf{x}_k - \mathbf{x}^*\|^2 - \|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 &\geq 2[\|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2 \\
 &\quad + \alpha_k \nabla f(\mathbf{x}_k)^\top (\mathbf{x}_{k+1} - \mathbf{x}_k)].
 \end{aligned}$$

Rearranging terms yields

$$\begin{aligned}
 \|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 &\leq \|\mathbf{x}_k - \mathbf{x}^*\|^2 - \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2 - 2\alpha_k \nabla f(\mathbf{x}_k)^\top (\mathbf{x}_{k+1} - \mathbf{x}_k) \\
 &\leq \|\mathbf{x}_k - \mathbf{x}^*\|^2 - 2\alpha_k \nabla f(\mathbf{x}_k)^\top (\mathbf{x}_{k+1} - \mathbf{x}_k) \\
 &= \|\mathbf{x}_k - \mathbf{x}^*\|^2 + \varepsilon_k,
 \end{aligned} \tag{12.11}$$

where

$$\varepsilon_k := 2\alpha_k \nabla f(\mathbf{x}_k)^\top (\mathbf{x}_k - \text{Proj}_X[\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)]), \quad k = 0, 1, \dots$$

Note that by the descent property of the algorithm, $\varepsilon_k \geq 0$ for all k .

In view of the Armijo rule,

$$\mu \nabla f(\mathbf{x}_k)^\top (\mathbf{x}_k - \text{Proj}_X[\mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)]) \leq f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}).$$

Combining the above two inequalities,

$$\varepsilon_k \leq \frac{2\alpha_k}{\mu} [f(\mathbf{x}_k) - f(\mathbf{x}_{k+1})] \leq \frac{2\bar{\alpha}}{\mu} [f(\mathbf{x}_k) - f(\mathbf{x}_{k+1})]. \tag{12.12}$$

By (12.12),

$$\sum_{k=0}^{\infty} \varepsilon_k \leq \frac{2\bar{\alpha}}{\mu} [f(\mathbf{x}_0) - f(\mathbf{x}^*)] < \infty.$$

We say that the sequence $\{\varepsilon_k\}$ is *summable*. The consequence for the inequality (12.11) will become apparent from the following lemma.

Lemma 12.5 (quasi-Fejér convergence) *Let $S \subset \mathbb{R}^n$ be nonempty and $\{\mathbf{a}_k\} \subset \mathbb{R}^n$ be a sequence such that for all $\mathbf{x} \in S$ and all k ,*

$$\|\mathbf{a}_{k+1} - \mathbf{x}\|^2 \leq \|\mathbf{a}_k - \mathbf{x}\|^2 + \varepsilon_k, \quad (12.13)$$

where $\{\varepsilon_k\}$ is a summable sequence in \mathbb{R}_+ .

- (a) Then, $\{\mathbf{a}_k\}$ is bounded; and
- (b) if a limit point $\bar{\mathbf{a}}$ of $\{\mathbf{a}_k\}$ belongs to S then $\mathbf{a}_k \rightarrow \bar{\mathbf{a}}$.

Proof. (a) Fix some $\mathbf{x} \in S$. Applying (12.13) iteratively yields, for some $C \in \mathbb{R}_+$,

$$\|\mathbf{a}_k - \mathbf{x}\|^2 \leq \|\mathbf{a}_0 - \mathbf{x}\|^2 + \sum_{j=0}^{k-1} \varepsilon_j \leq \|\mathbf{a}_0 - \mathbf{x}\|^2 + \sum_{j=0}^{\infty} \varepsilon_j \leq C < \infty.$$

Hence, $\{\mathbf{a}_k\}$ is bounded.

(b) Let $\bar{\mathbf{a}} \in S$ be a limit point of $\{\mathbf{a}_k\}$. Take $\delta > 0$. Let $\{\mathbf{a}_{l_k}\}$ be a subsequence of $\{\mathbf{a}_k\}$ which converges to $\bar{\mathbf{a}}$. Since $\{\varepsilon_k\}$ is a summable sequence, there exists k_0 such that $\sum_{j=k_0}^{\infty} \varepsilon_j \leq \delta/2$, and there exists k_1 such that $l_{k_1} \geq k_0$ and $\|\mathbf{a}_{l_k} - \bar{\mathbf{a}}\|^2 < \delta/2$ for any $k \geq k_1$. Then, for any $k > l_{k_1}$,

$$\|\mathbf{a}_k - \bar{\mathbf{a}}\|^2 \leq \|\mathbf{a}_{l_{k_1}} - \bar{\mathbf{a}}\|^2 + \sum_{j=l_{k_1}}^{k-1} \varepsilon_j \leq \|\mathbf{a}_{l_{k_1}} - \bar{\mathbf{a}}\|^2 + \sum_{j=l_{k_1}}^{\infty} \varepsilon_j < \frac{\delta}{2} + \frac{\delta}{2} = \delta.$$

We conclude that $\mathbf{a}_k \rightarrow \bar{\mathbf{a}}$. ■

By the above lemma, we conclude that $\{\mathbf{x}_k\}$ is convergent to a vector \mathbf{x}^∞ . This vector must be stationary, by Theorem 12.3, which means, by convexity, that it is also globally optimal. We are done. ■

Suppose now that $X = \mathbb{R}^n$. Then the gradient projection algorithm reduces to the steepest descent method in unconstrained optimization, and the Armijo step length rule along the projection arc reduces to the classic Armijo rule. The above result then states that the steepest descent algorithm converges to an optimal solution whenever f is convex and there exist optimal solutions (see Theorem 11.7).

Finally, we consider the problem of performing the Euclidean projection. This is a strictly convex quadratic programming problem of the form (4.12). We will show that we can utilize the phase I procedure of the simplex method (see Section 9.1.2) in order to solve this problem. We take a slightly more general viewpoint here, and present the algorithm for a general strictly convex quadratic program.

Consider the problem to

$$\begin{aligned} & \text{minimize } f(\mathbf{x}) := \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} - \mathbf{q}^T \mathbf{x}, \\ & \text{subject to } \mathbf{A} \mathbf{x} \leq \mathbf{b}, \\ & \mathbf{x} \geq \mathbf{0}^n, \end{aligned} \tag{12.14}$$

where $\mathbf{Q} \in \mathbb{R}^{n \times n}$ is symmetric and positive definite, $\mathbf{q} \in \mathbb{R}^n$, $\mathbf{A} \in \mathbb{R}^{m \times n}$, and $\mathbf{b} \in \mathbb{R}^m$. Assuming that there exists a feasible solution to this problem, its unique optimal solution \mathbf{x}^* is characterized by its KKT conditions as follows:

$$\mathbf{Q} \mathbf{x} + \mathbf{A}^T \mathbf{y} - \mathbf{v} = \mathbf{q}, \tag{12.15a}$$

$$\mathbf{A} \mathbf{x} \leq \mathbf{b}, \tag{12.15b}$$

$$\mathbf{y}^T (\mathbf{A} \mathbf{x} - \mathbf{b}) = 0, \tag{12.15c}$$

$$\mathbf{v}^T \mathbf{x} = 0, \tag{12.15d}$$

$$\mathbf{x}, \mathbf{y}, \mathbf{v} \geq \mathbf{0}, \tag{12.15e}$$

where \mathbf{y} and \mathbf{v} are the vectors of Lagrange multipliers for the constraints of (12.14). We introduce a slack variable vector \mathbf{s} in (12.15b), and can therefore write the above system equivalently as

$$\mathbf{Q} \mathbf{x} + \mathbf{A}^T \mathbf{y} - \mathbf{v} = \mathbf{q}, \tag{12.16a}$$

$$\mathbf{A} \mathbf{x} + \mathbf{I}^m \mathbf{s} = \mathbf{b}, \tag{12.16b}$$

$$\mathbf{y}^T \mathbf{s} = 0, \tag{12.16c}$$

$$\mathbf{v}^T \mathbf{x} = 0, \tag{12.16d}$$

$$\mathbf{x}, \mathbf{s}, \mathbf{y}, \mathbf{v} \geq \mathbf{0}. \tag{12.16e}$$

Disregarding the complementarity conditions (12.16c), (12.16d), this is a set of linear equations over nonnegative variables, a solution to which can be found by using the phase I procedure of the simplex method (see Section 9.1). We propose to take the conditions (12.16c), (12.16d) into account implicitly, in the following way.

Introducing artificial variables in the system (12.16a), (12.16b), multiplying before-hand any equation with a negative right-hand side q_j or b_i by -1 , we let the artificial variables define the starting BFS in the phase I problem of minimizing their sum. When deciding on the incoming variable, we then introduce the following rules which make sure that the conditions (12.16c), (12.16d) are always satisfied:

- (a) If a variable x_j (respectively, v_j), $j = 1, \dots, n$, is already in the basis, then the variable v_j (respectively, x_j) is *not* admissible to enter the basis.

- (b) If a variable s_i (respectively, y_i), $i = 1, \dots, m$, is already in the basis, then the variable y_i (respectively, s_i) is *not* admissible to enter the basis.

It is not straightforward to argue why it is possible to reach the optimal solution when we restrict the incoming rule in this way. The interested reader is referred to the classic linear programming text by Dantzig [Dan63, Section 24.4], in which the above method is proven to yield convergence in a finite number of iterations provided that \mathbf{Q} is positive semidefinite.

12.5 Application: traffic equilibrium

12.5.1 Model analysis

The *traffic equilibrium problem* is a mathematical model that describes the steady-state of traffic in a transportation network, wherein each traveler minimizes his/her own travel costs for reaching the desired destination, irrespective of the modes of travel used.³ Traffic equilibrium models are used frequently as simulation models, aiding in the design or improvement of transportation systems. We develop a classic model of traffic equilibrium, and show that it is naturally given a variational inequality formulation; under certain conditions, it is also possible to state and solve it as a (strictly) convex optimization problem. We illustrate the performance of the Frank–Wolfe and simplicial decomposition methods on a sample problem.

Given is a graph (or, network) $\mathcal{G} = (\mathcal{N}, \mathcal{L})$ of nodes and directed links, a set $(p, q) \in \mathcal{C}$ of pairs of origin–destination (OD) nodes with fixed demands d_{pq} of units of desired traffic. Wardrop’s [War52] *user equilibrium* principle states that for every OD pair $(p, q) \in \mathcal{C}$, the travel costs of the routes utilized are equal and minimal for each individual user.⁴ We denote by h_r the volume of traffic on route $r \in \mathcal{R}_{pq}$, and by $c_r(\mathbf{h})$, $r \in \mathcal{R}$, the travel cost on the route as experienced by an individual user given the volume \mathbf{h} . Wardrop’s condition can be written as follows:

³While travel costs mostly are composed by travel times, modern traffic equilibrium models also take into account the possible uses of congestion tolls and the differences in the travelers’ values of time.

⁴The rationale behind this principle is, roughly, that if the same network user every morning travels between the same OD pair and he/she is not travelling along the best route, then through a trial-and-error procedure he/she will eventually reach the best one and stick to it; if every traveler behaves in the same fashion, the steady-state that is reached eventually must be a user equilibrium.

$$\mathbf{0}^{|\mathcal{R}|} \leq \mathbf{h} \perp (\mathbf{c}(\mathbf{h}) - \mathbf{\Gamma}\boldsymbol{\pi}) \geq \mathbf{0}^{|\mathcal{R}|}, \quad (12.17a)$$

$$\mathbf{\Gamma}^T \mathbf{h} = \mathbf{d}, \quad (12.17b)$$

where the value of π_{pq} is interpreted as the minimal (that is, equilibrium) route cost in OD pair (p, q) , and where we introduced the matrix $\mathbf{\Gamma} \in \mathbb{R}^{|\mathcal{R}| \times |\mathcal{C}|}$ to be the route-OD pair incidence matrix (i.e., the element γ_{rk} is 1 if route r joins OD pair $k := (p, q) \in \mathcal{C}$, and 0 otherwise). The first condition essentially states that more costly routes are not used, and the second describes the demand condition.

Let $\boldsymbol{\Lambda} \in \{0, 1\}^{|\mathcal{L}| \times |\mathcal{R}|}$ be the link-route incidence matrix, whose element λ_{lr} equals one if route $r \in \mathcal{R}$ utilizes link $l \in \mathcal{L}$, and zero otherwise. Route r has an *additive* route cost $c_r(\mathbf{h})$ if it is the sum of the costs of using all the links defining it. In other words, $c_r(\mathbf{h}) = \sum_{l \in \mathcal{L}} \lambda_{lr} t_l(\mathbf{v})$, where $\mathbf{v} \in \mathbb{R}^{|\mathcal{L}|}$ is the total volume of traffic on the links, and $t_l : \mathbb{R}^{|\mathcal{L}|} \rightarrow \mathbb{R}$ is a function measuring the travel cost on link $l \in \mathcal{L}$ given the link volume \mathbf{v} . In matrix-vector notation, then, $\mathbf{c}(\mathbf{h}) = \boldsymbol{\Lambda}^T \mathbf{t}(\mathbf{v})$ holds. Also, implicit in this cost relationship is the assumption that the pair (\mathbf{h}, \mathbf{v}) is consistent, in the sense that \mathbf{v} equals the sum of the route volumes: $\mathbf{v} = \boldsymbol{\Lambda} \mathbf{h}$.

Consider the following variational inequality: find $\mathbf{v}^* \in \hat{F}$ such that

$$\mathbf{t}(\mathbf{v}^*)^T (\mathbf{v} - \mathbf{v}^*) \geq 0, \quad \mathbf{v} \in \hat{F}, \quad (12.18)$$

where $\hat{F} := \{ \mathbf{v} \in \mathbb{R}^{|\mathcal{L}|} \mid \exists \mathbf{h} \in \mathbb{R}_+^{|\mathcal{R}|} \text{ with } \mathbf{\Gamma}^T \mathbf{h} = \mathbf{d} \text{ and } \mathbf{v} = \boldsymbol{\Lambda} \mathbf{h} \}$ is the set of demand-feasible link volumes.

In the case where \mathbf{t} is *integrable*,⁵ the model (12.18) defines the first-order optimality conditions for an optimization problem; assuming, further, that \mathbf{t} is *separable*, that is, that t_l is a function only of v_l , $l \in \mathcal{L}$, the optimization problem has the form

$$\begin{aligned} & \underset{(\mathbf{h}, \mathbf{v})}{\text{minimize}} \quad f(\mathbf{v}) := \sum_{l \in \mathcal{L}} \int_0^{v_l} t_l(s) ds, \\ & \text{subject to} \quad \mathbf{\Gamma}^T \mathbf{h} = \mathbf{d}, \\ & \quad \mathbf{v} = \boldsymbol{\Lambda} \mathbf{h}, \\ & \quad \mathbf{h} \geq \mathbf{0}^{|\mathcal{R}|}. \end{aligned} \quad (12.19)$$

This is the classic *traffic assignment problem*.

Since the feasible set of the problem (12.19) is a bounded polyhedron there exists a nonempty and bounded set of optimal link and route

⁵If \mathbf{t} is continuously differentiable, then integrability is equivalent to the symmetry of its Jacobian matrix $\nabla \mathbf{t}(\mathbf{v})$ everywhere. Integrability is a more general property than this symmetry property, since \mathbf{t} need not be always be differentiable.

volumes, and the optimality conditions given by the respective variational inequality (or by the KKT conditions) are necessary for the local optimality of a pair of link and route volumes [cf. Proposition 4.23(b)].

Moreover, the optimality conditions are exactly the Wardrop conditions of user equilibrium. To see this, suppose that \mathbf{v}^* is a local minimum in (12.19) and consider the following problem, a solution to which necessarily then is \mathbf{v}^* [cf. (4.14)]:

$$\begin{aligned} & \underset{(\mathbf{h}, \mathbf{v})}{\text{minimize}} \quad \mathbf{t}(\mathbf{v}^*)^T \mathbf{v}, \\ & \text{subject to} \quad \mathbf{\Gamma}^T \mathbf{h} = \mathbf{d}, \\ & \quad \mathbf{v} - \mathbf{\Lambda} \mathbf{h} = \mathbf{0}^{|\mathcal{L}|}, \\ & \quad \mathbf{h} \geq \mathbf{0}^{|\mathcal{R}|}. \end{aligned} \tag{12.20}$$

Note that we are simply rephrasing (12.18). Introducing the LP dual variable vectors $\boldsymbol{\pi} \in \mathbb{R}^{|\mathcal{C}|}$ and $\boldsymbol{\alpha} \in \mathbb{R}^{|\mathcal{L}|}$, the LP dual to (12.20) is to

$$\begin{aligned} & \underset{(\boldsymbol{\pi}, \boldsymbol{\alpha})}{\text{maximize}} \quad \mathbf{d}^T \boldsymbol{\pi}, \\ & \text{subject to} \quad \mathbf{\Gamma} \boldsymbol{\pi} - \mathbf{\Lambda}^T \boldsymbol{\alpha} \leq \mathbf{0}^{|\mathcal{R}|}, \\ & \quad \boldsymbol{\alpha} = \mathbf{t}(\mathbf{v}^*). \end{aligned} \tag{12.21}$$

Eliminating $\boldsymbol{\alpha}$ through $\boldsymbol{\alpha} = \mathbf{t}(\mathbf{v}^*)$ the primal–dual optimality conditions are, precisely, the Wardrop conditions (12.17), together with the consistency condition $\mathbf{v} = \mathbf{\Lambda} \mathbf{h}$.

Suppose, in addition, that each link cost function t_l is increasing; this is a natural assumption considering that congestion on a link, that is, the travel time, increases with its volume. According to Theorem 3.40(b) this means that the function f is convex, and therefore the problem (12.19) is a convex one. Therefore also the optimality conditions stated in the variational inequality (12.18) or the (equivalent) Wardrop conditions (12.17) are both necessary and sufficient for a volume \mathbf{v} to be an equilibrium one.

If further t_l is *strictly* increasing for every $l \in \mathcal{L}$ then the solution \mathbf{v}^* is unique (cf. Proposition 4.11).

12.5.2 Algorithms and a numerical example

When solving this problem by using the Frank–Wolfe or the simplicial decomposition method the search directions in Step 1 correspond to a very special problem: if the current traffic volume is \mathbf{v}_k then the solution \mathbf{y}_k to the LP problem is found by assigning, for each OD pair $(p, q) \in \mathcal{C}$,

the demand d_{pq} onto a *shortest route* between the origin and destination node given the fixed link cost vector $\mathbf{t}(\mathbf{v}_k)$, and then aggregating these route volumes through the relation $\mathbf{v} = \mathbf{\Lambda}\mathbf{h}$; for each OD pair, the shortest route is found by using, for example, Dijkstra's algorithm. Doing this, we in fact need not store any route information at all, which saves computer storage.

A Matlab implementation of the two algorithms was in [Jos03] devised and tested on a classic traffic assignment problem, modelling the small city of Sioux Falls in South Dakota, USA, whose traffic network representation has 24 nodes, 76 links, and 528 OD pairs.

In the simplicial decomposition algorithm, we tested three algorithms for the restricted master problems (RSMPs)—a Newton method and two gradient projection methods. In Figure 12.3 we illustrate the solution times necessary for reaching a given accuracy; accuracy is here measured in terms of the relative error stemming from the lower and upper bounds on the optimal value.

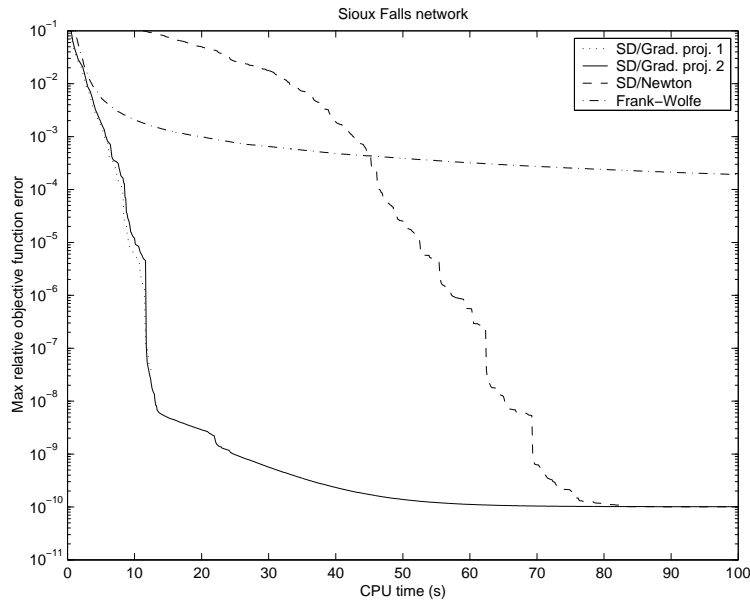


Figure 12.3: The performance of SD vs. FW on the Sioux Falls network.

It is clear that the Frank–Wolfe algorithm is extremely slow in comparison with the simplicial decomposition algorithm, regardless of the method used for solving the restricted master problems. An analysis of the algorithm's behaviour reveals that very small step lengths are taken

already quite early on in the solution process. The basis for this behaviour is the congestion effects that imply that several OD pairs need more than one route to have a positive volume at the solution; this means that the optimal link volume is not an extreme point, and the solutions to (12.2) will zig-zag between assigning the total volume onto these routes.

12.6 Notes and further reading

Algorithms for linearly constrained optimization problems are disappearing from modern text books on optimization. It is perhaps a sign of maturity, as we are now better at solving optimization problem with general constraints, and therefore do no longer have to especially consider the class of linearly constrained optimization problems. Nevertheless we feel that it provides a link between linear programming and nonlinear optimization problems with general constraints, being a subclass of nonlinear optimization problems for which primal feasibility can be retained throughout the procedure.

The Frank–Wolfe method was developed for QP problems in [FrW56], and later for more general problems, including non-polyhedral sets, in [Gil66] and [PsD78, Section III.3], among others. The latter source includes several convergence results for the method under different step length rules, assuming that ∇f is Lipschitz continuous, for example a Newton-type step length rule. The convergence Theorem 12.1 for the Frank–Wolfe algorithm was taken from [Pat98, Theorem 5.8]. The convergence result for convex problems given in Theorem 12.2 is due to Dunn and Harshbarger [DuH78]. The version of the Frank–Wolfe algorithm produced by the selection $\alpha_k := 1/k$ is known as the *method of successive averages* (MSA).

The simplicial decomposition algorithm was developed in [vHo77]. Restricted simplicial decomposition methods have been developed in [HLV87, Pat98].

The gradient projection method presented here was first given in [Gol64, LeP66]; see also the textbook [Ber99]. Theorem 12.4 is due to [Ius03], while Lemma 12.5 is due to [BGIS95].

The traffic equilibrium models of Section 12.5 are described and analyzed more fully in [She85, Pat94].

Apart from the algorithms developed here, there are other classical algorithms for linearly constrained problems, including the reduced gradient method, Rosen’s gradient projection method, active set methods, and other sub-manifold methods. They are not treated here, as some of

them have fallen out of popularity. Reduced gradient methods still constitute the main building block of some commercial software, however.

12.7 Exercises

Exercise 12.1 (extensions of the Frank–Wolfe algorithm to unbounded sets) Develop an extension to the Frank–Wolfe algorithm applicable to cases where X is unbounded. Which steps need to be changed? What can go wrong?

Exercise 12.2 (convergence of the sequence $\{z_k(\mathbf{y}_k)\}$) We are interested in the convergence of the sequence $\{z_k(\mathbf{y}_k)\} := \{\nabla f(\mathbf{x}_k)^\top \mathbf{p}_k\} := \{\nabla f(\mathbf{x}_k)^\top (\mathbf{y}_k - \mathbf{x}_k)\}$ of optimal values of the linear objective function in Step 1 of the Frank–Wolfe algorithm. In addition to the properties of the problem (12.1) and the execution of the Frank–Wolfe algorithm assumed in Theorem 12.1, suppose that

$$\nabla f \text{ is Lipschitz continuous on } \mathbb{R}^n.$$

Establish that $z_k(\mathbf{y}_k) \rightarrow 0$ holds.

[Hint: Utilize the following technical lemma, which is of independent interest and therefore given a proof:

Lemma 12.6 (descent lemma) Suppose that ∇f is Lipschitz continuous on \mathbb{R}^n , with modulus L . Let \mathbf{x}, \mathbf{p} both lie in \mathbb{R}^n . Then,

$$f(\mathbf{x} + \mathbf{p}) - f(\mathbf{x}) \leq \nabla f(\mathbf{x})^\top \mathbf{p} + \frac{L}{2} \|\mathbf{p}\|^2$$

holds.

Proof. Let $\ell \in \mathbb{R}$ and $g(\ell) := f(\mathbf{x} + \ell \mathbf{p})$. The chain rule yields $\frac{dg}{d\ell}(\ell) = \mathbf{p}^\top \nabla f(\mathbf{x} + \ell \mathbf{p})$. Then,

$$\begin{aligned} f(\mathbf{x} + \mathbf{p}) - f(\mathbf{x}) &= g(1) - g(0) = \int_0^1 \frac{dg}{d\ell}(\ell) d\ell = \int_0^1 \mathbf{p}^\top \nabla f(\mathbf{x} + \ell \mathbf{p}) d\ell \\ &\leq \int_0^1 \mathbf{p}^\top \nabla f(\mathbf{x}) d\ell + \left| \int_0^1 \mathbf{p}^\top [\nabla f(\mathbf{x} + \ell \mathbf{p}) - \nabla f(\mathbf{x})] d\ell \right| \\ &\leq \int_0^1 \mathbf{p}^\top \nabla f(\mathbf{x}) d\ell + \int_0^1 \|\mathbf{p}\| \|\nabla f(\mathbf{x} + \ell \mathbf{p}) - \nabla f(\mathbf{x})\| d\ell \\ &\leq \mathbf{p}^\top \nabla f(\mathbf{x}) + \|\mathbf{p}\| \int_0^1 L\ell d\ell \\ &\leq \mathbf{p}^\top \nabla f(\mathbf{x}) + \frac{L}{2} \|\mathbf{p}\|^2. \end{aligned}$$

We are done. ■

Apply this result to the inequality resulting from applying the Armijo rule at a given iteration k , with \mathbf{x} replaced by \mathbf{x}_k and \mathbf{p} replaced by $\alpha_k \mathbf{p}_k$.

Summing all these inequalities and utilizing that $\{\|p_k\|\}$ is a bounded sequence thanks to the boundedness of X , conclude that the sum $\sum_{k=0}^{\infty} [z_k(\mathbf{y}_k)]^2$ must be convergent and therefore $z_k(\mathbf{y}_k) \rightarrow 0$ must hold.]

Exercise 12.3 (numerical example of the Frank–Wolfe algorithm) Consider the problem to

$$\begin{aligned} & \text{minimize } f(\mathbf{x}) := (x_1 + 2x_2 - 6)^2 + (2x_1 - x_2 - 2)^2, \\ & \text{subject to } 2x_1 + 3x_2 \leq 9, \\ & \quad x_1 \leq 3, \\ & \quad x_1, x_2 \geq 0. \end{aligned}$$

- (a) Show that the problem is convex.
- (b) Apply one step of the Frank–Wolfe algorithm, starting at the origin. Provide an interval where f^* lies.

Exercise 12.4 (numerical example of the Frank–Wolfe algorithm) Consider the problem to

$$\begin{aligned} & \text{maximize } f(\mathbf{x}) := -x_1^2 - 4x_2^2 + 16x_1 + 24x_2, \\ & \text{subject to } x_1 + x_2 \leq 6, \\ & \quad x_1 - x_2 \leq 3, \\ & \quad x_1, x_2 \geq 0. \end{aligned}$$

- (a) Show that the problem is convex.
- (b) Solve the problem by using the Frank–Wolfe algorithm, starting at the origin.

Exercise 12.5 (numerical example of the Frank–Wolfe algorithm) Consider the problem to

$$\begin{aligned} & \text{minimize } f(\mathbf{x}) := \frac{1}{2} \left(x_1 - \frac{1}{2} \right)^2 + \frac{1}{2} x_2^2, \\ & \text{subject to } x_1 \leq 1, \\ & \quad x_2 \leq 1, \\ & \quad x_1, x_2 \geq 0. \end{aligned}$$

Apply two iterations of the Frank–Wolfe algorithm, starting at $\mathbf{x}_0 := (1, 1)^T$. Give upper and lower bounds on the optimal value.

Exercise 12.6 (convergence of a gradient projection algorithm) Establish Theorem 12.3.

Exercise 12.7 (numerical example of the simplicial decomposition algorithm) Solve the problem in Exercise 12.3 by using the simplicial decomposition algorithm.

Optimization over convex sets

Exercise 12.8 (numerical example of the simplicial decomposition algorithm)
Solve the problem in Exercise 12.4 by using the simplicial decomposition algorithm.

Exercise 12.9 (numerical example of the simplicial decomposition algorithm)
On the problem in Exercise 12.5 apply two iterations of the simplicial decomposition algorithm. Is \mathbf{x}_2 optimal? Why/why not?

Constrained optimization

XIII

In this chapter, we will discuss the conversion of nonlinear programming problems with inequality and equality constraints into (in some sense) equivalent unconstrained problems or problems with simple constraints. In practice, a sequence of such equivalent (or, approximating) problems is solved because of computational considerations.

13.1 Penalty methods

Let us consider a general optimization problem:

$$\begin{aligned} &\text{minimize } f(\mathbf{x}), \\ &\text{subject to } \mathbf{x} \in S, \end{aligned} \tag{13.1}$$

where $S \subseteq \mathbb{R}^n$ is a nonempty, closed set and $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a given differentiable function. The basic idea behind all penalty algorithms is to replace the problem (13.1) with the equivalent unconstrained one:

$$\text{minimize } f(\mathbf{x}) + \chi_S(\mathbf{x}), \tag{13.2}$$

where

$$\chi_S(\mathbf{x}) = \begin{cases} 0, & \text{if } \mathbf{x} \in S, \\ +\infty, & \text{otherwise.} \end{cases}$$

The role of χ_S , which in the optimization community is known as the *indicator function* of the set S , is to make sure that feasibility is top priority, and only when achieving feasibility do we concentrate on optimizing the function f . Of course, the so defined χ_S is rather bizarre from the computational point of view: it is non-differentiable, discontinuous, and even not finite (though it is convex provided S is). Thus, from the

practical point of view we would like to replace the additional term χ_S with a numerically better behaving function.

There are two alternative approaches achieving this. The first is called the *penalty*, or the *exterior penalty* method, in which we add a penalty to the objective function for points not lying in the feasible set and thus violating some of the constraints. This method typically generates a sequence of infeasible points, approaching optimal solutions to the original problem from the outside (exterior) of the feasible set, whence the name of the method. The function χ_S is approximated “from below” in these methods.

Alternatively, in the *barrier*, or *interior point* methods, we add a continuous barrier term that equals $+\infty$ everywhere except in the interior of the feasible set and thus ensure that globally optimal solutions to the approximating unconstrained problems do not escape the feasible set of the original constrained problem. The method thus generates a sequence of interior points, whose limit is an optimal solution to the original constrained problem. The function χ_S is approximated “from above” in these methods.

Clearly we would like to transfer “nice” properties of original constrained problems, such as convexity, smoothness, to penalized problems as well. We easily achieve this by carefully choosing penalty functions; use Exercises 13.1 and 13.2 to verify that convexity may be easily transferred to penalized problems.

13.1.1 Exterior penalty methods

We assume that the feasible set S of the optimization problem (13.1) is given by a system of inequality and equality constraints:

$$S := \{ \mathbf{x} \in \mathbb{R}^n \mid g_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m; \\ h_j(\mathbf{x}) = 0, \quad j = 1, \dots, \ell \}, \quad (13.3)$$

where $g_i \in C(\mathbb{R}^n)$, $i = 1, \dots, m$, $h_j \in C(\mathbb{R}^n)$, $j = 1, \dots, \ell$. In this case, we can choose a continuous function $\psi : \mathbb{R} \rightarrow \mathbb{R}_+$ such that $\psi(s) = 0$ if and only if $s = 0$ (typical examples of $\psi(\cdot)$ will be $\psi_1(s) = |s|$, or $\psi_2(s) = s^2$), and try the approximation:

$$\chi_S(\mathbf{x}) \approx \nu \tilde{\chi}_S(\mathbf{x}) := \nu \left(\sum_{i=1}^m \psi(\max\{0, g_i(\mathbf{x})\}) + \sum_{j=1}^{\ell} \psi(h_j(\mathbf{x})) \right), \quad (13.4)$$

where the real number $\nu > 0$ is called the *penalty parameter*. The different treatment of inequality and equality constraints in the equation (13.4) stems from the fact that equality constraints are violated at

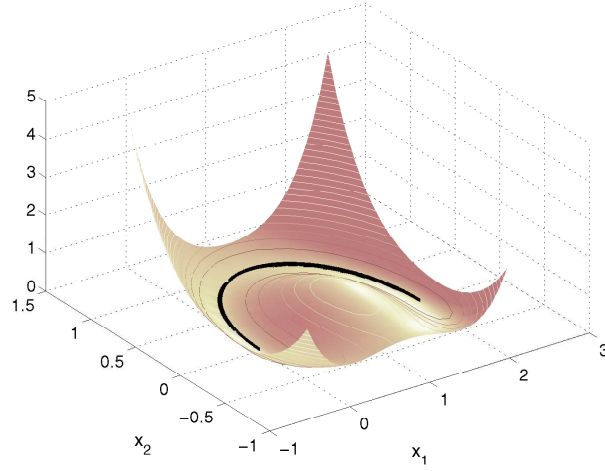


Figure 13.1: The graph of $\check{\chi}_S$ and the feasible set S (black).

\mathbf{x} whenever $h_j(\mathbf{x}) \neq 0$ for some $j = 1, \dots, \ell$, while inequality constraints are violated only when $g_i(\mathbf{x}) > 0$ for some $i = 1, \dots, m$; the latter fact can be equivalently expressed as $\max\{0, g_i(\mathbf{x})\} \neq 0$.

Example 13.1 We repeat the settings of Example 5.7. Let $S := \{\mathbf{x} \in \mathbb{R}^2 \mid -x_2 \leq 0; (x_1 - 1)^2 + x_2^2 = 1\}$. Let $\psi(s) = s^2$. Then, in this example,

$$\check{\chi}_S(\mathbf{x}) := [\max\{0, -x_2\}]^2 + [(x_1 - 1)^2 + x_2^2 - 1]^2.$$

The graph of the function $\check{\chi}_S$, together with the feasible set S , is shown in Figure 13.1. ■

To exclude trivial cases, we assume that the original constrained problem (13.1), and thus its equivalent reformulation (13.2), has an optimal solution \mathbf{x}^* . Furthermore, we assume that for every $\nu > 0$ the approximating optimization problem to

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{x}) + \nu \check{\chi}_S(\mathbf{x}) \quad (13.5)$$

has at least one optimal solution \mathbf{x}_ν^* .

Clearly, $\check{\chi}_S$ is non-negative and, furthermore, $\check{\chi}_S(\mathbf{x}) = 0$ if and only if $\mathbf{x} \in S$. Therefore, from the Relaxation Theorem (Theorem 6.1) we know that the inequality $f(\mathbf{x}_{\nu_1}^*) + \nu_1 \check{\chi}_S(\mathbf{x}_{\nu_1}^*) \leq f(\mathbf{x}_{\nu_2}^*) + \nu_2 \check{\chi}_S(\mathbf{x}_{\nu_2}^*) \leq$

Constrained optimization

$f(\mathbf{x}^*) + \chi_S(\mathbf{x}^*) = f(\mathbf{x}^*)$ holds for every positive $\nu_1 \leq \nu_2$. In fact, we can establish an even stronger inequality, which will be used later; see the following lemma.

Lemma 13.2 (penalization constitutes a relaxation) *For every positive $\nu_1 \leq \nu_2$ it holds that $f(\mathbf{x}_{\nu_1}^*) \leq f(\mathbf{x}_{\nu_2}^*)$.*

Proof. The claim is trivial for $\nu_1 = \nu_2$, thus we assume that $\nu_1 < \nu_2$. Since $\mathbf{x}_{\nu_1}^*$ minimizes $f(\mathbf{x}) + \nu_1 \check{\chi}_S(\mathbf{x})$, and $\mathbf{x}_{\nu_2}^*$ is feasible in this (unconstrained) optimization problem, it holds that

$$f(\mathbf{x}_{\nu_1}^*) + \nu_1 \check{\chi}_S(\mathbf{x}_{\nu_1}^*) \leq f(\mathbf{x}_{\nu_2}^*) + \nu_1 \check{\chi}_S(\mathbf{x}_{\nu_2}^*). \quad (13.6)$$

Similarly,

$$f(\mathbf{x}_{\nu_2}^*) + \nu_2 \check{\chi}_S(\mathbf{x}_{\nu_2}^*) \leq f(\mathbf{x}_{\nu_1}^*) + \nu_2 \check{\chi}_S(\mathbf{x}_{\nu_1}^*).$$

Adding the two inequalities, we conclude that

$$(\nu_2 - \nu_1)[\check{\chi}_S(\mathbf{x}_{\nu_2}^*) - \check{\chi}_S(\mathbf{x}_{\nu_1}^*)] \leq 0,$$

which, substituted into (13.6), implies the claim, because $\nu_2 - \nu_1 > 0$. ■

Now we are ready to show that every limit point of the sequence $\{\mathbf{x}_\nu^*\}$, as ν converges to infinity, is optimal in the problem (13.1). Thus, the family of problems (13.5) is indeed an approximation of the original problem (13.1), and setting ν to a “large enough” value we can solve the problem (13.5) in place of (13.1).

Theorem 13.3 (global convergence of a penalty method) *Assume that the original constrained problem (13.1) possesses optimal solutions. Then, every limit point of the sequence $\{\mathbf{x}_\nu^*\}$, $\nu \rightarrow +\infty$, of globally optimal solutions to (13.5) is globally optimal in the problem (13.1).*

In other words,

$$\left. \begin{array}{l} \mathbf{x}_\nu^* \text{ globally optimal in (13.5)} \\ \mathbf{x}_\nu^* \rightarrow \mathbf{x}^* \text{ as } \nu \rightarrow +\infty \end{array} \right\} \implies \mathbf{x}^* \text{ globally optimal in (13.1)}.$$

Proof. Let \mathbf{x}^* denote an arbitrary globally optimal solution to (13.1).

From the inequality (cf. the Relaxation Theorem 6.1)

$$f(\mathbf{x}_\nu^*) + \nu \check{\chi}_S(\mathbf{x}_\nu^*) \leq f(\mathbf{x}^*), \quad (13.7)$$

and Lemma 13.2, we obtain uniform bounds on the penalty term $\nu \check{\chi}_S(\mathbf{x}_\nu^*)$ for all $\nu \geq 1$:

$$0 \leq \nu \check{\chi}_S(\mathbf{x}_\nu^*) \leq f(\mathbf{x}^*) - f(\mathbf{x}_1^*).$$

Thus, $\check{\chi}_S(\mathbf{x}_\nu^*)$ converges to zero as ν converges to $+\infty$, and, owing to the continuity of $\check{\chi}_S$, every limit point of the sequence $\{\mathbf{x}_\nu^*\}$ must be feasible in (13.1).

Now, let $\hat{\mathbf{x}}$ denote an arbitrary limit point of $\{\mathbf{x}_\nu^*\}$, that is,

$$\lim_{k \rightarrow \infty} \mathbf{x}_{\nu_k}^* = \hat{\mathbf{x}},$$

for some sequence $\{\nu_k\}$ converging to infinity. Then, we have the following chain of inequalities:

$$f(\hat{\mathbf{x}}) = \lim_{k \rightarrow +\infty} f(\mathbf{x}_{\nu_k}^*) \leq \lim_{k \rightarrow +\infty} \{f(\mathbf{x}_{\nu_k}^*) + \nu_k \check{\chi}_S(\mathbf{x}_{\nu_k}^*)\} \leq f(\mathbf{x}^*),$$

where the last inequality follows from (13.7). However, owing to the feasibility of $\hat{\mathbf{x}}$ in (13.1) the reverse inequality $f(\mathbf{x}^*) \leq f(\hat{\mathbf{x}})$ must also hold. The two inequalities combined imply the required claim. ■

We emphasize that Theorem 13.3 establishes the convergence of *globally* optimal solutions only; the result may therefore be of limited practical value for nonconvex nonlinear programs. However, assuming more regularity of the stationary points, such as LICQ (see Definition 5.41), and using specific continuously differentiable penalty functions, such as $\psi(s) := s^2$, we can show that every limit point of sequences of stationary points of (13.5) also is stationary (i.e., a KKT point) in (13.1). Furthermore, we easily obtain estimates of the corresponding Lagrange multipliers $(\hat{\mu}, \hat{\lambda})$.

Theorem 13.4 (convergence of a penalty method) *Let the objective function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and the functions $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$, $i = 1, \dots, m$, and $h_j : \mathbb{R}^n \rightarrow \mathbb{R}$, $j = 1, \dots, \ell$, defining the inequality and equality constraints of (13.1) be in $C^1(\mathbb{R}^n)$. Further assume that the penalty function $\psi : \mathbb{R} \rightarrow \mathbb{R}_+$ is in C^1 and that $\psi'(s) \geq 0$ for all $s \geq 0$.*

Consider a sequence $\{\mathbf{x}_k\}$ of points that are stationary for the sequence of problems (13.5), for some positive sequence of penalty parameters $\{\nu_k\}$ converging to $+\infty$. Assume that $\lim_{k \rightarrow +\infty} \mathbf{x}_k = \hat{\mathbf{x}}$, and that LICQ holds at $\hat{\mathbf{x}}$. Then, if $\hat{\mathbf{x}}$ is feasible in (13.1) it must also verify the KKT conditions.

In other words,

$$\left. \begin{array}{l} \mathbf{x}_k \text{ stationary in (13.5)} \\ \mathbf{x}_k \rightarrow \hat{\mathbf{x}} \text{ as } k \rightarrow +\infty \\ \text{LICQ holds at } \hat{\mathbf{x}} \\ \hat{\mathbf{x}} \text{ feasible in (13.1)} \end{array} \right\} \implies \hat{\mathbf{x}} \text{ stationary in (13.1)}.$$

Proof. [Sketch] Owing to the optimality conditions (4.14) for unconstrained optimization we know that every point \mathbf{x}_k , $k = 1, 2, \dots$, necessarily satisfies the equation

$$\nabla[f(\mathbf{x}_k) + \nu_k \check{\chi}_S(\mathbf{x}_k)] = \nabla f(\mathbf{x}_k) \quad (13.8a)$$

$$+ \sum_{i=1}^m \nu_k \psi'[\max\{0, g_i(\mathbf{x}_k)\}] \nabla g_i(\mathbf{x}_k) \quad (13.8b)$$

$$+ \sum_{j=1}^{\ell} \nu_k \psi'[h_j(\mathbf{x}_k)] \nabla h_j(\mathbf{x}_k) = \mathbf{0}^n. \quad (13.8c)$$

Let, as before, $\mathcal{I}(\hat{\mathbf{x}})$ denote the index set of active inequality constraints at $\hat{\mathbf{x}}$. If $i \notin \mathcal{I}(\hat{\mathbf{x}})$ then $g_i(\mathbf{x}_k) < 0$ for all large k , and the terms corresponding to this index do not contribute to (13.8).

Since LICQ holds at $\hat{\mathbf{x}}$, we know that the vectors $\{\nabla g_i(\hat{\mathbf{x}}), \nabla h_j(\hat{\mathbf{x}}) \mid i \in \mathcal{I}(\hat{\mathbf{x}}), j = 1, \dots, \ell\}$ are linearly independent. Therefore, we can easily show that the sequence $\{\nu_k \psi'[\max\{0, g_i(\mathbf{x}_k)\}]\}$ must converge to some limit $\hat{\mu}_i$ as $k \rightarrow +\infty$ for all $i \in \mathcal{I}(\hat{\mathbf{x}})$. Similarly, $\lim_{k \rightarrow +\infty} \nu_k \psi'[h_j(\mathbf{x}_k)] = \hat{\lambda}_j$, $j = 1, \dots, \ell$. At last, since $\nu_k \psi'[\max\{0, g_i(\mathbf{x}_k)\}] \geq 0$ for all $k = 1, 2, \dots$, $i \in \mathcal{I}(\hat{\mathbf{x}})$ it follows that $\hat{\boldsymbol{\mu}} \geq \mathbf{0}^{|\mathcal{I}(\hat{\mathbf{x}})|}$.

Passing to the limit as $k \rightarrow +\infty$ in (13.8) we deduce that

$$\nabla f(\hat{\mathbf{x}}) + \sum_{i \in \mathcal{I}(\hat{\mathbf{x}})} \hat{\mu}_i \nabla g_i(\hat{\mathbf{x}}) + \sum_{j=1}^{\ell} \hat{\lambda}_j \nabla h_j(\hat{\mathbf{x}}) = \mathbf{0}^n,$$

i.e., $\hat{\mathbf{x}}$ is a KKT point for (13.1) with Lagrange multipliers $(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\lambda}})$. ■

Notice that if the original problem (13.1) is convex and verifies LICQ, and if every penalized problem is also convex (cf. Exercise 13.1), then Theorems 13.3 and 13.4 essentially work with the same sequences: under convexity and LICQ globally optimal solutions are KKT points and vice versa. Therefore, in this case we automatically get feasibility of limit points in Theorem 13.3, as well as expressions for estimating Lagrange multipliers in Theorem 13.4.

13.1.2 Interior penalty methods

While the idea behind exterior penalty functions is to nicely approximate χ_S on the whole of \mathbb{R}^n , interior penalty, or *barrier*, function methods construct approximations only inside the feasible set and set a barrier against leaving it. If a globally optimal solution to (13.1) happens to be

located on the boundary of the feasible region, then the method generates a sequence of interior points that converges to it.

In this section we assume that the feasible set S of the optimization problem (13.1) has the following form:

$$S := \{ \mathbf{x} \in \mathbb{R}^n \mid g_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m \}. \quad (13.9)$$

For the method to work, we need to assume that there exists a *strictly feasible* point $\hat{\mathbf{x}} \in \mathbb{R}^n$, that is, such that $g_i(\hat{\mathbf{x}}) < 0, i = 1, \dots, m$. Thus, in contrast with the exterior penalty algorithms, we cannot include equality constraints into the penalty term. While it is possible to extend the discussion to allow for equality constraints, we prefer to keep the notation simple and assume that equality constraints are not present.

To formulate a *barrier problem*, we consider the following approximation of χ_S :

$$\chi_S(\mathbf{x}) \approx \nu \hat{\chi}_S(\mathbf{x}) := \begin{cases} \nu \sum_{i=1}^m \phi[g_i(\mathbf{x})], & \text{if } g_i(\mathbf{x}) < 0, i = 1, \dots, m, \\ +\infty, & \text{otherwise,} \end{cases} \quad (13.10)$$

and the function $\phi : \mathbb{R}_- \rightarrow \mathbb{R}_+$ is a continuous nonnegative function such that $\phi(s_k) \rightarrow +\infty$ for all *negative* sequences $\{s_k\}$ converging to zero. Typical examples of ϕ are $\phi_1(s) := -s^{-1}$, and $\phi_2(s) := -\log[\min\{1, -s\}]$. Note that ϕ_2 is not differentiable at the point $s = -1$. However, dropping the nonnegativity requirement on ϕ , the famous differentiable *logarithmic barrier function* $\phi_2(s) := -\log(-s)$ gives rise to the same convergence theory as we are going to present.

Example 13.5 Consider the simple one-dimensional set $S := \{x \in \mathbb{R} \mid -x \leq 0\}$. Choosing $\phi = \phi_1 = -s^{-1}$, the graph of the barrier function $\nu \hat{\chi}_S$ is shown in Figure 13.2 for various values of ν . Note how $\nu \hat{\chi}_S$ converges towards χ_S as $\nu \downarrow 0$. ■

Having chosen the function ϕ and a penalty parameter $\nu > 0$, we are going to solve the following problem:

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{x}) + \nu \hat{\chi}_S(\mathbf{x}). \quad (13.11)$$

Similarly to the case of exterior penalty functions discussed in the previous section, we can prove the convergence to globally optimal solutions (however, in this case we need to assume the regularity assumption $S = \text{cl}\{\mathbf{x} \in \mathbb{R}^n \mid g_i(\mathbf{x}) < 0, i = 1, \dots, m\}$). Rather, we proceed directly to establish a convergence result for stationary points, similar to Theorem 13.4. Not only is this result more practical than the one concerning

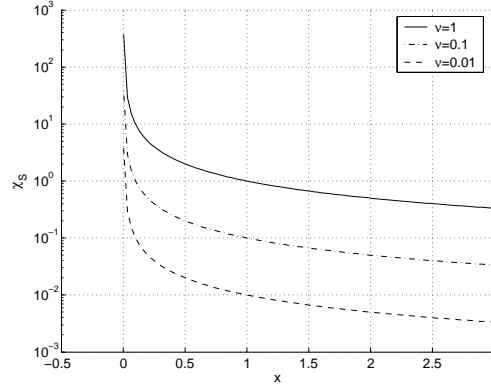


Figure 13.2: The graph of $\nu\chi_S$ for various choices of ν . Note the logarithmic scale.

globally optimal solutions, but also the interior point algorithms are most often applied to convex optimization problems, and thus stationarity implies global optimality (see Section 5.8). The reason is that interior point algorithms are especially efficient both practically and theoretically for convex optimization problems. In fact, one can show that the number of computational steps an interior point algorithm needs in order to achieve a prescribed accuracy $\varepsilon > 0$ is bounded by a polynomial function of the “size” of the problem (that is, the number of variables and constraints) and ε^{-1} . For non-convex problems, on the contrary, it is known that the number of steps necessary can grow exponentially. For other algorithms that can be applied to convex optimization problems, for example, exterior penalty methods, no well-developed complexity theory exists.

The proof of the general convergence theorem for barrier methods goes in parallel with the corresponding result for exterior penalty methods. An important difference, though, is that now the constrained problem (13.1) is the relaxation of (13.11) for every $\nu > 0$, and the convergence is studied as $\nu \downarrow 0$.

Theorem 13.6 (convergence of an interior point algorithm) *Let the objective function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and the functions g_i , $i = 1, \dots, m$, defining the inequality constraints of (13.9) be in $C^1(\mathbb{R}^n)$. Further assume that the barrier function $\phi : \mathbb{R} \rightarrow \mathbb{R}_+$ is in C^1 and that $\phi'(s) \geq 0$ for all $s < 0$.*

Consider a sequence $\{\mathbf{x}_k\}$ of points that are stationary for the sequence of problems (13.11) with $\nu = \nu_k$, for some positive sequence of penalty parameters $\{\nu_k\}$ converging to 0. Assume that $\lim_{k \rightarrow +\infty} \mathbf{x}_k = \hat{\mathbf{x}}$, and that LICQ holds at $\hat{\mathbf{x}}$. Then, $\hat{\mathbf{x}}$ is a KKT point of (13.1).

In other words,

$$\left. \begin{array}{l} \mathbf{x}_k \text{ stationary in (13.11)} \\ \mathbf{x}_k \rightarrow \hat{\mathbf{x}} \text{ as } k \rightarrow +\infty \\ \text{LICQ holds at } \hat{\mathbf{x}} \end{array} \right\} \implies \hat{\mathbf{x}} \text{ stationary in (13.1).}$$

Proof. [Sketch] Owing to the optimality conditions (4.14) for unconstrained optimization we know that every point \mathbf{x}_k , $k = 1, 2, \dots$, necessarily satisfies the equation

$$\begin{aligned} \nabla[f(\mathbf{x}_k) + \nu_k \chi_S(\mathbf{x}_k)] &= \\ \nabla f(\mathbf{x}_k) + \sum_{i=1}^m \nu_k \phi'[g_i(\mathbf{x}_k)] \nabla g_i(\mathbf{x}_k) &= \mathbf{0}^n. \end{aligned} \quad (13.12)$$

Because every point \mathbf{x}_k is *strictly* feasible in (13.1), the limit $\hat{\mathbf{x}}$ is clearly feasible in (13.1). Let $\mathcal{I}(\hat{\mathbf{x}})$ denote the index set of active inequality constraints at $\hat{\mathbf{x}}$.

If $i \notin \mathcal{I}(\hat{\mathbf{x}})$ then $g_i(\mathbf{x}_k) < 0$ for all large k , and $\nu_k \phi'[g_i(\mathbf{x}_k)] \rightarrow 0$ as $\nu_k \downarrow 0$.

Since LICQ holds at $\hat{\mathbf{x}}$, we know that the vectors $\{\nabla g_i(\hat{\mathbf{x}}) \mid i \in \mathcal{I}(\hat{\mathbf{x}})\}$ are linearly independent. Therefore, we can easily show that the sequence $\{\nu_k \phi'[g_i(\mathbf{x}_k)]\}$ must converge to some limit $\hat{\mu}_i$ as $k \rightarrow +\infty$ for all $i \in \mathcal{I}(\hat{\mathbf{x}})$. At last, since $\nu_k \phi'[g_i(\mathbf{x}_k)] \geq 0$ for all $k = 1, 2, \dots$, $i \in \mathcal{I}(\hat{\mathbf{x}})$, it follows that $\hat{\boldsymbol{\mu}} \geq \mathbf{0}^{|\mathcal{I}(\hat{\mathbf{x}})|}$.

Passing to the limit as $k \rightarrow +\infty$ in (13.12) we deduce that

$$\nabla f(\hat{\mathbf{x}}) + \sum_{i \in \mathcal{I}(\hat{\mathbf{x}})} \hat{\mu}_i \nabla g_i(\hat{\mathbf{x}}) = \mathbf{0}^n,$$

that is, $\hat{\mathbf{x}}$ is a KKT point for (13.1) with Lagrange multiplier vector $\hat{\boldsymbol{\mu}}$. ■

For example, if we use $\phi(s) := \phi_1(s) := -1/s$, then $\phi'(s) = 1/s^2$ in Theorem 13.6, and the sequence $\{\nu_k/g_i^2(\mathbf{x}_k)\}$ converges to the Lagrange multiplier $\hat{\mu}_i$ corresponding to the constraint i ($i = 1, \dots, m$).

13.1.3 Computational considerations

As the penalty parameter increases in the exterior penalty methods, or decreases in the interior penalty methods, the approximating problem (13.5) [respectively, (13.11)] becomes more and more ill-conditioned. Therefore, a typical computational strategy is to start from “safe” values of the penalty parameter (relatively small for exterior penalties, or

relatively large for barriers), and then proceed step after step slightly modifying the penalty parameter (e.g., multiplying it with some number close to 1).

It is natural to use the optimal solution $\mathbf{x}_{\nu_k}^*$ as a starting point for an iterative algorithm used to solve the approximating problem corresponding to the next value ν_{k+1} of the penalty parameter. The idea behind such a “warm start” is that, *typically*, $\nu_k \approx \nu_{k+1}$ implies $\mathbf{x}_{\nu_k}^* \approx \mathbf{x}_{\nu_{k+1}}^*$.

In fact, in many cases we can perform only few (maybe, only one) steps of an iterative algorithm starting at \mathbf{x}_{ν_k} to obtain a satisfactory approximation $\mathbf{x}_{\nu_{k+1}}$ of an optimal solution corresponding to the penalty parameter ν_{k+1} , and still preserve the convergence $\mathbf{x}_{\nu_k} \rightarrow \mathbf{x}^*$, as $k \rightarrow +\infty$, towards optimal solutions of the original constrained problem (13.1). This technique is especially applicable to convex optimization problems, and all the complexity estimates for interior penalty algorithms depend on this fact.

13.1.4 Applications and examples

Example 13.7 (exterior penalty) Consider the problem to

$$\begin{aligned} &\text{minimize } f(\mathbf{x}) := \frac{1}{2}(x_1^2 + x_2^2) + 2x_2, \\ &\text{subject to } x_2 = 0. \end{aligned} \quad (13.13)$$

The problem is convex with affine constraints; therefore, the KKT conditions are both necessary and sufficient for the global optimality. The KKT system in this case reduces to: $x_2 = 0$ and

$$\begin{pmatrix} x_1 \\ x_2 + 2 \end{pmatrix} + \lambda \begin{pmatrix} 0 \\ 1 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}.$$

The only solution to this system is $\mathbf{x} = \mathbf{0}^2$, $\lambda = -2$.

Let us use the exterior penalty method with quadratic penalty $\psi(s) := s^2$ to solve this problem. That is, we want to

$$\text{minimize } \frac{1}{2}(x_1^2 + x_2^2) + 2x_2 + \nu x_2^2,$$

where $\nu > 0$ is a penalty parameter. This problem is convex as well, so that stationarity is both necessary and sufficient for global optimality:

$$\begin{pmatrix} x_1 \\ (1 + 2\nu)x_2 + 2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix},$$

which has the unique solution $\mathbf{x}_{\nu}^* = (0, -2/(1 + 2\nu))^T$ for every $\nu > 0$. Note that $\lim_{\nu \rightarrow +\infty} \mathbf{x}_{\nu}^* = \mathbf{0}^2$ is a globally optimal solution to (13.13),

and that

$$\lim_{\nu \rightarrow +\infty} \nu \psi'[(\mathbf{x}_\nu^*)_2] = \lim_{\nu \rightarrow +\infty} \frac{-4\nu}{1+2\nu} = -2 = \lambda,$$

where λ is the Lagrange multiplier corresponding to the equality constraint $x_2 = 0$. ■

Example 13.8 (interior penalty) Consider the problem to

$$\begin{aligned} & \text{minimize } f(\mathbf{x}) := x_1^2 + x_2, \\ & \text{subject to } x_1^2 + x_2^2 - 1 \leq 0. \end{aligned} \quad (13.14)$$

The problem is convex and verifies Slater's CQ (see Definition 5.38); therefore, the KKT conditions are both necessary and sufficient for global optimality. The KKT system in this case reduces to: $x_1^2 + x_2^2 \leq 1$ and

$$\begin{aligned} \begin{pmatrix} 2x_1 \\ 1 \end{pmatrix} + \mu \begin{pmatrix} 2x_1 \\ 2x_2 \end{pmatrix} &= \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \\ \mu &\geq 0, \\ \mu(x_1^2 + x_2^2 - 1) &= 0. \end{aligned}$$

After easy calculations, which the reader is encouraged to perform, we can see that the only solution to this system is $\mathbf{x}^* = (0, -1)^T$, $\mu = 1/2$.

Now, let us use the barrier method with the barrier function $\phi(s) := -\log(-s)$. That is, we want to

$$\text{minimize } x_1^2 + x_2 - \nu \log(1 - x_1^2 - x_2^2),$$

where $\nu > 0$ is a penalty parameter. This problem is convex as well (verify this!), so that stationarity (restricted to the interior of the feasible set, $\{\mathbf{x} \in \mathbb{R}^2 \mid x_1^2 + x_2^2 < 1\}$) is both necessary and sufficient for global optimality:

$$\begin{pmatrix} 2x_1 \\ 1 \end{pmatrix} + \frac{\nu}{1 - x_1^2 - x_2^2} \begin{pmatrix} 2x_1 \\ 2x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}.$$

Again, after some transformations we can verify that this system has two solutions $\mathbf{x}_\nu = (0, \nu - \sqrt{\nu^2 + 1})^T$ and $\mathbf{y}_\nu = (0, \nu + \sqrt{\nu^2 + 1})^T$, out of which only \mathbf{x}_ν is (strictly) feasible. We can easily see that $\lim_{\nu \rightarrow +0} \mathbf{x}_\nu = (0, -1)^T$ is a globally optimal solution to (13.14), and that

$$\begin{aligned} \lim_{\nu \rightarrow +0} \nu \phi'[(\mathbf{x}_\nu)_1^2 + (\mathbf{x}_\nu)_2^2 - 1] &= \lim_{\nu \rightarrow +0} \frac{\nu}{1 - (\nu - \sqrt{\nu^2 + 1})^2} \\ &= \lim_{\nu \rightarrow +0} \frac{1}{2\sqrt{\nu^2 + 1} - 2\nu} = \frac{1}{2} = \mu, \end{aligned}$$

where μ is the Lagrange multiplier corresponding to the inequality constraint $x_1^2 + x_2^2 - 1 \leq 0$. ■

Example 13.9 (linear programming) Consider an LP problem of the following form:

$$\begin{aligned} & \text{maximize } \mathbf{b}^T \mathbf{y}, \\ & \text{subject to } \mathbf{A}^T \mathbf{y} \leq \mathbf{c}, \end{aligned} \quad (13.15)$$

where $\mathbf{b}, \mathbf{y} \in \mathbb{R}^m$, $\mathbf{c} \in \mathbb{R}^n$, and $\mathbf{A} \in \mathbb{R}^{m \times n}$. Using standard linear programming theory (see Theorem 10.15), we can write the primal–dual optimality conditions for this problem in the form:

$$\begin{aligned} & \mathbf{A}^T \mathbf{y} \leq \mathbf{c}, \\ & \mathbf{A} \mathbf{x} = \mathbf{b}, \\ & \mathbf{x} \geq \mathbf{0}^n, \\ & \mathbf{x}^T (\mathbf{c} - \mathbf{A}^T \mathbf{y}) = 0, \end{aligned} \quad (13.16)$$

where $\mathbf{x} \in \mathbb{R}^n$ is the vector of Lagrange multipliers for the inequality constraints, or just a vector of dual variables as it is customary called in the linear programming literature.

Introducing a vector $\mathbf{s} \in \mathbb{R}^n$ of slack variables for the inequality constraints into the problem (13.15), it assumes the form

$$\begin{aligned} & \text{maximize } \mathbf{b}^T \mathbf{y}, \\ & \text{subject to } \mathbf{A}^T \mathbf{y} + \mathbf{s} = \mathbf{c}, \\ & \mathbf{s} \geq \mathbf{0}^n, \end{aligned} \quad (13.17)$$

and the corresponding system of optimality conditions will be:

$$\begin{aligned} & \mathbf{A}^T \mathbf{y} + \mathbf{s} = \mathbf{c}, \\ & \mathbf{A} \mathbf{x} = \mathbf{b}, \\ & \mathbf{x} \geq \mathbf{0}^n, \mathbf{s} \geq \mathbf{0}^n, \mathbf{x}^T \mathbf{s} = 0. \end{aligned} \quad (13.18)$$

Now, let us apply the barrier method to the optimization problem (13.17). It has equality constraints, which we do not move into the penalty function, but rather leave them as they are. Thus, we consider the following problem with equality constraints only:

$$\begin{aligned} & \text{minimize } -\mathbf{b}^T \mathbf{y} - \nu \sum_{j=1}^n \log(s_j), \\ & \text{subject to } \mathbf{A}^T \mathbf{y} + \mathbf{s} = \mathbf{c}, \end{aligned} \quad (13.19)$$

where we use the logarithmic barrier function, $\nu > 0$ is a penalty parameter, and we have multiplied the original objective function with

-1 to change the maximization problem into a minimization one. The problem (13.19) is convex with affine constraints, therefore the KKT conditions are both necessary and sufficient for the global optimality. The KKT system in this case is: $\mathbf{A}^T \mathbf{y} + \mathbf{s} = \mathbf{c}$, and

$$\begin{pmatrix} -\mathbf{b} \\ -\nu/s_1 \\ \vdots \\ -\nu/s_n \end{pmatrix} + \begin{pmatrix} \mathbf{A} \\ \mathbf{I}^n \end{pmatrix} \mathbf{x} = \begin{pmatrix} \mathbf{0}^m \\ \mathbf{0}^n \end{pmatrix}, \quad (13.20)$$

where $\mathbf{x} \in \mathbb{R}^n$ is a vector of Lagrange multipliers for the equality constraints in the problem (13.19). Further, the system (13.20) can be rewritten in the following more convenient form:

$$\begin{aligned} \mathbf{A}^T \mathbf{y} + \mathbf{s} &= \mathbf{c}, \\ \mathbf{A} \mathbf{x} &= \mathbf{b}, \\ x_j s_j &= \nu, \quad j = 1, \dots, n. \end{aligned} \quad (13.21)$$

Recalling that due to the presence of the barrier the vector \mathbf{s} must be strictly feasible, that is, $\mathbf{s} > \mathbf{0}^n$, and that the penalty parameter ν is positive, the last equation in (13.21) does in fact imply the strict inequality $\mathbf{x} > \mathbf{0}^n$.

Therefore, comparing (13.21) and (13.18) we see that for linear programs the introduction of a logarithmic barrier amounts to a small perturbation of the complementarity condition. Namely, instead of the requirement

$$\mathbf{x} \geq \mathbf{0}^n, \quad \mathbf{s} \geq \mathbf{0}^n, \quad x_j s_j = 0, \quad j = 1, \dots, n,$$

we get a similar one (for small $\nu > 0$):

$$\mathbf{x} > \mathbf{0}^n, \quad \mathbf{s} > \mathbf{0}^n, \quad x_j s_j = \nu, \quad j = 1, \dots, n.$$

For the case $n = 1$ the difference between the two is shown in Figure 13.3. Note the smoothing effect on the feasible set introduced by the interior penalty algorithm. We can use Newton's method to solve the system of nonlinear equations (13.21), but not (13.18). ■

13.2 Sequential quadratic programming

13.2.1 Introduction

We begin by studying the equality constrained problem to

$$\text{minimize } f(\mathbf{x}), \quad (13.22a)$$

$$\text{subject to } h_j(\mathbf{x}) = 0, \quad j = 1, \dots, \ell, \quad (13.22b)$$

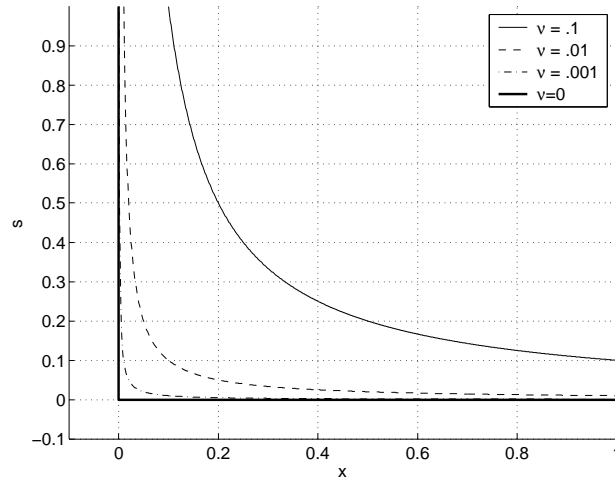


Figure 13.3: The approximation of the complementarity constraint resulting from the use of logarithmic barrier functions in linear programming.

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $h_j : \mathbb{R}^n \rightarrow \mathbb{R}$ are given functions in C^1 on \mathbb{R}^n . The Karush–Kuhn–Tucker conditions for this problem state (cf. Theorem 5.33) that at a local minimum \mathbf{x}^* of f over the feasible set, where \mathbf{x}^* satisfies some constraint qualification, there exists a vector $\boldsymbol{\lambda}^* \in \mathbb{R}^\ell$ such that

$$\nabla_{\mathbf{x}} L(\mathbf{x}^*, \boldsymbol{\lambda}^*) := \nabla f(\mathbf{x}^*) + \sum_{j=1}^{\ell} \lambda_j^* \nabla h_j(\mathbf{x}^*) = \mathbf{0}^n, \quad (13.23a)$$

$$\nabla_{\boldsymbol{\lambda}} L(\mathbf{x}^*, \boldsymbol{\lambda}^*) := \mathbf{h}(\mathbf{x}^*) = \mathbf{0}^\ell. \quad (13.23b)$$

It is an appealing idea to find such a point by directly attacking this system of nonlinear equations, which has $n + \ell$ unknowns as well as equations. Newton's method is then the natural choice. Let us see what the Newton subproblem looks like. We now assume, for the time being, that f and h_j , $j = 1, \dots, \ell$, are in C^2 on \mathbb{R}^n . Suppose that $(\mathbf{x}_k, \boldsymbol{\lambda}_k) \in \mathbb{R}^n \times \mathbb{R}^\ell$. Then, since Newton's method takes a unit step in the direction towards the approximate problem's solution, we obtain the following characterization of the next iterate $(\mathbf{x}_{k+1}, \boldsymbol{\lambda}_{k+1})$: $(\mathbf{x}_{k+1}, \boldsymbol{\lambda}_{k+1}) := (\mathbf{x}_k, \boldsymbol{\lambda}_k) + (\mathbf{p}_k, \mathbf{v}_k)$, where $(\mathbf{p}_k, \mathbf{v}_k) \in \mathbb{R}^n \times \mathbb{R}^\ell$ solves the second-order approximation of the stationary point condition for the

Lagrange function:

$$\nabla^2 L(\mathbf{x}_k, \boldsymbol{\lambda}_k) \begin{pmatrix} \mathbf{p}_k \\ \mathbf{v}_k \end{pmatrix} = -\nabla L(\mathbf{x}_k, \boldsymbol{\lambda}_k),$$

that is,

$$\begin{bmatrix} \nabla_{xx}^2 L(\mathbf{x}_k, \boldsymbol{\lambda}_k) & \nabla \mathbf{h}(\mathbf{x}_k) \\ \nabla \mathbf{h}(\mathbf{x}_k)^T & \mathbf{0}^{\ell \times \ell} \end{bmatrix} \begin{pmatrix} \mathbf{p}_k \\ \mathbf{v}_k \end{pmatrix} = \begin{pmatrix} -\nabla_x L(\mathbf{x}_k, \boldsymbol{\lambda}_k) \\ -\mathbf{h}(\mathbf{x}_k) \end{pmatrix}, \quad (13.24)$$

where the matrix $\nabla \mathbf{h}(\mathbf{x}_k)^T$ is the Jacobian of \mathbf{h} at \mathbf{x}_k , comprised of the rows $\nabla h_j(\mathbf{x}_k)^T$ for $j = 1, \dots, \ell$.

This system of linear equations has a nice interpretation, namely as the KKT system corresponding to the quadratic programming problem to

$$\underset{\mathbf{p}}{\text{minimize}} \quad \frac{1}{2} \mathbf{p}^T \nabla_{xx}^2 L(\mathbf{x}_k, \boldsymbol{\lambda}_k) \mathbf{p} + \nabla_x L(\mathbf{x}_k, \boldsymbol{\lambda}_k) \mathbf{p}, \quad (13.25a)$$

$$\text{subject to } h_j(\mathbf{x}_k) + \nabla h_j(\mathbf{x}_k)^T \mathbf{p} = 0, \quad j = 1, \dots, \ell. \quad (13.25b)$$

This approximate problem has as its objective a second-order approximation of the Lagrange function with respect to the primal variables \mathbf{x} , and the original constraints have been replaced by their first-order approximations at \mathbf{x}_k . The Lagrange multiplier vector \mathbf{v}_k appearing in (13.24) is the vector of Lagrange multipliers for the constraints (13.25b).

As for Newton methods in unconstrained optimization, convergence to a stationary point of the Lagrangian in $\mathbb{R}^n \times \mathbb{R}^\ell$ requires (unless some sort of line search is introduced) that we start the algorithm close to such a point and where also the Hessian of the Lagrangian is invertible so that the algorithm is well-defined. Under the additional conditions that the stationary point \mathbf{x}^* is a strict minimum of f over the feasible set, that it satisfies the linear independence constraint qualification LICQ (see Definition 5.41), and that it together with the KKT multiplier vector $\boldsymbol{\lambda}^*$ satisfies a second-order sufficient condition (cf. Theorem 4.17), the sequence $\{(\mathbf{x}_k, \boldsymbol{\lambda}_k)\}$ converges towards the KKT point $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$ with a superlinear rate (cf. Section 11.10).

We remark that the convergence theory presented for the above rudimentary Newton method is far from satisfactory, for several reasons:

- Convergence is only *local*, which means that the algorithm must be combined with an algorithm that converges to a KKT point from any starting vector, that is, a *global* algorithm.
- The algorithm requires strong assumptions about the problem, such as that the functions f and h_j are in C^2 and that the Hessian of the Lagrangian is positive definite, in order for the solution to (13.25) to be well-defined.

In the next section, we will therefore develop a modification of the above algorithm, which is globally convergent to stationary points. Moreover, we will work also with inequality constraints, which is not immediate to incorporate into the above Newton-like framework.

13.2.2 A penalty-function based SQP algorithm

In order to introduce a penalty function into the discussion, let us consider first the following one:

$$P(\mathbf{x}, \boldsymbol{\lambda}) := \|\nabla_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda})\|^2 + \|\mathbf{h}(\mathbf{x})\|^2. \quad (13.26)$$

This is an *exact penalty function*, because its unconstrained minima are (or, strongly relate to) optimal solutions and/or Lagrange multipliers of the constrained problem. The exact penalty function (13.26) has been used extensively in cases where $\ell = n$ and the problem is to find a solution to $h_j(\mathbf{x}) = 0$ for all j . The function has significant drawbacks, however: it does not distinguish between local minima and maxima, and it may have local minima that are not global and even do not correspond to vectors where the value of P is zero; in other words, it may have local minima that are even infeasible in the original problem.

The above case of penalty function is differentiable; the more popular penalty functions are non-differentiable. We present such a one next.

Consider the constrained optimization problem to

$$\text{minimize } f(\mathbf{x}), \quad (13.27a)$$

$$\text{subject to } g_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m, \quad (13.27b)$$

$$h_j(\mathbf{x}) = 0, \quad j = 1, \dots, \ell, \quad (13.27c)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$, and $h_j : \mathbb{R}^n \rightarrow \mathbb{R}$ are given functions in C^1 on \mathbb{R}^n . We introduce the l_1 penalty function [cf. (13.4)]

$$\check{\chi}_S(\mathbf{x}) := \sum_{i=1}^m \text{maximum} \{0, g_i(\mathbf{x})\} + \sum_{j=1}^{\ell} |h_j(\mathbf{x})|, \quad (13.28)$$

and the associated exact penalty function

$$P_e(\mathbf{x}) := f(\mathbf{x}) + \nu \check{\chi}_S(\mathbf{x}), \quad (13.29)$$

where $\nu > 0$.

Proposition 13.10 (an exact penalty function) *Suppose that \mathbf{x}^* satisfies the KKT conditions (5.17) of the problem (13.27), together with*

Lagrange multipliers $(\boldsymbol{\mu}^*, \boldsymbol{\lambda}^*)$. Suppose further that the functions f and g_i , $i \in \mathcal{I}(\mathbf{x}^*)$, all are convex, and that h_j , $j = 1, \dots, \ell$, are affine. If the value of ν is large enough such that

$$\nu \geq \text{maximum}\{\mu_i^*, i \in \mathcal{I}(\mathbf{x}^*); |\lambda_j^*|, j = 1, \dots, \ell\},$$

then the vector \mathbf{x}^* is also a global minimum of the function P_e .

Proof. [Sketch] Consider the problem of minimizing P_e over \mathbb{R}^n . We can rewrite this problem as follows:

$$\text{minimize} \quad f(\mathbf{x}) + \nu \left(\sum_{i=1}^m y_i + \sum_{j=1}^{\ell} z_j \right), \quad (13.30a)$$

$$\text{subject to} \quad y_i \geq g_i(\mathbf{x}) \text{ and } y_i \geq 0, \quad i = 1, \dots, m, \quad (13.30b)$$

$$z_j \geq h_j(\mathbf{x}) \text{ and } z_j \geq -h_j(\mathbf{x}), \quad j = 1, \dots, \ell. \quad (13.30c)$$

Analyzing the KKT conditions for this problem, we can construct multipliers for the problem (13.30) from the multiplier vectors $(\boldsymbol{\mu}^*, \boldsymbol{\lambda}^*)$ and show that \mathbf{x}^* is a globally optimal solution to it (note the convexity assumptions). ■

There are similar results also for more general, non-convex, problems that establish that if \mathbf{x}^* is a (strict) local minimum to (13.27) then it is also a (strict) local minimum of the exact penalty function.

We must note, however, that the implication is in a somewhat unsatisfactory direction: there may exist local minima of P_e that do not correspond to constrained local minima in the original problem, for any value of ν . The theory is much more satisfactory in the convex case.

We develop a penalty SQP algorithm, known as the MSQP method (as in Merit SQP, *merit function* being synonymous with objective function), for solving the general problem (13.27). Given an iterate $\mathbf{x}_k \in \mathbb{R}^n$ and a vector $(\boldsymbol{\mu}_k, \boldsymbol{\lambda}_k) \in \mathbb{R}_+^m \times \mathbb{R}^\ell$, suppose we choose a positive definite, symmetric matrix $\mathbf{B}_k \in \mathbb{R}^{n \times n}$; for example, it can be an approximation of $\nabla_{xx}^2 L(\mathbf{x}_k, \boldsymbol{\mu}_k, \boldsymbol{\lambda}_k)$. We then solve the following subproblem:

$$\text{minimize}_{\mathbf{p}} \quad \frac{1}{2} \mathbf{p}^T \mathbf{B}_k \mathbf{p} + \nabla f(\mathbf{x}_k)^T \mathbf{p}, \quad (13.31a)$$

$$\text{subject to} \quad g_i(\mathbf{x}_k) + \nabla g_i(\mathbf{x}_k)^T \mathbf{p} \leq 0, \quad i = 1, \dots, m, \quad (13.31b)$$

$$h_j(\mathbf{x}_k) + \nabla h_j(\mathbf{x}_k)^T \mathbf{p} = 0, \quad j = 1, \dots, \ell. \quad (13.31c)$$

Note that if we were to utilize $\mathbf{B}_k := \nabla_{xx}^2 L(\mathbf{x}_k, \boldsymbol{\mu}_k, \boldsymbol{\lambda}_k)$ then the problem (13.31) would be the optimization problem associated with a second-order approximation of the KKT conditions for the original problem

(13.27); the close connection to quasi-Newton methods in unconstrained optimization should be obvious.

We also took the liberty to replace the term $\nabla_x L(\mathbf{x}_k, \boldsymbol{\mu}_k, \boldsymbol{\lambda}_k)^T \mathbf{p}$ by the term $\nabla f(\mathbf{x}_k)^T \mathbf{p}$. This is without any loss of generality, as the KKT conditions for the problem (13.31) imply that they can be interchanged without any loss of generality—the only difference in the two objectives lies in the constant term which plays no role in the optimization.

A quasi-Newton type method based on the subproblem (13.31) followed by a unit step and a proper update of the matrix \mathbf{B}_k , as in the BFGS algorithm, is locally convergent with a superlinear speed, just as in the unconstrained case. But we are still interested in a globally convergent version, whence we develop the theory of an algorithm that utilizes the exact penalty function (13.29) in a line search rather than taking a unit step. Our first result shows when the subproblem solution provides a descent direction with respect to this function.

Lemma 13.11 (a descent property) *Given $\mathbf{x}_k \in \mathbb{R}^n$ consider the strictly convex quadratic problem (13.31), where $\mathbf{B}_k \in \mathbb{R}^{n \times n}$ is symmetric and positive definite. Suppose that \mathbf{p}_k solves this problem together with the multipliers $\boldsymbol{\mu}$ and $\boldsymbol{\lambda}$. Assume that $\mathbf{p}_k \neq \mathbf{0}^n$. If*

$$\nu \geq \text{maximum} \{ \mu_1, \dots, \mu_m, |\lambda_1|, \dots, |\lambda_\ell| \}$$

then the vector \mathbf{p}_k is a direction of descent with respect to the l_1 penalty function (13.29) at $(\mathbf{x}_k, \boldsymbol{\mu}_k, \boldsymbol{\lambda}_k)$.

Proof. Using the KKT conditions of the problem (13.31) we obtain that

$$\begin{aligned} \nabla f(\mathbf{x}_k)^T \mathbf{p} &= -\mathbf{p}^T \mathbf{B}_k \mathbf{p} - \sum_{i=1}^m \mu_i \nabla g_i(\mathbf{x}_k)^T \mathbf{p} - \sum_{j=1}^{\ell} \lambda_j \nabla h_j(\mathbf{x}_k)^T \mathbf{p} \\ &= -\mathbf{p}^T \mathbf{B}_k \mathbf{p} + \sum_{i=1}^m \mu_i g_i(\mathbf{x}_k) + \sum_{j=1}^{\ell} \lambda_j h_j(\mathbf{x}_k) \\ &\leq -\mathbf{p}^T \mathbf{B}_k \mathbf{p} + \sum_{i=1}^m \mu_i \text{maximum} \{0, g_i(\mathbf{x}_k)\} + \sum_{j=1}^{\ell} |\lambda_j| |h_j(\mathbf{x}_k)| \\ &\leq -\mathbf{p}^T \mathbf{B}_k \mathbf{p} + \nu \left(\sum_{i=1}^m \text{maximum} \{0, g_i(\mathbf{x}_k)\} + \sum_{j=1}^{\ell} |h_j(\mathbf{x}_k)| \right). \end{aligned}$$

In order to investigate the descent properties of \mathbf{p}_k with respect to

P_e at \mathbf{x}_k , we next note that

$$\begin{aligned} P_e(\mathbf{x}_k) - P_e(\mathbf{x}_k + \alpha \mathbf{p}_k) &= f(\mathbf{x}_k) - f(\mathbf{x}_k + \alpha \mathbf{p}_k) \\ &\quad + \nu \sum_{i=1}^m [\max\{0, g_i(\mathbf{x}_k)\} - \max\{0, g_i(\mathbf{x}_k + \alpha \mathbf{p}_k)\}] \\ &\quad + \nu \sum_{j=1}^{\ell} [|h_j(\mathbf{x}_k)| - |h_j(\mathbf{x}_k + \alpha \mathbf{p}_k)|]. \end{aligned}$$

Let O_i (respectively, O_j) denote functions from \mathbb{R} to \mathbb{R} each function O_r having the property that $\lim_{\alpha \rightarrow 0} O_r(\alpha) = 0$, and specially chosen such that the below identities follow. Then, for $\alpha > 0$ small enough,

$$f(\mathbf{x}_k + \alpha \mathbf{p}_k) = f(\mathbf{x}_k) + \alpha \nabla f(\mathbf{x}_k)^T \mathbf{p}_k + \alpha O_0(\alpha).$$

Also, $g_i(\mathbf{x}_k + \alpha \mathbf{p}_k) = g_i(\mathbf{x}_k) + \alpha \nabla g_i(\mathbf{x}_k)^T \mathbf{p}_k + \alpha O_i(\alpha) \leq g_i(\mathbf{x}_k) - \alpha g_i(\mathbf{x}_k) + \alpha O_i(\alpha)$ holds by the KKT conditions of the problem (13.31). Hence,

$$\text{maximum}\{0, g_i(\mathbf{x}_k + \alpha \mathbf{p}_k)\} \leq (1 - \alpha) \text{maximum}\{0, g_i(\mathbf{x}_k)\} + \alpha |O_i(\alpha)|.$$

Similarly we obtain that $h_j(\mathbf{x}_k + \alpha \mathbf{p}_k) = h_j(\mathbf{x}_k) + \alpha \nabla h_j(\mathbf{x}_k)^T \mathbf{p}_k + \alpha O_j(\alpha) = (1 - \alpha)h_j(\mathbf{x}_k) + \alpha O_j(\alpha)$, and hence

$$|h_j(\mathbf{x}_k + \alpha \mathbf{p}_k)| \leq (1 - \alpha)|h_j(\mathbf{x}_k)| + \alpha |O_j(\alpha)|.$$

Using these three expressions in the expression for $P_e(\mathbf{x}_k) - P_e(\mathbf{x}_k + \alpha \mathbf{p}_k)$ we obtain that for small enough $\alpha > 0$, $P_e(\mathbf{x}_k) - P_e(\mathbf{x}_k + \alpha \mathbf{p}_k) \geq \alpha [-\nabla f(\mathbf{x}_k)^T \mathbf{p}_k + \nu \sum_{i=1}^m \text{maximum}\{0, g_i(\mathbf{x}_k)\} + \nu \sum_{j=1}^{\ell} |h_j(\mathbf{x}_k)| + O(\alpha)]$. Hence, we obtain that

$$P_e(\mathbf{x}_k) - P_e(\mathbf{x}_k + \alpha \mathbf{p}_k) \geq \alpha [\mathbf{p}_k^T \mathbf{B}_k \mathbf{p}_k + O(\alpha)] > 0$$

for every $\alpha > 0$ small enough, due to the positive definiteness of the matrix \mathbf{B}_k . We are done. \blacksquare

The MSQP algorithm then works as follows. At some iteration k , we have at hand a vector \mathbf{x}_k . Select a symmetric and positive definite matrix $\mathbf{B}_k \in \mathbb{R}^{n \times n}$. Solve the QP problem (13.31) in order to obtain the vector \mathbf{p}_k and multipliers $(\boldsymbol{\mu}_{k+1}, \boldsymbol{\lambda}_{k+1})$. If $\mathbf{p}_k = \mathbf{0}^n$ we stop with \mathbf{x}_k being a KKT point for the original problem (13.27) together with the multipliers $(\boldsymbol{\mu}_{k+1}, \boldsymbol{\lambda}_{k+1})$. (Why?) Otherwise, find $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$ where α_k minimizes $P_e(\mathbf{x}_k + \alpha \mathbf{p}_k)$ over $\alpha \geq 0$. Increase k by one and repeat.

Convergence of this rudimentary algorithm follows below.

Theorem 13.12 (convergence of the MSQP method) *The algorithm MSQP either terminates finitely at a KKT point for the problem (13.27) or it produces an infinite sequence $\{\mathbf{x}_k\}$. In the latter case, we assume that $\{\mathbf{x}_k\}$ lies in a compact set $X \subset \mathbb{R}^n$ and that for every $\mathbf{x} \in X$ and symmetric and positive definite matrix \mathbf{B}_k the QP (13.31) has a unique solution, and also unique multiplier vectors $\boldsymbol{\mu}$ and $\boldsymbol{\lambda}$ satisfying $\nu \geq \max\{\mu_1, \dots, \mu_m, |\lambda_1|, \dots, |\lambda_\ell|\}$, where $\nu > 0$ is the penalty parameter. Furthermore, assume that the sequence $\{\mathbf{B}_k\}$ of matrices is bounded and that every limit point of this sequence is positive definite (or, the sequence $\{\mathbf{B}_k^{-1}\}$ of matrices is bounded). Then, every limit point of $\{\mathbf{x}_k\}$ is a KKT point for the problem (13.27).*

Proof. [Sketch] Clearly, the algorithm stops precisely at KKT points, so we concentrate on the case where $\{\mathbf{x}_k\}$ is an infinite sequence. We can consider an iteration as a descent step wherein we first construct a descent direction \mathbf{p}_k , followed by a line search in the continuous function P_e , and followed by an update of the matrix \mathbf{B}_k . By the properties stated, each of these steps is well defined.

Since the sequence $\{\mathbf{x}_k\}$ is bounded, it has a limit point, say \mathbf{x}^∞ . Consider from now on this subsequence. By the assumptions stated, also the sequence $\{\mathbf{p}_k\}$ must be bounded. (Why?) Suppose that \mathbf{p}^∞ is a limit point of $\{\mathbf{p}_k\}$ within the above-mentioned subsequence. Suppose that it is non-zero. By assumption the sequence $\{\mathbf{B}_k\}$ also has limit points within this subsequence, all of which are positive definite. Suppose \mathbf{B}^∞ is one such matrix. Then, by Lemma 13.11 the vector \mathbf{p}^∞ must define a descent direction for P_e . This contradicts the assumption that \mathbf{x}^∞ is a limit point. (Why?) Therefore, it must be the case that $\mathbf{p}^\infty = \mathbf{0}^n$, in which case \mathbf{x}^∞ is stationary, that is, a KKT point. We are done. ■

Note that we here have not described any rules for selecting the value of ν . Clearly, this is a difficult task, which must be decided upon from experiments including the results from the above line searches with respect to the merit function P_e . Further, we have no guarantees that the QP subproblems (13.31) are feasible; in the above theorem we *assumed* that the problem is well-defined. Further still, P_e is only continuous and directionally differentiable, whence we cannot utilize several of the step length rules devised in Section 11.3. Local superlinear or quadratic convergence of this algorithm can actually be impaired due to the use of this merit function, as it is possible to construct examples where a unit step does not reduce its value even very close to an optimal solution. (This is known as the *Maratos effect*, after [Mar78].) The Notes Section 13.4 lead to further reading on these issues.

13.2.3 A numerical example on the MSQP algorithm

Consider the two-dimensional optimization problem to

$$\text{minimize } f(\mathbf{x}) := 2x_1^2 + 2x_2^2 - 2x_1x_2 - 4x_1 - 6x_2, \quad (13.32a)$$

$$\text{subject to } g_1(\mathbf{x}) := 2x_1^2 - x_2 \leq 0, \quad (13.32b)$$

$$g_2(\mathbf{x}) := x_1 + 5x_2 - 5 \leq 0, \quad (13.32c)$$

$$g_3(\mathbf{x}) := -x_1 \leq 0, \quad (13.32d)$$

$$g_4(\mathbf{x}) := -x_2 \leq 0. \quad (13.32e)$$

Check that the vector $(\frac{7}{3}, \frac{8}{3})^T$ is an “unconstrained” globally optimal solution, which however is infeasible.

Suppose we wish to utilize the MSQP algorithm for solving this problem. We choose $\nu := 10$, \mathbf{B}_k to always be the partial Hessian $\nabla_{xx}^2 L(\mathbf{x}_k, \boldsymbol{\mu}_k)$ of the Lagrangian (notice that it always is positive definite due to the convexity properties of the problem), and the starting point $\mathbf{x}_0 := (0, 1)^T$, which is feasible. Hence, $f(\mathbf{x}_0) = P_e(\mathbf{x}_0) = -4$. We also select $\boldsymbol{\mu}_0 := \mathbf{0}^4$. Setting up the first QP subproblem accordingly, we obtain the problem to

$$\text{minimize } \frac{1}{2}(4p_1^2 + 4p_2^2 - 4p_1p_2) - 6p_1 - 2p_2, \quad (13.33a)$$

$$\text{subject to } -1 - p_2 \leq 0, \quad (13.33b)$$

$$p_1 + 5p_2 \leq 0, \quad (13.33c)$$

$$-p_1 \leq 0, \quad (13.33d)$$

$$-1 - p_2 \leq 0. \quad (13.33e)$$

Solving this problem yields the solution $\mathbf{p}_1 = (\frac{35}{31}, -\frac{7}{31})^T$ and the multiplier vector $\boldsymbol{\mu}_1 \approx (0, 1.032258, 0, 0)^T$.

We next perform a line search in the exact penalty function:

$$\begin{aligned} \text{minimize}_{\alpha \geq 0} P_e(\mathbf{x}_0 + \alpha \mathbf{p}_1) &= 3.1612897\alpha^2 - 6.3225804\alpha - 4 \\ &\quad + 10 \max\{0, 2.5494274\alpha^2 + 0.2258064\alpha - 1\} \\ &\quad + 10 \max\{0, 0\} + 10 \max\{0, -1.1290322\alpha\} \\ &\quad + 10 \max\{0, -1 + 0.2258064\alpha\}. \end{aligned}$$

We obtain that $\alpha_1 \approx 0.5835726$. (Note that the unconstrained minimum of $\alpha \mapsto f(\mathbf{x}_0 + \alpha \mathbf{p}_1)$ is $\alpha = 1$, which however leads to an infeasible point having a too high penalty.)

This produces the next iterate, $\mathbf{x}_1 \approx (0.6588722, 0.8682256)^T$.

We ask the reader to confirm that this is a near-optimal solution by checking the KKT conditions, and to confirm that the next QP problem verifies this.

We were able to find the optimal solution this quickly, due to the facts that the problem is quadratic and that the value $\nu = 10$ is large enough. (Check the value of the Lagrange multipliers.)

13.2.4 On recent developments in SQP algorithms

We have seen that the SQP algorithm above has an inherent decision problem, namely to choose the right value of the penalty parameter ν . In recent years, there has been a development of algorithms where the penalty parameter is avoided altogether. We call such methods *filter-SQP methods*.

In such methods we borrow a term from *multi-objective optimization*, and say that \mathbf{x}^1 *dominates* \mathbf{x}^2 if $\tilde{\chi}(\mathbf{x}^1) \leq \tilde{\chi}(\mathbf{x}^2)$ and $f(\mathbf{x}^1) \leq f(\mathbf{x}^2)$ [where $\tilde{\chi} = \tilde{\chi}_S$ is our measure of infeasibility], that is, if \mathbf{x}^1 is at least as good, both in terms of feasibility and optimality. A *filter* is a list of pairs $(\tilde{\chi}_i, f_i)$ such that $\tilde{\chi}_i < \tilde{\chi}_j$ or $f_i < f_j$ for all $j \neq i$ in the list. By adding elements to the filter, we build up an *efficient frontier*, that is, the *Pareto set* in the bi-criterion problem of simultaneously finding low objective values and reduce the infeasibility. The filter is used in place of the penalty function, when the standard Newton-like step cannot be computed, for example because the subproblem is infeasible.

This algorithm class is quickly becoming popular, and has already been found to be among the best general algorithms in nonlinear programming, especially because it does not rely on any parameters that need to be estimated.

13.3 A summary and comparison

Quite a few algorithms of the penalty and SQP type exist, of which only a small number could be summarized here. Which are the relative strengths and weaknesses of these methods?

First, we may contrast the two types of methods with regards to their ill-conditioning. The barrier methods of Section 13.1.2 solve a sequence of unconstrained optimization problems that become more and more ill-conditioned. In contrast, exact penalty methods need not be ill-conditioned and moreover only one approximate problem is, at least in principle, enough to solve the original problem. However, it is known at least for linear and quadratic programming problems that the inherent ill-conditioning of barrier methods can be eliminated (we say the

ill-conditioning is *benign*), because of the special structure of these problems and their optimality conditions.

Among the features of SQP methods is that they can deal with very large classes of problems, including those with nonlinear equality constraints, and they do not rely on the existence of second-order derivatives—although they can make good use of them. While it is known from practice that the number of quadratic subproblems can be rather small before reaching a near-locally optimal solution, these subproblems can be costly to solve. A major development has been made of specialized quadratic programming methods for solving and re-optimizing large-scale quadratic SQP subproblems, and the most recent codes are quite robust. Still, which methods to prefer depend on many factors.

The solver `fmincon` in the MATLAB Optimization Toolbox is an SQP method.

13.4 Notes and further reading

Exterior and interior penalty methods were popularized by the book *Nonlinear Programming: Sequential Unconstrained Minimization Technique* by Fiacco and McCormick [FiM68], although barrier methods had been presented already in 1961 and exterior penalty methods were developed by Courant much earlier still (in 1943). Their name for many years was “SUMT”, after the title of the book. These methods lost popularity when the classes of SQP and augmented Lagrangian methods (see Exercise 13.8 below) had begun to mature, but following the discovery of the polynomial complexity of certain interior point methods for LP, and in particular the discovery that some of them could be derived as special barrier methods where the barrier parameter is updated in a special way, made them popular again. Most text books on nonlinear optimization concentrate on these two classes of methods; see [Lue84, Fle87, BSS93, Ber99, NoW99].

In linear programming, several recent text books offer descriptions of interior point methods; see, e.g., [Pad99, Van01].

Sequential quadratic programming (SQP) methods were first developed by Wilson [Wil63]. Sequential linear programming (SLP) methods (cf. Exercise 13.7 below), which are based on first-order approximations of the KKT conditions, were developed by staff in the chemical (especially oil) industry; one reason why SLP methods are effective in such applications is that some important blending problems are only mildly nonlinear. The MSQP method described here stems from [Han75, Pow78]. A formal proof of Theorem 13.12 is given in [BSS93, Theorem 10.4.2].

The issue of feasibility of the SQP subproblems is taken up in [Fle87]. The boundedness of the subproblem solution is often ensured by combining SQP with a trust region method (cf. Section 11.7), such that the QP subproblem is further constrained. The Maratos effect has been overcome during the last decade of research; cf. [PaT91, Fac95]. An excellent paper which addresses most of the computational issues within an SQP algorithm and provides a very good compromise in the form of the SNOPT software is [GMS05]. Filter-SQP algorithms offer a substantial development over the standard SQP methods. Good references to this rapidly developing class of methods are [FLT02, UUV04].

We recommend a visit to the *NEOS Server for Optimization* at <http://www-neos.mcs.anl.gov/neos/> for a continuously updated list of optimization solvers, together with an excellent software guide for several types of classes of optimization models.

13.5 Exercises

Exercise 13.1 (convexity, exterior penalty method) Assume that the problem (13.1) is convex. Show that with the choice $\chi(s) := s^2$ [where χ enters the definition of the penalty function via (13.4)], for every $\nu > 0$ the problem (13.5) is convex.

Exercise 13.2 (convexity, interior penalty method) Assume that the problem (13.1) is convex. Show that with the choice $\phi(s) := -\log(-s)$ [where ϕ enters the definition of the penalty function via (13.10)], for every $\nu > 0$ the problem (13.11) is convex.

Exercise 13.3 (numerical example, exterior penalty method) Consider the problem to

$$\begin{aligned} &\text{minimize } f(\mathbf{x}) := \frac{1}{2} (x_1^2 + x_2^2), \\ &\text{subject to } x_1 = 1. \end{aligned}$$

Apply the exterior penalty method with the standard quadratic penalty function.

Exercise 13.4 (numerical example, logarithmic barrier method) Consider the problem to

$$\begin{aligned} &\text{minimize } f(\mathbf{x}) := \frac{1}{2} (x_1^2 + x_2^2), \\ &\text{subject to } x_1 \leq 1. \end{aligned}$$

Apply the interior penalty method with a logarithmic penalty function on the constraint.

Exercise 13.5 (logarithmic barrier, exam 990827) Consider the problem to

$$\begin{aligned} & \text{minimize } f(\mathbf{x}) := \frac{1}{2}x_1^2 + x_2^2, \\ & \text{subject to } x_1 + 2x_2 \geq 10. \end{aligned}$$

Attack this problem with a logarithmic barrier method. Describe explicitly the trajectory the method follows, as a function of the barrier parameter. Confirm that the limit point of the trajectory solves the problem.

Exercise 13.6 (logarithmic barrier method in linear programming) Consider the linear programming problem to

$$\begin{aligned} & \text{minimize } f(\mathbf{x}) := -y_1 + y_2, \\ & \text{subject to } y_2 \leq 1, \\ & \quad -y_1 \leq -1, \\ & \quad \mathbf{y} \geq \mathbf{0}^2. \end{aligned}$$

Apply the interior penalty method with a logarithmic penalty function on the non-negativity restrictions on the slack variables.

Exercise 13.7 (sequential linear programming) Consider the optimization problem to

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}), \tag{13.34a}$$

$$\text{subject to } g_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m, \tag{13.34b}$$

$$h_j(\mathbf{x}) = 0, \quad j = 1, \dots, \ell, \tag{13.34c}$$

where the functions f , g_i , $i = 1, \dots, m$, and h_j , $j = 1, \dots, \ell$, all are continuously differentiable on \mathbb{R}^n .

Suppose that $\bar{\mathbf{x}}$ is feasible in the problem (13.34). Prove the following statement by using linear programming duality: $\bar{\mathbf{x}}$ satisfies the KKT conditions if and only if the following LP problem has the optimal value zero:

$$\begin{aligned} & \text{minimize}_{\mathbf{p}} \quad \nabla f(\bar{\mathbf{x}})^T \mathbf{p}, \\ & \text{subject to } g_i(\bar{\mathbf{x}}) + \nabla g_i(\bar{\mathbf{x}})^T \mathbf{p} \leq 0, \quad i = 1, \dots, m, \\ & \quad h_j(\bar{\mathbf{x}}) + \nabla h_j(\bar{\mathbf{x}})^T \mathbf{p} = 0, \quad j = 1, \dots, \ell. \end{aligned}$$

Describe briefly how this LP problem could be used to devise an iterative method for the problem (13.34).

[Note: Algorithms in this class of methods are referred to as *Sequential Linear Programming* (SLP) methods.]

Exercise 13.8 (augmented Lagrangian) Consider the problem

$$\begin{aligned} f^* &:= \infimum_{\mathbf{x}} f(\mathbf{x}), \\ & \text{subject to } \mathbf{x} \in X, \end{aligned} \tag{P}$$

Constrained optimization

and

$$\begin{aligned} l^* &:= \infimum l(\mathbf{x}), \\ &\text{subject to } \mathbf{x} \in G. \end{aligned} \tag{R}$$

If $X \subseteq G$ and $l(\mathbf{x}) \leq f(\mathbf{x})$ for all $\mathbf{x} \in X$ we say that (R) is a *relaxation* of (P); cf. Section 6.1. Conversely, (P) is then a *restrification* of (R).

Consider the problem of the form

$$\begin{aligned} f^* &:= \infimum f(\mathbf{x}), \\ &\text{subject to } g_i(\mathbf{x}) = 0, \quad i = 1, \dots, m, \end{aligned}$$

where f and g_i , $i = 1, \dots, m$, are continuous functions on \mathbb{R}^n . Let μ_i , $i = 1, \dots, m$, be multipliers for the constraints and let $P : \mathbb{R}^m \rightarrow \mathbb{R}_+$ be a continuous *exterior penalty function*, that is, a function such that

$$P(\mathbf{y}) \begin{cases} = 0, & \text{if } \mathbf{y} = \mathbf{0}^m, \\ > 0, & \text{if } \mathbf{y} \neq \mathbf{0}^m. \end{cases}$$

Consider the penalized problem

$$\theta^* := \infimum_{\mathbf{x} \in \mathbb{R}^n} \theta(\mathbf{x}) := f(\mathbf{x}) + \sum_{i=1}^m \mu_i g_i(\mathbf{x}) + \nu P(\mathbf{g}(\mathbf{x})),$$

where $\mathbf{g}(\mathbf{x})$ is the m -vector of $g_i(\mathbf{x})$ and where $\nu > 0$. Show that this problem is a relaxation of the original one.

[Note: Algorithms based on the relaxation (R)—which linearly combines the Lagrangian and a penalty function—are known as *augmented Lagrangian methods*, and the function θ is known as the *augmented Lagrangian function*. They constitute an alternative to exact penalty methods, in that they also can be made convergent without having to let the penalty parameter tend to infinity, in this case because of the Lagrangian term; in augmented Lagrangian algorithms the multiplier $\boldsymbol{\mu}$ plays a much more active role than in SQP methods.]

Part VI

Appendix

Answers to the exercises



Chapter 1: Modelling and classification

Exercise 1.1 Variables:

x_j = number of units produced in process j , $j = 1, 2$;
 y = number of half hours hiring the model.

Optimization model:

$$\begin{aligned} &\text{maximize } f(x, y) := 50(3x_1 + 5x_2) - 3(x_1 + 2x_2) - 2(2x_1 + 3x_2) - 5000y, \\ &\text{subject to } \quad x_1 + 2x_2 \leq 20,000, \\ &\quad \quad \quad 2x_1 + 3x_2 \leq 35,000, \\ &\quad \quad \quad 3x_1 + 5x_2 \leq 1,000 + 200y \\ &\quad \quad \quad x_1 \geq 0, \\ &\quad \quad \quad x_2 \geq 0, \\ &\quad \quad \quad 0 \leq y \leq 6. \end{aligned}$$

Exercise 1.2 Variables:

x_j = number of trainees trained during month j , $j = 1, \dots, 5$;
 y_j = number of technicians available at the beginning of month j , $j = 1, \dots, 5$.

Optimization model:

Answers to the exercises

$$\begin{aligned}
 & \text{minimize } z = \sum_{j=1}^5 (15000y_j + 7500x_j) \\
 & \text{subject to } \begin{aligned}
 & 160y_1 - 50x_1 \geq 6000 \\
 & 160y_2 - 50x_2 \geq 7000 \\
 & 160y_3 - 50x_3 \geq 8000 \\
 & 160y_4 - 50x_4 \geq 9500 \\
 & 160y_5 - 50x_5 \geq 11,500 \\
 & 0.95y_1 + x_1 = y_2 \\
 & 0.95y_2 + x_2 = y_3 \\
 & 0.95y_3 + x_3 = y_4 \\
 & 0.95y_4 + x_4 = y_5 \\
 & y_1 = 50 \\
 & y_j, x_j \in \mathbb{Z}_+, \quad j = 1, \dots, 5.
 \end{aligned}
 \end{aligned}$$

Exercise 1.3 We declare the following indices:

- $i, i = 1, \dots, 3$: Work place,
- $k, k = 1, \dots, 2$: Connection point,

and variables

- (x_i, y_i) : Coordinates for work place i ;
- $t_{i,k}$: Indicator variable; its value is defined as 1 if work place i is connected to the connection point k , and as 0 otherwise;
- z : The longest distance to the window.

The problem to minimize the maximum distance to the window is that to

$$\text{minimize } z, \quad (\text{A.1})$$

subject to the work spaces being inside the rectangle:

$$\frac{d}{2} \leq x_i \leq l - \frac{d}{2}, \quad i = 1, \dots, 3, \quad (\text{A.2})$$

$$\frac{d}{2} \leq y_i \leq b - \frac{d}{2}, \quad i = 1, \dots, 3, \quad (\text{A.3})$$

that the work spaces do not overlap:

$$(x_i - x_j)^2 + (y_i - y_j)^2 \geq d^2, \quad i = 1, \dots, 3, \quad j = 1, \dots, 3, \quad i \neq j, \quad (\text{A.4})$$

that the cables are long enough:

$$t_{1,k} \left[(x_i - \frac{l}{2})^2 + (y_i - 0)^2 \right] \leq a_i^2, \quad i = 1, \dots, 3, \quad (\text{A.5})$$

$$t_{2,k} \left[(x_i - l)^2 + (y_i - \frac{b}{2})^2 \right] \leq a_i^2, \quad i = 1, \dots, 3, \quad (\text{A.6})$$

that each work space must be connected to a connection point:

$$t_{i,1} + t_{i,2} = 1, \quad i = 1, \dots, 3, \quad (\text{A.7})$$

$$t_{i,k} \in \{0, 1\}, \quad i = 1, \dots, 3, \quad k = 1, 2, \quad (\text{A.8})$$

and finally that the value of z is at least as high as the longest distance to the window:

$$b - y_i \geq z, \quad i = 1, \dots, 3. \quad (\text{A.9})$$

The problem hence is to minimize the objective function in (A.1) under the constraints (A.2)–(A.9).

Exercise 1.4 We declare the following indices:

- i : Warehouses ($i = 1, \dots, 10$),
- j : Department stores ($j = 1, \dots, 30$),

and variables:

- x_{ij} : portion (between 0 and 1) of the total demand at department store j which is served from warehouse i ,
- y_i : Indicator variable; its value is defined as 1 if warehouse i is built, and 0 otherwise.

We also need the following constants, describing the department stores that are within the specified maximum distance from a warehouse:

$$a_{ij} := \begin{cases} 1, & \text{if } d_{ij} \leq D, \\ 0, & \text{otherwise,} \end{cases} \quad i = 1, \dots, 10, \quad j = 1, \dots, 30.$$

(a) The problem becomes:

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^{10} c_i y_i, \\ & \text{subject to} && x_{ij} \leq a_{ij} y_i, \quad i = 1, \dots, 10, \quad j = 1, \dots, 30, \\ & && \sum_{j=1}^{30} e_j x_{ij} \leq k_i y_i, \quad i = 1, \dots, 10, \\ & && \sum_{i=1}^{10} x_{ij} = 1, \quad j = 1, \dots, 30, \\ & && x_{ij} \geq 0, \quad j = 1, \dots, 30, \\ & && y_i \in \{0, 1\}, \quad i = 1, \dots, 10. \end{aligned}$$

The first constraint makes sure that only warehouses that are built and which lie sufficiently close to a department store can supply any goods to it.

The second constraint describes the capacity of each warehouse, and the demand at the various department stores.

The third and fourth constraints describe that the total demand at a department store must be a non-negative (in fact, convex) combination of the contributions from the different warehouses.

(b) Additional constraints: $x_{ij} \in \{0, 1\}$ for all i and j .

Chapter 3: Convexity

Exercise 3.1 Use the definition of convexity (Definition 3.1).

Exercise 3.2 (a) S is a polyhedron. It is the parallelogram with the corners $a_1 + a_2, a_1 - a_2, -a_1 + a_2, -a_1 - a_2$, that is, $S = \text{conv} \{a_1 + a_2, a_1 - a_2, -a_1 + a_2, -a_1 - a_2\}$ which is a polytope and hence a polyhedron.

(b) S is a polyhedron.

(c) S is not a polyhedron. Note that although S is defined as an intersection of halfspaces it is not a polyhedron, since we need infinitely many halfspaces.

(d) $S = \{\mathbf{x} \in \mathbb{R}^n \mid -\mathbf{1}^n \leq \mathbf{x} \leq \mathbf{1}^n\}$, that is, a polyhedron.

(e) S is a polyhedron. By squaring both sides of the inequality, it follows that $-2(\mathbf{x}^0 - \mathbf{x}^1)^T \mathbf{x} \leq \|\mathbf{x}^1\|_2^2 - \|\mathbf{x}^0\|_2^2$, so S is in fact a halfspace.

(f) S is a polyhedron. Similarly as in e) above it follows that S is the intersection of the halfspaces

$$-2(\mathbf{x}^0 - \mathbf{x}^i)^T \mathbf{x} \leq \|\mathbf{x}^i\|_2^2 - \|\mathbf{x}^0\|_2^2, \quad i = 1, \dots, k.$$

Exercise 3.3 (a) \mathbf{x}^1 is not an extreme point.

(b) \mathbf{x}^2 is an extreme point. This follows by checking the rank of the equality subsystem and then using Theorem 3.17.

Exercise 3.4 Let

$$\mathbf{D} := \begin{pmatrix} \mathbf{A} \\ -\mathbf{A} \\ -\mathbf{I}^n \end{pmatrix}, \quad \mathbf{d} := \begin{pmatrix} \mathbf{b} \\ -\mathbf{b} \\ \mathbf{0}^n \end{pmatrix}.$$

Then P is defined by $\mathbf{D}\mathbf{x} \leq \mathbf{d}$. Further, P is nonempty, so let $\tilde{\mathbf{x}} \in P$. Now, if $\tilde{\mathbf{x}}$ is not an extreme point of P , then the rank of equality subsystem is lower than n . By using this it is possible to construct an $\mathbf{x}' \in P$ such that the rank of the equality subsystem of \mathbf{x}' is at least one larger than the rank of the equality subsystem of $\tilde{\mathbf{x}}$. If this argument is used repeatedly we end up with an extreme point of P .

Exercise 3.5 We have that

$$\begin{pmatrix} 1 \\ 1 \end{pmatrix} = 0.5 \begin{pmatrix} 0 \\ 1 \end{pmatrix} + 0.5 \begin{pmatrix} 1 \\ 0 \end{pmatrix} + 0.5 \begin{pmatrix} 1 \\ 1 \end{pmatrix},$$

and since $(0, 1)^T, (1, 0)^T \in Q$ and $(1, 1)^T \in C$ we are done.

Exercise 3.6 Assume that $a_1, a_2, a_3, b \in \mathbb{R}$ satisfy

$$a_1 x_1 + a_2 x_2 + a_3 x_3 \leq b, \quad \mathbf{x} \in A, \tag{A.10}$$

$$a_1 x_1 + a_2 x_2 + a_3 x_3 \geq b, \quad \mathbf{x} \in B. \tag{A.11}$$

From (A.10) follows that $a_2 = 0$ and that $a_3 \leq b$. Further, since $(1/n, n, 1)^T \in B$ for all $n > 0$, from (A.11) we have that $a_3 \geq b$. Hence, it holds that $a_3 = b$. Since $(0, 0, 0)^T, (1, n^2, n)^T \in B$ for all $n \geq 0$, inequality (A.11) shows that $b \leq 0$ and $a_3 \geq 0$. Hence $a_2 = a_3 = b = 0$, and it follows that $H := \{x \in \mathbb{R}^3 \mid x_1 = 0\}$ is the only hyperplane that separates A and B . Finally, $A \subseteq H$ and $(0, 0, 0)^T \in H \cap B$, so H meets both A and B .

Exercise 3.7 Let B be the intersection of all closed halfspaces in \mathbb{R}^n containing A . It follows easily that $A \subseteq B$. In order to show that $B \subseteq A$, show that $A^c \subseteq B^c$ by using the Separation Theorem 3.24.

Exercise 3.8 Assume that $P \neq \emptyset$. Then, by using Farkas' Lemma (Theorem 3.30), show that there exists a $p \neq 0^m$ such that $p \geq 0^m$ and $Bp \geq 0^m$. From this it follows that P is unbounded and hence not compact.

Exercise 3.9 —

Exercise 3.10 The function is strictly convex on \mathbb{R}^2 .

Exercise 3.11 (a) Not convex; (b)–(f) strictly convex.

Exercise 3.12 (a)–(f) Strictly convex.

Exercise 3.13 (a)

$$f(x, y) = \frac{1}{2}(x, y) \begin{bmatrix} 4 & -2 \\ -2 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} + (3, -1) \begin{bmatrix} x \\ y \end{bmatrix}.$$

(b) Yes. (c) Yes.

Exercise 3.14 (a) Non-convex; (b) convex; (c) non-convex; (d) convex; (e) convex.

Exercise 3.15 Yes.

Exercise 3.16 Yes.

Exercise 3.17 We will try to apply Definition 3.45. It is clear that the objective function can be written as the minimization of a (strictly) convex function. The constraints are analyzed thus: the first and third, taken together and applying also Example 3.37(c), describe a closed and convex set; the second and fourth constraint describes a (convex) polyhedron. By Proposition 3.3 we therefore are done. The answer is Yes.

Exercise 3.18 The first constraint is redundant; the feasible set hence is a nonempty polyhedron. Regarding the objective function, it is defined only for positive x_1 ; the objective function is strictly convex on \mathbb{R}_{++} , since its second

derivative there equals $1/x_1 > 0$ [cf. Theorem 3.41(b)]. We may extend the definition of $x_1 \ln x_1$ to a continuous (in fact convex) function, on the whole of \mathbb{R}_+ by defining $0 \ln 0 = 0$. With this classic extension, together with the constraint, we see that it is the problem of maximizing a convex function over a closed convex set. This is not a convex problem. The answer is No.

Chapter 4: An introduction to optimality conditions

Exercise 4.1 (a) —

(b) Argue by contradiction and utilize the convexity of f . The proof is similar in form to that of Theorem 4.3, and utilizes that moving slightly from \mathbf{x}^* still makes the removed constraint feasible, thus falsifying the initial claim that \mathbf{x}^* is optimal in the problem (4.32).

Exercise 4.2 Investigating the Hessian matrix yields that $a \in (-4, 2)$ and $b \in \mathbb{R}$ implies that the objective function is strictly convex (in fact, strongly convex, because it is quadratic).

[Note: It would be a mistake to here perform a classic transformation, namely to observe that the problem is symmetric in x_1 and x_2 and utilize this to eliminate one of the variables through the identification $x_1^* = x_2^*$. Suppose we do so. We then reduce the problem to that of minimizing the one-dimensional function $x \mapsto (4+a)x^2 - 2x + b$ over \mathbb{R} . The condition for this function to be strictly convex, and therefore have a unique solution (see the above remark on strong convexity), is that $a > -4$, which is a milder condition than the above. However, if the value of a is larger than 2 the *original* problem has no solution! Indeed, suppose we look at the direction $\mathbf{x} \in \mathbb{R}^2$ in which $x_1 = -x_2 = p$. Then, the function $f(\mathbf{x})$ behaves like $(2-a)p^2 - 2p + b$ which clearly tends to minus infinity whenever $|p|$ tends to infinity, whenever $a > 2$. It is important to notice that the transformation works *when the problem has a solution*; otherwise, it is not.]

Exercise 4.3 Let $\rho(\mathbf{x}) := \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}}$. Stationarity for ρ at \mathbf{x} means that

$$\frac{2}{\mathbf{x}^T \mathbf{x}} (\mathbf{A} \mathbf{x} - \rho(\mathbf{x}) \cdot \mathbf{x}) = \mathbf{0}^n.$$

If $\mathbf{x}_i \neq \mathbf{0}^n$ is an eigenvector of \mathbf{A} , corresponding to the eigenvalue λ_i , then $\rho(\mathbf{x}_i) = \lambda_i$ holds. From the above two equations follow that for $\mathbf{x} \neq \mathbf{0}^n$ to be stationary it is both necessary and sufficient that \mathbf{x} is an eigenvector.

The global minimum is therefore an arbitrary nonzero eigenvector, corresponding to the minimal eigenvalue λ_i of \mathbf{A} .

Exercise 4.4 (a) The proof is by contradiction, so suppose that $\bar{\mathbf{x}}$ is a local optimum, \mathbf{x}^* is a global optimum, and that $f(\bar{\mathbf{x}}) < f(\mathbf{x}^*)$ holds. We first note

that by the local optimality of $\bar{\mathbf{x}}$ and the affine nature of the constraints, it must hold that

$$\nabla f(\bar{\mathbf{x}})^T \mathbf{p} = \mathbf{0}^m, \quad \text{for all vectors } \mathbf{p} \text{ with } \mathbf{A}\mathbf{p} = \mathbf{0}^m.$$

We will especially look at the vector $\mathbf{p} := \mathbf{x}^* - \bar{\mathbf{x}}$.

Next, by assumption, $f(\bar{\mathbf{x}}) < f(\mathbf{x}^*)$, which implies that $(\bar{\mathbf{x}} - \mathbf{x}^*)^T \mathbf{Q}(\bar{\mathbf{x}} - \mathbf{x}^*) < 0$ holds. We utilize this strict inequality together with the above to last establish that, for every $\gamma > 0$,

$$f(\bar{\mathbf{x}} + \gamma(\bar{\mathbf{x}} - \mathbf{x}^*)) < f(\bar{\mathbf{x}}),$$

which contradicts the local optimality of $\bar{\mathbf{x}}$. We are done.

(b) —

Exercise 4.5 Utilize the variational inequality characterization of the projection operation.

Exercise 4.6 Utilize Proposition 4.23(b) for this special case of feasible set. We obtain the following necessary conditions for $\mathbf{x}^* \geq \mathbf{0}^n$ to be local minimum:

$$0 \leq x_j^* \perp \frac{\partial f(\mathbf{x}^*)}{\partial x_j} \geq 0, \quad j = 1, 2, \dots, n,$$

where (for real values a and b) $a \perp b$ means the condition that $a \cdot b = 0$ holds. In other words, if $x_j^* > 0$ then the partial derivative of f at \mathbf{x}^* with respect to x_j must be zero; conversely, if this partial derivative is non-zero then the value of x_j^* must be zero. (This is called complementarity.)

Exercise 4.7 By a logarithmic transformation, we may instead maximize the function $f(\mathbf{x}) = \sum_{j=1}^n a_j \log x_j$. The optimal solution is

$$x_j^* = \frac{a_j}{\sum_{i=1}^n a_i}, \quad j = 1, \dots, n.$$

(Check the optimality conditions for a problem defined over a simplex.)

We confirm that it is a unique optimal solution by checking that the objective function is strictly concave where it is defined.

Exercise 4.8 —

Exercise 4.9 —

Exercise 4.10 —

Chapter 5: Optimality conditions

Exercise 5.1 $(1, 2)^T$ is a KKT point for this problem with KKT multipliers $(1, 0)^T$. Since the problem is convex, this is also a globally optimal solution (cf. Theorem 5.45). Slater's CQ (and, in fact, LICQ as well) is verified.

Exercise 5.2 (a) The feasible set of the problem consists of countably many isolated points $x_k = -\pi/2 + 2\pi k$, $k = 1, 2, \dots$, each of which is thus a locally optimal solution. The globally optimal solution is $x^* = -\pi/2$. KKT conditions are not satisfied at the points of local minimum and therefore they are not necessary for optimality in this problem. (The reason is of course that CQs are not verified.)

(b) It is easy to verify that FJ conditions are satisfied (as they should be, cf. Theorems 5.8 and 5.15).

(c) The point $(x, y)^T = (0, 0)^T$ is a FJ point, but it has nothing to do with points of local minimum.

Exercise 5.3 KKT system:

$$\begin{aligned} Ax &\geq b, \\ \lambda &\geq 0, \\ c - A^T \lambda &= 0, \\ \lambda^T (Ax - b) &= 0. \end{aligned}$$

Combining the last two equations we obtain $c^T x = b^T \lambda$.

Exercise 5.4 (a) Clearly, the two problems are equivalent. On the other hand, $\nabla \{\sum_{i=1}^m [h_i(x)]^2\} = 2 \sum_{i=1}^m h_i(x) \nabla h_i(x) = 0$ at every feasible solution. Therefore, MFCQ is violated at every feasible point of the problem (5.22) (even though Slater's CQ, LICQ, or at least MFCQ might hold for the original problem).

(b) The objective function is non-differentiable. Therefore, we rewrite the problem as

$$\begin{aligned} &\text{minimize } z, \\ &\text{subject to } f_1(x) - z \leq 0, \\ &\quad f_2(x) - z \leq 0, \end{aligned}$$

The problem verifies MFCQ (e.g., the direction $(0^T, 1)^T \in \overset{\circ}{G}(x, z)$ for all feasible points $(x^T, z)^T$). Therefore, the KKT conditions are necessary for local optimality; these conditions are exactly what we need.

Exercise 5.5 The problem is convex and a CQ is fulfilled, so we need to find an arbitrary KKT point. The KKT system is as follows:

$$\begin{aligned} x + A^T \lambda &= 0, \\ Ax &= b. \end{aligned}$$

$\mathbf{A}\mathbf{x} + \mathbf{A}\mathbf{A}^T\boldsymbol{\lambda} = \mathbf{0}$ and $\mathbf{A}\mathbf{A}^T\boldsymbol{\lambda} = -\mathbf{b}$ yields $\mathbf{x} = \mathbf{A}^T(\mathbf{A}\mathbf{A}^T)^{-1}\mathbf{b}$.

Exercise 5.6 (a) —

(b) Show that the KKT multiplier λ is positive at every optimal solution. It means that $\sum_{j=1}^n x_j^2 = 1$ is satisfied at every optimal solution; use convexity to conclude that there may be only one optimal solution.

Exercise 5.7 (a) Locally and globally optimal solutions may be found using geometrical considerations; $(x, y) = (2, 0)$ gives us a local min, $(x, y) = (3/2, 3/2)$ is a globally optimal solution. KKT system incidentally has two [in the space (x, y)] solutions, but at every point there are infinitely many KKT multipliers. Therefore, in this particular problem KKT conditions are both necessary and sufficient for local optimality.

(b) The gradients of the constraints are linearly dependent at every feasible point; thus LICQ is violated.

The feasible set is a union of two convex sets $\mathcal{F}_1 := \{(x, y)^T \mid y = 0; x - y \geq 0\}$ and $\mathcal{F}_2 := \{(x, y)^T \mid y \geq 0; x - y = 0\}$. Thus we can solve two convex optimization problems to minimize f over \mathcal{F}_1 , and to minimize f over \mathcal{F}_2 ; then simply choose the best solution.

(c) The feasible set may be split into 2^n convex parts \mathcal{F}_I , $I \subseteq \{1, \dots, n\}$, where

$$\begin{aligned} \mathbf{a}_i^T \mathbf{x} &= b_i, \text{ and } x_i \geq 0, & i \in I, \\ \mathbf{a}_i^T \mathbf{x} &\geq b_i, \text{ and } x_i = 0, & i \notin I. \end{aligned}$$

Thus we (in principle) have reduced the original non-convex problem that violates LICQ to 2^n convex problems.

Exercise 5.8 Use the KKT conditions (convex problem + Slater's CQ). $c \leq -1$.

Exercise 5.9 Slater's CQ implies that the KKT conditions are necessary for optimality. Prove that $x_j^* > 0$; then $x_j^* = Dc_j / \sum_{k=1}^n c_k$, $j = 1, \dots, n$.

Chapter 6: Lagrangian duality

Exercise 6.1 (a) For $\mu > 0$, the Lagrangian $L(\cdot, \mu)$ is strictly convex; $x_1(\mu) = 1/\sqrt{\mu}$ and $x_2(\mu) = 4/\sqrt{\mu}$ uniquely. For $\mu = 0$, $q(\mu) = -\infty$.

(b) $\mu^* = 9/16$.

(c) The dual problem is to maximize $_{\mu \geq 0} q(\mu) = 6\sqrt{\mu} - 4\mu$.

(d) $\mathbf{x}^* = (4/3, 8/3)^T$.

(e) $f^* = q^* = 9/4$.

Exercise 6.2 —

Answers to the exercises

Exercise 6.3 $\mathbf{x}^* = (4/2, 2/3)^T$; $\mu^* = 8/3$; $f^* = q^* = 22/9$.

Exercise 6.4 From $L(\mathbf{x}, \boldsymbol{\lambda}) := \frac{1}{2}\|\mathbf{x} - \mathbf{y}\|^2 + \boldsymbol{\lambda}^T \mathbf{A}\mathbf{x}$ we get that $\nabla_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{0}^n$ when $\mathbf{x} = \mathbf{y} - \mathbf{A}^T \boldsymbol{\lambda}$ [uniquely by the strict convexity and coercivity of $L(\cdot, \boldsymbol{\lambda})$]. Inserted into the Lagrangian we obtain $q(\boldsymbol{\lambda}) = \mathbf{y}^T \mathbf{A}^T \boldsymbol{\lambda} - \frac{1}{2}\|\mathbf{A}^T \boldsymbol{\lambda}\|^2$. From $\nabla q(\boldsymbol{\lambda}) = \mathbf{0}^m$ we obtain that $\boldsymbol{\lambda}^* = (\mathbf{A}\mathbf{A}^T)^{-1} \mathbf{A}\mathbf{y}$ (uniquely), which yields the formula sought.

Exercise 6.5 (a) The Slater CQ is verified since the problem is convex (even linear), and there is a strictly feasible point [e.g., $(x, y)^T = (3, 1)^T$].

Introducing Lagrange multipliers μ_1 and μ_2 we calculate the Lagrangian dual function q :

$$\begin{aligned} q(\mu_1, \mu_2) &= \min_{(\mu_1, \mu_2) \in \mathbb{R}_+^2} \{x - 1/2y + \mu_1(-x + y + 1) + \mu_2(-2x + y + 2)\} \\ &= \mu_1 + 2\mu_2 + \min_{x \geq 0} (1 - \mu_1 - 2\mu_2)x + \min_{y \geq 0} (-1/2 + \mu_1 + \mu_2)y \\ &= \begin{cases} \mu_1 + 2\mu_2, & \text{if } \mu_1 + 2\mu_2 \leq 1 \text{ and } \mu_1 + \mu_2 \geq 1/2, \\ -\infty, & \text{otherwise.} \end{cases} \end{aligned}$$

Therefore, the set of optimal Lagrange multipliers is $\{(\mu_1, \mu_2) \mid \mu_1 \geq 0; \mu_2 \geq 0; \mu_1 + 2\mu_2 = 1; \mu_1 + \mu_2 \geq 1/2\}$, which is clearly convex and bounded (e.g., you may illustrate this graphically) as it should be in the presence of Slater's CQ.

(b) At $(\mu_1, \mu_2)^T = (1/4, 1/3)^T$ the set of optimal solutions to the Lagrangian relaxed problem is the singleton $\{(0, 0)^T\}$. Hence, the Lagrangian function is differentiable at this point and its gradient equals the value of the vector of constraint functions evaluated at the optimal solution to the relaxed problem, i.e., $(-0 + 0 + 1, -2 \cdot 0 + 0 + 2)^T = (1, 2)^T$. Alternatively, differentiate q at a given point to obtain the result.

At $(\mu_1, \mu_2)^T = (1, 0)^T$ the set of optimal solutions to the Lagrangian relaxed problem is not a singleton: it equals $\{(x, 0)^T \mid x \geq 0\}$. Hence, the dual function is not differentiable, and the set of subgradients is obtained by evaluating the constraint functions at the optimal solutions to the relaxed problem, i.e., $\partial q(1, 0) = \{(-x + 1, -2x + 2)^T \mid x \geq 0\}$.

Exercise 6.6 Introduce the multipliers μ_j , $j = 1, \dots, n$, and λ_i , $i = 1, \dots, m$. We obtain that the minimum of the Lagrangian is obtained at $x_{ij}(\mu_j, \lambda_i) = e^{-(1+\mu_j+\lambda_i)}$. Inserted into the Lagrangian yields the dual objective function

$$q(\boldsymbol{\mu}, \boldsymbol{\lambda}) = - \sum_{j=1}^n \sum_{i=1}^m e^{-(1+\mu_j+\lambda_i)} - \sum_{j=1}^n b_j \mu_j - \sum_{i=1}^m a_i \lambda_i,$$

which is to be maximized over $(\boldsymbol{\mu}, \boldsymbol{\lambda}) \in \mathbb{R}^n \times \mathbb{R}^m$.

Exercise 6.7

$$\begin{aligned}\lambda = 1 &\implies x_1 = 1, \quad x_2 = 2, && \text{infeasible,} && q(1) = 6; \\ \lambda = 2 &\implies x_1 = 1, \quad x_2 = 5/2, && \text{infeasible,} && q(2) = 43/4; \\ \lambda = 3 &\implies x_1 = 3, \quad x_2 = 3, && \text{feasible,} && q(3) = 9.\end{aligned}$$

Further, $f(3, 3) = 21$, so $43/4 \leq f^* \leq 21$.

Exercise 6.8 (a) The value of the Lagrangian dual function is given by $q(\boldsymbol{\mu}) := \text{minimum}_{\mathbf{x}^p \in \{1, \dots, P\}} \{f(\mathbf{x}^p) + \boldsymbol{\mu}^T \mathbf{g}(\mathbf{x}^p)\}$, which is the point-wise minimum of P affine functions. Therefore, q is piece-wise linear, with no more than P pieces; the number is less if for some value(s) of $\boldsymbol{\mu}$ more than one element \mathbf{x}^p attains the minimum value of the Lagrangian.

(b) —

(c) $q(\boldsymbol{\mu}) := \text{minimum}_{i \in \{1, \dots, I\}} \{f(\mathbf{x}^i) + \boldsymbol{\mu}^T \mathbf{g}(\mathbf{x}^i)\}$, where \mathbf{x}^i , $i \in \{1, \dots, I\}$, are the extreme points of the polytope X . The number of pieces of the dual function is bounded by the number I of extreme points of X .

Exercise 6.9 The dual problem is to maximize $q(\mu)$ over $\mu \in \mathbb{R}_+$, where

$$q(\mu) := 5\mu + \text{minimum}_{x_1 \in \{0, 1, \dots, 4\}} (2 - \mu)x_1 + \text{minimum}_{x_2 \in \{0, 1, \dots, 4\}} (1 - \mu)x_2.$$

$$\begin{aligned}\mu = 0 &\text{ yields } \mathbf{x}(\mu) = (0, 0)^T; \quad q(0) = 0; \\ \mu = 1 &\text{ yields } \mathbf{x}(\mu) = (0, 0)^T \text{ (for example); } \quad q(1) = 5; \\ \mu = 2 &\text{ yields } \mathbf{x}(\mu) = (0, 4)^T \text{ (for example); } \quad q(2) = 6; \\ \mu = 3 &\text{ yields } \mathbf{x}(\mu) = (4, 4)^T; \quad q(3) = 3.\end{aligned}$$

Since $\mathbf{x}(\mu)$ is feasible for $\mu = 3$, we also gain access to a primal feasible solution; $f((4, 4)^T) = 12$.

We conclude that $f^* \in [6, 12]$.

Exercise 6.10 (a) The feasible set of (S) includes that of (P) . The result then follows from the Relaxation Theorem 6.1.

(b) —

(c) The Relaxation Theorem 6.1 applies; the objective function in the Lagrangian minorizes f and the former problem's feasible set is larger.

Chapter 8: Linear programming models

Exercise 8.1 (a) Introduce the new variables $\mathbf{y} \in \mathbb{R}^m$. Then the problem is equivalent to the linear program

$$\begin{aligned}\text{minimize} \quad & \sum_{i=1}^m y_i, \\ \text{subject to} \quad & -\mathbf{y} \leq \mathbf{A}\mathbf{x} - \mathbf{b} \leq \mathbf{y}, \\ & -\mathbf{1}^n \leq \mathbf{x} \leq \mathbf{1}^n.\end{aligned}$$

Answers to the exercises

(b) Introduce the new variables $\mathbf{y} \in \mathbb{R}^m$ and $t \in \mathbb{R}$. Then the problem is equivalent to the linear program

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^m y_i + t, \\ & \text{subject to} && -\mathbf{y} \leq \mathbf{A}\mathbf{x} - \mathbf{b} \leq \mathbf{y}, \\ & && -t\mathbf{1}^n \leq \mathbf{x} \leq t\mathbf{1}^n. \end{aligned}$$

Exercise 8.2 (a) Let

$$\mathbf{B} := \begin{pmatrix} -(\mathbf{v}^1)^\top & 1 \\ \vdots & \vdots \\ -(\mathbf{v}^k)^\top & 1 \\ (\mathbf{w}^1)^\top & -1 \\ \vdots & \vdots \\ (\mathbf{w}^l)^\top & -1 \end{pmatrix}, \quad \mathbf{x} := \begin{pmatrix} a \\ b \end{pmatrix}.$$

Then from the rank assumption it follows that $\text{rank } \mathbf{B} = n + 1$, which means that $\mathbf{x} \neq \mathbf{0}^{n+1}$ implies that $\mathbf{B}\mathbf{x} \neq \mathbf{0}^{k+l}$. Hence the problem can be solved by solving the linear program

$$\begin{aligned} & \text{minimize} && (\mathbf{0}^{n+1})^\top \mathbf{x}, \\ & \text{subject to} && \mathbf{B}\mathbf{x} \geq \mathbf{0}^{k+l}, \\ & && (\mathbf{1}^{k+l})^\top \mathbf{B}\mathbf{x} = 1. \end{aligned}$$

(b) Let $\alpha = R^2 - \|\mathbf{x}^c\|_2^2$. Then the problem can be solved by solving the linear program

$$\begin{aligned} & \text{minimize} && (\mathbf{0}^n)^\top \mathbf{x}^c + 0\alpha, \\ & \text{subject to} && \|\mathbf{v}^i\|_2^2 - 2(\mathbf{v}^i)^\top \mathbf{x}^c \leq \alpha, \quad i = 1, \dots, k, \\ & && \|\mathbf{w}^i\|_2^2 - 2(\mathbf{w}^i)^\top \mathbf{x}^c \geq \alpha, \quad i = 1, \dots, l, \end{aligned}$$

and compute R as $R = \sqrt{\alpha + \|\mathbf{x}^c\|_2^2}$ (from the first set of inequalities in the LP above it follows that $\alpha + \|\mathbf{x}^c\|_2^2 \geq 0$ so this is well defined).

Exercise 8.3 Since P is bounded there exists no $\mathbf{y} \neq \mathbf{0}^n$ such that $\mathbf{A}\mathbf{y} \leq \mathbf{0}^m$. Hence there exists no feasible solution to the system

$$\begin{aligned} & \mathbf{A}\mathbf{y} \leq \mathbf{0}^m, \\ & \mathbf{d}^\top \mathbf{y} = 1, \end{aligned}$$

which implies that $z > 0$ in every feasible solution to (8.11).

Further, let (\mathbf{y}^*, z^*) be a feasible solution to (8.11). Then $z^* > 0$ and $\mathbf{x}^* = \mathbf{y}^*/z^*$ is feasible to (8.10), and $f(\mathbf{x}^*) = g(\mathbf{y}^*, z^*)$. Conversely, let \mathbf{x}^* be a feasible solution to (8.10). Then by the hypothesis $\mathbf{d}^T \mathbf{x}^* + \beta > 0$. Let $z^* = 1/(\mathbf{d}^T \mathbf{x}^* + \beta)$ and $\mathbf{y}^* = z^* \mathbf{x}^*$. Then (\mathbf{y}^*, z^*) is a feasible solution to (8.11) and $g(\mathbf{y}^*, z^*) = f(\mathbf{x}^*)$. These facts imply the assertion.

Exercise 8.4 The problem can be transformed into the standard form:

$$\begin{aligned} \text{minimize } z' &= x_1' - 5x_2^+ + 5x_2^- - 7x_3^+ + 7x_3^-, \\ \text{subject to } & 5x_1' - 2x_2^+ + 2x_2^- + 6x_3^+ - 6x_3^- - s_1 = 15, \\ & 3x_1' + 4x_2^+ - 4x_2^- - 9x_3^+ + 9x_3^- = 9, \\ & 7x_1' + 3x_2^+ - 3x_2^- + 5x_3^+ - 5x_3^- + s_2 = 23, \\ & x_1', x_2^+, x_2^-, x_3^+, x_3^-, s_1, s_2 \geq 0, \end{aligned}$$

where $x_1' = x_1 + 2$, $x_2 = x_2^+ - x_2^-$, $x_3 = x_3^+ - x_3^-$, and $z' = z - 2$.

Exercise 8.5 (a) The first equality constraint gives that

$$x_3 = \frac{1}{6}(11 - 2x_1 - 4x_2).$$

Now, by substituting x_3 with this expression in the objective function and the second equality constraint the problem is in standard form and x_3 is eliminated.

(b) If $x_3 \geq 0$, then we must add the constraint $(11 - 2x_1 - 4x_2)/6 \geq 0$ to the problem. But this is an inequality, so in order to transform the problem into standard form we must add a slack variable.

Exercise 8.6 Assume that the column in the constraint matrix corresponding to the variable x_j^+ is \mathbf{a}_j . Then the column in the constraint matrix corresponding to the variable x_j^- is $-\mathbf{a}_j$. The statement follows from the definition of a BFS, since \mathbf{a}_j and $-\mathbf{a}_j$ are linearly dependent.

Exercise 8.7 Let P be the set of feasible solutions to (8.12) and Q be the set of feasible solutions to (8.13). Obviously $P \subseteq Q$. In order to show that $Q \subseteq P$ assume that there exists an $\mathbf{x} \in Q$ such that $\mathbf{x} \notin P$ and derive a contradiction.

Chapter 9: The simplex method

Exercise 9.1 The phase I problem becomes

$$\begin{aligned} \text{minimize } w &= a_1 + a_2, \\ \text{subject to } & -3x_1 - 2x_2 + x_3 - s_1 + a_1 = 3, \\ & x_1 + x_2 - 2x_3 - s_2 + a_2 = 1, \\ & x_1, x_2, x_3, s_1, s_2, a_1, a_2 \geq 0. \end{aligned}$$

Answers to the exercises

From the equality constraints follow that $a_1 + a_2 \geq 4$ for all $x_1, x_2, x_3, s_1, s_2 \geq 0$. In particular, it follows that $w \geq 4$ for all feasible solutions to the phase I problem, which means that the original problem is infeasible.

Exercise 9.2 (a) The standard form is given by

$$\begin{aligned} \text{minimize} \quad & 3x_1 + 2x_2 + x_3, \\ \text{subject to} \quad & 2x_1 + x_3 - s_1 = 3, \\ & 2x_1 + 2x_2 + x_3 = 5, \\ & x_1, x_2, x_3, s_1 \geq 0. \end{aligned}$$

By solving the phase I problem with the simplex algorithm we get the feasible basis $\mathbf{x}_B = (x_1, x_2)^T$. Then by solving the phase II problem with the simplex algorithm we get the optimal solution $\mathbf{x}^* = (x_1, x_2, x_3)^T = (0, 1, 3)^T$.

(b) No, the set of all optimal solution is given by the set

$$\{\mathbf{x} \in \mathbb{R}^3 \mid \lambda(0, 1, 3)^T + (1 - \lambda)(0, 0, 5)^T; \quad \lambda \in [0, 1]\}.$$

Exercise 9.3 The reduced cost for all the variables except for x_j must be greater than or equal to 0. Hence it follows that the current basis is optimal to the problem that arises if x_j is fixed to zero. The assertion then follows from the fact that the current basis is non-degenerate.

Exercise 9.4 —

Chapter 10: LP duality and sensitivity analysis

Exercise 10.1 The linear programming dual is given by

$$\begin{aligned} \text{minimize} \quad & 11y_1 + 23y_2 + 12y_3, \\ \text{subject to} \quad & 4y_1 + 3y_2 + 7y_3 \geq 6, \\ & 3y_1 + 2y_2 + 4y_3 \geq -3, \\ & -8y_1 + 7y_2 + 3y_3 \leq -2, \\ & 7y_1 + 6y_2 + 2y_3 = 5, \\ & y_2 \leq 0, \\ & y_3 \geq 0. \end{aligned}$$

Exercise 10.2 (a) The linear programming dual is given by

$$\begin{aligned} \text{maximize} \quad & \mathbf{b}^T \mathbf{y}^1 + \mathbf{l}^T \mathbf{y}^2 + \mathbf{u}^T \mathbf{y}^3, \\ \text{subject to} \quad & \mathbf{A}^T \mathbf{y}^1 + \mathbf{I}^n \mathbf{y}^2 + \mathbf{I}^n \mathbf{y}^3 = \mathbf{c}, \\ & \mathbf{y}^2 \geq \mathbf{0}^n, \\ & \mathbf{y}^3 \leq \mathbf{0}^n. \end{aligned}$$

(b) A feasible solution to the linear programming dual is given by

$$\begin{aligned} \mathbf{y}^1 &= \mathbf{0}^m, \\ \mathbf{y}^2 &= (\max\{0, c_1\}, \dots, \max\{0, c_n\})^T, \\ \mathbf{y}^3 &= (\min\{0, c_1\}, \dots, \min\{0, c_n\})^T. \end{aligned}$$

Exercise 10.3 Use the Weak and Strong Duality Theorems.

Exercise 10.4 The LP dual is infeasible. Hence, from the Weak and Strong Duality Theorems it follows that the primal problem is either infeasible or unbounded.

Exercise 10.5 By using the Strong Duality Theorem we get the following polyhedron:

$$\begin{aligned} \mathbf{A}\mathbf{x} &\geq \mathbf{b}, \\ \mathbf{A}^T \mathbf{y} &\leq \mathbf{c}, \\ \mathbf{c}^T \mathbf{x} &= \mathbf{b}^T \mathbf{y}, \\ \mathbf{x} &\geq \mathbf{0}^n, \\ \mathbf{y} &\leq \mathbf{0}^m. \end{aligned}$$

Exercise 10.6 From the Strong Duality Theorem it follows that $\mathbf{c}^T \mathbf{x}^* = \mathbf{b}^T \mathbf{y}^*$. Use this to establish the statement.

Exercise 10.7 The dual problem only contains two variables and hence can be solved graphically. We get the optimal solution $\mathbf{y}^* = (-2, 0)^T$. The complementary slackness conditions then implies that $x_1 = x_2 = x_3 = x_5 = 0$. Hence, let $\mathbf{x}_B = (x_4, x_6)^T$. The optimal solution is $\mathbf{x}^* = (x_1, x_2, x_3, x_4, x_5, x_6)^T = (0, 0, 0, 3, 0, 1)^T$.

Exercise 10.8 From the complementary slackness conditions and the fact

Answers to the exercises

that $c_1/a_1 \geq \dots \geq c_n/a_n$ it follows that

$$\begin{aligned}u &= \frac{c_r}{a_r}, \\y_j &= c_j - \frac{c_r}{a_r} a_j, \quad j = 1, \dots, r-1, \\y_j &= 0, \quad j = r, \dots, n,\end{aligned}$$

is a dual feasible solution which together with the given primal solution fulfil the LP primal–dual optimality conditions.

Exercise 10.9 —

Exercise 10.10 —

Exercise 10.11 —

Exercise 10.12 —

Exercise 10.13 The basis $\mathbf{x}_B := (x_1, x_2)^T$ is optimal as long as $c_3 \leq 5$ and $c_4 \geq 8$.

Exercise 10.14 b) The basis $\mathbf{x}_B := (x_1, x_3)^T$ is optimal for all $\delta \geq -6.5$.
c) The basis $\mathbf{x}_B := (x_1, x_3)$ is not primal feasible for $\delta = -7$, but it is dual feasible, so by using the dual simplex method it follows that $\mathbf{x}_B := (x_1, x_5)^T$ is an optimal basis.

Exercise 10.15 —

Exercise 10.16 The problem is that $\min \leq \max$ is not true if the feasible set is empty. For example, suppose we take $\max 2x_1 + x_2 : x_1 + x_2 \leq 1; \mathbf{x} \geq \mathbf{0}^2$. The dual is $\min y : y \geq 2; y \geq 1; y \geq 0$. A few lines down, where we require $\mathbf{x} \leq \mathbf{0}^2$, the primal is infeasible; that's where the sequence fails.

The conclusion is however true for $\mathbf{b} = \mathbf{0}^m$ and $\mathbf{c} = \mathbf{0}^n$.

Chapter 11: Unconstrained optimization

Exercise 11.1 —

Exercise 11.2 The directional derivative is 13; the answer is No.

Exercise 11.3 (a) The search direction is not a descent direction, for example because the Hessian matrix is indefinite or negative definite.

(b) The linear system is unsolvable, for example because the Hessian matrix is indefinite. [Note: Even for indefinite Hessians, the search direction might exist for *some* right-hand sides.]

(c) Use the Levenberg–Marquardt modification.

Exercise 11.4 Let $y_1 := x_1 - 2$ and $y_2 := \sqrt{5}(x_2 + 6)$. We then get $f(\mathbf{x}) = g(\mathbf{y}) = y_1^2 + y_2^2$. At every $\mathbf{y} \in \mathbb{R}^2$ the negative gradient points towards the optimum!

Exercise 11.5 (a) $\mathbf{x}_1 = (1/2, 1)^T$.
(b) The Hessian matrix is

$$\nabla^2 f(\mathbf{x}_1) = \begin{pmatrix} 10 & -4 \\ -4 & 2 \end{pmatrix}.$$

The answer is Yes.

(c) The answer is Yes.

Exercise 11.6 (a) $\mathbf{x}_1 = (2, 1/2)^T$.
(b) The answer is No. The gradient is zero.
(c) The answer is Yes.

Exercise 11.7 (a) —
(b) $\mu \in (0, 0.6)$.

Exercise 11.8 (a) $\mathbf{f}(\mathbf{x}_0) = (-1, -2)^T \implies \|\mathbf{f}(\mathbf{x}_0)\| = \sqrt{5}$; $\mathbf{x}_1 = (4/3, 2/3)^T \implies \|\mathbf{f}(\mathbf{x}_1)\| = 16/27$.

(b) If \mathbf{f} is the gradient of a C^2 function $f : \mathbb{R}^n \mapsto \mathbb{R}$ we obtain that $\nabla \mathbf{f} = \nabla^2 f$, that is, Newton's method for unconstrained optimization is obtained.

Exercise 11.9 (a) $\mathbf{x}^* = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}$.

(b) The objective function is convex, since the Hessian is $\mathbf{A}^T \mathbf{A}$ (which is always positive semidefinite; check!). Therefore, the normal solution in (a) is globally optimal.

Exercise 11.10 The reason is that we do not wish to allow for a bad direction \mathbf{p}_k to be compensated for by simply allowing it to be longer.

Exercise 11.11 —

Exercise 11.12 (a) We have that

$$\begin{aligned} \nabla f(\mathbf{y}) + \gamma(\mathbf{y} - \mathbf{x}_k) = \mathbf{0}^n &\iff \mathbf{Q}\mathbf{y} + \mathbf{q} + \gamma(\mathbf{y} - \mathbf{x}_k) = \mathbf{0}^n \iff \\ (\mathbf{Q} + \gamma\mathbf{I}^n)\mathbf{y} &= \gamma\mathbf{x}_k - \mathbf{q} \iff (\mathbf{Q} + \gamma\mathbf{I}^n)(\mathbf{y} - \mathbf{x}_k). \end{aligned}$$

Further,

$$(\mathbf{Q} + \gamma\mathbf{I}^n)(\mathbf{y} - \mathbf{x}_k) = \gamma\mathbf{x}_k - \mathbf{q} - (\mathbf{Q} + \gamma\mathbf{I}^n)\mathbf{x}_k = -(\mathbf{Q}\mathbf{x}_k + \mathbf{q}).$$

(b) If $\{\mathbf{x}_k\}$ converges to \mathbf{x}^∞ then $\{\mathbf{p}_k\} = \{\mathbf{x}_{k+1} - \mathbf{x}_k\}$ must converge to zero. From the updating formula we obtain that $\mathbf{p}_k = (\mathbf{Q} + \gamma\mathbf{I}^n)^{-1} \nabla f(\mathbf{x}_k)$

Answers to the exercises

for every k . The sequence $\{\nabla f(\mathbf{x}_k)\}$ converges to $\nabla f(\mathbf{x}^\infty)$, since $f \in C^1$. If $\nabla f(\mathbf{x}^\infty) \neq \mathbf{0}^n$ then $\{\mathbf{p}_k\}$ would converge to $(\mathbf{Q} + \gamma \mathbf{I}^n)^{-1} \nabla f(\mathbf{x}^\infty) \neq \mathbf{0}^n$, since $(\mathbf{Q} + \gamma \mathbf{I}^n)^{-1}$ is positive definite when $\mathbf{Q} + \gamma \mathbf{I}^n$ is. This leads to a contradiction. Hence, $\nabla f(\mathbf{x}^\infty) = \mathbf{0}^n$. Since f is convex \mathbf{x}^∞ is a global minimum of f over \mathbb{R}^n .

Exercise 11.13 Case I: $\{\nabla f(\mathbf{x}_k)\} \rightarrow \mathbf{0}^n$; $\{\mathbf{x}_k\}$ and $\{f(\mathbf{x}_k)\}$ diverge.

Example: $f(x) = -\log x$; $\{x_k\} \rightarrow \infty$; $\{f(x_k)\} \rightarrow -\infty$; $\{f'(x_k)\} \rightarrow 0$.

Case II: $\{\nabla f(\mathbf{x}_k)\} \rightarrow \mathbf{0}^n$; $\{\mathbf{x}_k\}$ diverges; $\{f(\mathbf{x}_k)\}$ converges.

Example: $f(x) = 1/x$; $\{x_k\} \rightarrow \infty$; $\{f(x_k)\} \rightarrow 0$; $\{f'(x_k)\} \rightarrow 0$.

Case III: $\{\nabla f(\mathbf{x}_k)\} \rightarrow \mathbf{0}^n$; $\{\mathbf{x}_k\}$ is bounded; $\{f(\mathbf{x}_k)\}$ is bounded.

Example: $f(x) = \frac{1}{3}x^3 - x$; $x_k = \begin{cases} 1 + 1/k, & k \text{ even} \\ -1 - 1/k, & k \text{ odd} \end{cases}$

$\{x_k\}$ has two limit points: ± 1 ; $\{f(x_k)\}$ has two limit points: $\pm 2/3$.

Case IV: $\{\nabla f(\mathbf{x}_k)\} \rightarrow \mathbf{0}^n$; $\{\mathbf{x}_k\}$ is bounded; $\{f(\mathbf{x}_k)\}$ converges.

Example: $f(x) = x^2 - 1$; x_k as above; $\{f(x_k)\} \rightarrow 0$.

Case V: $\{\nabla f(\mathbf{x}_k)\} \rightarrow \mathbf{0}^n$; $\{\mathbf{x}_k\}$ and $\{f(\mathbf{x}_k)\}$ converge.

Example: f as in Case IV; $x_k = 1 + 1/k$.

Exercise 11.14 —

Exercise 11.15 —

Exercise 11.16 —

Exercise 11.17 —

Chapter 12: Optimization over convex sets

Exercise 12.1 If the LP problem (12.2) has an unbounded solution, we must replace the search direction towards an extreme point with the extreme direction identified in the last iteration of the simplex method; that is, we choose \mathbf{p}_k to be a direction of S in which the objective value in (12.2) tends to $-\infty$. The feasible set in the line search is \mathbb{R}_+ .

Under a compactness assumption on the intersection of S with the level set of f at \mathbf{x}_0 the convergence properties in Theorem 12.1 can be reached also for this more general problem and algorithm, otherwise not necessarily.

Exercise 12.2 —

Exercise 12.3 (a) —

(b) $\mathbf{x}_1 = (12/5, 4/5)^T$; UBD = $f(\mathbf{x}_1) = 8$. The LP problem defined at \mathbf{x}_0 gives LBD = 0. Hence, $f^* \in [0, 8]$.

Exercise 12.4 (a) —
 (b) $\mathbf{x}^* = (4, 2)^T$; $f^* = 80$.

Exercise 12.5 $\mathbf{x}_0 = (1, 1)^T$; $f(\mathbf{x}_0) = 5/8$; $\mathbf{y}_0 = (0, 0)^T$; $z(\mathbf{y}_0) = -7/8$; $\mathbf{x}_1 = (1/4, 1/4)^T$; $f(\mathbf{x}_1) = 1/16$; $\mathbf{y}_1 = (1, 0)^T$; $z(\mathbf{y}_1) = -3/16$; $\mathbf{x}_2 = (13/20, 5/20)^T$; $f(\mathbf{x}_2) = 1/80$.
 $f^* \in [-3/16, 1/80]$.

Exercise 12.6 —

Exercise 12.7 —

Exercise 12.8 —

Exercise 12.9 The answer is yes. The extreme points visited and stored are $(1, 1)^T$ (if we start at the same place as in Exercise 12.5), $(0, 0)^T$, and $(1, 0)^T$, which are the same as in the Frank-Wolfe algorithm. Using the simplicial decomposition method, the optimal solution $\mathbf{x}^* = (1/2, 0)^T$ is found in the convex hull of these points.

Chapter 13: Constrained optimization

Exercise 13.1 —

Exercise 13.2 —

Exercise 13.3 —

Exercise 13.4 —

Exercise 13.5 For a given parameter value $\nu > 0$ the unconstrained problem to

$$\underset{\mathbf{x} \in \mathbb{R}^2}{\text{minimize}} \quad f(\mathbf{x}) - \nu \cdot \log(x_1 + 2x_2 - 10)$$

uniquely solvable:

$$x_1 - \frac{\nu}{x_1 + 2x_2 - 10} = 0; \quad 2x_2 - \frac{2\nu}{x_1 + 2x_2 - 10} = 0$$

yields that $x_1 = x_2$ must hold; the resulting quadratic equation $3x_1^2 - 10x_1 - \nu = 0$ has two roots, of which $x_1(\nu) = 5/3 + \sqrt{25/9 + \nu/3}$ is strictly feasible. As $\nu \rightarrow 0$, $x_1(\nu) = x_2(\nu)$ tends to $10/3$.

One then shows that $\mathbf{x}^* = (\frac{10}{3}, \frac{10}{3})^T$ is a KKT point. The constraint is binding, and $\mu^* = 10/3 \geq 0$. Since the problem is convex, \mathbf{x}^* is optimal.

Exercise 13.6 —

Exercise 13.7 Let us first rewrite the LP problem into the following equivalent form, and note that $h_j(\bar{\mathbf{x}}) = 0$ for all j , since $\bar{\mathbf{x}}$ is feasible:

$$\begin{aligned} & \underset{\mathbf{p}}{\text{minimize}} && \nabla f(\bar{\mathbf{x}})^T \mathbf{p}, \\ & \text{subject to} && -\nabla g_i(\bar{\mathbf{x}})^T \mathbf{p} \geq g_i(\bar{\mathbf{x}}), \quad i = 1, \dots, m, \\ & && -\nabla h_j(\bar{\mathbf{x}})^T \mathbf{p} = 0, \quad j = 1, \dots, \ell. \end{aligned}$$

Letting $\boldsymbol{\mu} \geq \mathbf{0}^m$ and $\boldsymbol{\lambda} \in \mathbb{R}^\ell$ be the dual variable vector for the inequality and equality constraints, respectively, we obtain the following dual program:

$$\begin{aligned} & \underset{(\boldsymbol{\mu}, \boldsymbol{\lambda})}{\text{maximize}} && \sum_{i=1}^m \mu_i g_i(\bar{\mathbf{x}}), \\ & \text{subject to} && -\sum_{i=1}^m \mu_i \nabla g_i(\bar{\mathbf{x}}) - \sum_{j=1}^{\ell} \lambda_j \nabla h_j(\bar{\mathbf{x}}) = \nabla f(\bar{\mathbf{x}}), \\ & && \mu_i \geq 0, \quad i = 1, \dots, m. \end{aligned}$$

LP duality now establishes the result sought: First, suppose that the optimal value of the above primal problem over \mathbf{p} is zero. Then, the same is true for the dual problem. Hence, by the sign conditions $\mu_i \geq 0$ and $g_i(\bar{\mathbf{x}}) \leq 0$, each term in the sum must be zero. Hence, we established that *complementarity* holds. Next, the two constraints in the dual problem are precisely the *dual feasibility* conditions, which hence are fulfilled. Finally, *primal feasibility* of $\bar{\mathbf{x}}$ was assumed. It follows that this vector indeed is a KKT point.

Conversely, if $\bar{\mathbf{x}}$ is a KKT point, then the dual problem above has a feasible solution given by any KKT multiplier vector $(\boldsymbol{\mu}, \boldsymbol{\lambda})$. The dual objective is upper bounded by zero, since each term in the sum is non-positive. On the other hand, there is a feasible solution with the objective value 0, namely any KKT point! So, each KKT point must constitute an optimal solution to this dual LP problem! It then follows by duality theory that the dual of this problem, which is precisely the primal problem in \mathbf{p} above, has a finite optimal solution, whose optimal value must then be zero. We are done.

[Note: The LP problem given in the exercise is essentially the subproblem in the *Sequential Linear Programming* (SLP) algorithm. By the above analysis, the optimal value must be negative if $\bar{\mathbf{x}}$ is not a KKT point, and it must therefore also be negative (since a zero value is given by setting $\mathbf{p} = \mathbf{0}^n$). The optimal value of \mathbf{p} , if one exists, is therefore a descent direction with respect to f at $\bar{\mathbf{x}}$. A convergent SLP method introduces additional box constraints on \mathbf{p} in the LP subproblem to make sure that the solution is finite, and the update is made according to a line search with respect to some penalty function.]

Exercise 13.8 —

References

- [Aba67] J. ABADIE, *On the Kuhn–Tucker theorem*, in Nonlinear Programming (NATO Summer School, Menton, 1964), North-Holland, Amsterdam, 1967, pp. 19–36.
- [AMO93] R. K. AHUJA, T. L. MAGNANTI, AND J. B. ORLIN, *Network Flows: Theory, Algorithms, and Applications*, Prentice-Hall, Englewood Cliffs, NJ, 1993.
- [Arm66] L. ARMIJO, *Minimization of functions having Lipschitz continuous first partial derivatives*, Pacific Journal of Mathematics, 16 (1966), pp. 1–3.
- [AHU58] K. J. ARROW, L. HURWICZ, AND H. UZAWA, eds., *Studies in Linear and Non-Linear Programming*, Stanford University Press, Stanford, CA, 1958.
- [AHU61] K. J. ARROW, L. HURWICZ, AND H. UZAWA, *Constraint qualifications in maximization problems*, Naval Research Logistics Quarterly, 8 (1961), pp. 175–191.
- [Avr76] M. AVRIEL, *Nonlinear Programming: Analysis and Methods*, Prentice Hall Series in Automatic Computation, Prentice Hall, Englewood Cliffs, NJ, 1976.
- [AvG96] M. AVRIEL AND B. GOLANY, eds., *Mathematical Programming for Industrial Engineers*, vol. 20 of Industrial Engineering, Marcel Dekker, New York, NY, 1996.
- [Ban22] S. BANACH, *Sur les opérations dans les ensembles abstraits et leur application aux équations intégrales*, Fundamenta Mathematicae, 3 (1922), pp. 133–181.
- [Bar71] R. H. BARTELS, *A stabilization of the simplex method*, Numerische Mathematik, 16 (1971), pp. 414–434.
- [BaG69] R. H. BARTELS AND G. H. GOLUB, *The simplex method of linear programming using LU-decomposition*, Communications of the ACM, 12 (1969), pp. 266–268 and 275–278.
- [BSS93] M. S. BAZARAA, H. D. SHERALI, AND C. M. SHETTY, *Nonlinear Programming: Theory and Algorithms*, John Wiley & Sons, New York, NY, second ed., 1993.

References

- [Ben62] J. F. BENDERS, *Partitioning procedures for solving mixed variables programming problems*, Numerische Mathematik, 4 (1962), pp. 238–252.
- [Ber99] D. P. BERTSEKAS, *Nonlinear Programming*, Athena Scientific, Belmont, MA, second ed., 1999.
- [Ber04] ———, *Lagrange multipliers with optimal sensitivity properties in constrained optimization*, Report LIDS 2632, Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA, 2004.
- [BNO03] D. P. BERTSEKAS, A. NEDIĆ, AND A. E. OZDAGLAR, *Convex Analysis and Optimization*, Athena Scientific, Belmont, MA, 2003.
- [BeT89] D. P. BERTSEKAS AND J. N. TSITSIKLIS, *Parallel and Distributed Computation: Numerical Methods*, Prentice Hall, London, U.K., 1989.
- [BeT00] D. P. BERTSEKAS AND J. N. TSITSIKLIS, *Gradient convergence in gradient methods with errors*, SIAM Journal on Optimization, 10 (2000), pp. 627–642.
- [Bla77] R. G. BLAND, *New finite pivoting rules for the simplex method*, Mathematics of Operations Research, 2 (1977), pp. 103–107.
- [BIO72] E. BLUM AND W. OETTLI, *Direct proof of the existence theorem in quadratic programming*, Operations Research, 20 (1972), pp. 165–167.
- [BGLS03] J. F. BONNANS, J. C. GILBERT, C. LEMARÉCHAL, AND C. A. SAGASTIZÁBAL, *Numerical Optimization: Theoretical and Practical Aspects*, Universitext, Springer-Verlag, Berlin, 2003. Translated from the original French edition, published by Springer-Verlag 1997.
- [BoS00] J. F. BONNANS AND A. SHAPIRO, *Perturbation Analysis of Optimization Problems*, Springer Series in Operations Research, Springer-Verlag, New York, NY, 2000.
- [BoL00] J. M. BORWEIN AND A. S. LEWIS, *Convex Analysis and Nonlinear Optimization: Theory and Examples*, CMS Books in Mathematics, Springer-Verlag, New York, NY, 2000.
- [BHM77] S. P. BRADLEY, A. C. HAX, AND T. L. MAGNANTI, *Applied Mathematical Programming*, Addison-Wesley, Reading, MA, 1977.
- [Bre73] R. P. BRENT, *Algorithms for Minimization Without Derivatives*, Prentice Hall Series in Automatic Computation, Prentice Hall, Englewood Cliffs, NJ, 1973. Reprinted by Dover Publications, Inc., Mineola, NY, 2002.
- [Bro09] L. E. J. BROUWER, *On continuous vector distributions on surfaces*, Amsterdam Proceedings, 11 (1909).
- [Bro12] ———, *Über Abbildung von Mannigfaltigkeiten*, Mathematische Annalen, 71 (1912), pp. 97–115.
- [Bro70] C. G. BROYDEN, *The convergence of single-rank quasi-Newton methods*, Mathematics of Computation, 24 (1970), pp. 365–382.

- [BGIS95] R. BURACHIK, L. M. G. DRUMMOND, A. N. IUSEM, AND B. F. SVAITER, *Full convergence of the steepest descent method with inexact line searches*, Optimization, 32 (1995), pp. 137–146.
- [BuF91] J. V. BURKE AND M. C. FERRIS, *Characterization of solution sets of convex programs*, Operations Research Letters, 10 (1991), pp. 57–60.
- [Car07] C. CARATHÉODORY, *Über den Variabilitätsbereich der Koeffizienten von Potenzreihen, die gegebene Werte nicht annehmen*, Mathematische Annalen, 64 (1907), pp. 95–115.
- [Car11] ———, *Über den Variabilitätsbereich der Fourier’schen Konstanten von positiven harmonischen Funktionen*, Rendiconti del Circolo Matematico di Palermo, 32 (1911), pp. 193–217.
- [Casetal02] E. CASTILLO, A. J. CONEJO, P. R. G. PEDREGAL, AND N. ALGUACIL, *Building and Solving Mathematical Programming Models in Engineering and Science*, Pure and Applied Mathematics, John Wiley & Sons, New York, NY, 2002.
- [Cau1847] A. CAUCHY, *Méthode générale pour la résolution des systèmes d’équations simultanées*, Comptes Rendus Hebdomadaires des Séances de l’Académie des Sciences (Paris), Série A, 25 (1847), pp. 536–538.
- [Cha52] A. CHARNES, *Optimality and degeneracy in linear programming*, Econometrica, 20 (1952), pp. 160–170.
- [Chv83] V. CHVÁTAL, *Linear Programming*, Freeman, New York, NY, 1983.
- [CGT00] A. R. CONN, N. I. M. GOULD, AND PH. L. TOINT, *Trust-Region Methods*, vol. 1 of MPS/SIAM Series on Optimization, SIAM and Mathematical Programming Society, Philadelphia, PA, 2000.
- [Cro36] H. CROSS, *Analysis of flow in networks of conduits or conductors*, Bulletin 286, Engineering Experiment Station, University of Illinois, Urbana, IL, 1936.
- [Dan51] G. B. DANTZIG, *Maximization of a linear function of variables subject to linear inequalities*, in Activity Analysis of Production and Allocation, Tj. C. Koopmans, ed., New York, NY, 1951, John Wiley & Sons, pp. 339–347.
- [Dan53] ———, *Computational algorithm of the revised simplex method*, Report RM 1266, The Rand Corporation, Santa Monica, CA, 1953.
- [Dan57] ———, *Concepts, origins, and use of linear programming*, in Proceedings of the First International Conference on Operational Research, Oxford, 1957, M. Davies, R. T. Eddison, and T. Page, eds., London, U.K., 1957, The English Universities Press, pp. 100–108.
- [Dan63] ———, *Linear Programming and Extensions*, Princeton University Press, Princeton, NJ, 1963.
- [DaM05] G. B. DANTZIG AND N. T. MUKUND, *Linear programming 3: Implementation*, Springer Series in Operations Research, Springer-Verlag, New York, NY, 2005.

References

- [DaO53] G. B. DANTZIG AND A. ORDEN, *Notes on linear programming: Part 2, duality theorems*, technical report RM-1265, The Rand Corporation, Santa Monica, CA, 1953.
- [DOW55] G. B. DANTZIG, A. ORDEN, AND P. WOLFE, *The generalized simplex method for minimizing a linear form under linear inequality restraints*, Pacific Journal of Mathematics, 5 (1955), pp. 183–195.
- [DaT97] G. B. DANTZIG AND M. N. THAPA, *Linear programming 1: Introduction*, Springer Series in Operations Research, Springer-Verlag, New York, NY, 1997.
- [DaT03] ———, *Linear programming 2: Theory and Extensions*, Springer Series in Operations Research, Springer-Verlag, New York, NY, 2003.
- [DaW60] G. B. DANTZIG AND P. WOLFE, *Decomposition principle for linear programs*, Operations Research, 8 (1960), pp. 101–111.
- [dAu47] A. D’AURIAC, *A propos de l’unicité de solution dans les problèmes de réseaux maillés*, La Houille Blanche, 2 (1947), pp. 209–211.
- [Dav59] W. C. DAVIDON, *Variable metric method for minimization*, Report ANL-5990 Rev, Argonne National Laboratories, Argonne, IL, 1959. Also published in SIAM Journal on Optimization, 1 (1991), pp. 1–17.
- [DeF49] B. DE FINETTI, *Sulla stratificazioni convesse*, Annali di Matematica Pura ed Applicata, 30 (1949), pp. 173–183.
- [Den59] J. B. DENNIS, *Mathematical Programming and Electrical Networks*, John Wiley & Sons, New York, NY, 1959.
- [DeS83] J. E. DENNIS AND R. E. SCHNABEL, *Numerical Methods for Unconstrained Optimization and Nonlinear Equations*, Prentice Hall, Englewood Cliffs, NJ, 1983.
- [DiJ79] Y. M. I. DIRICKX AND L. P. JENNERGREN, *Systems Analysis by Multilevel Methods: With Applications to Economics and Management*, vol. 6 of International Series on Applied Systems Analysis, John Wiley & Sons, Chichester, U.K., 1979.
- [Duf46] R. J. DUFFIN, *Nonlinear networks, I*, Bulletin of the American Mathematical Society, 52 (1946), pp. 833–838.
- [Duf47] ———, *Nonlinear networks, IIa*, Bulletin of the American Mathematical Society, 53 (1947), pp. 963–971.
- [DuH78] J. C. DUNN AND S. HARSHBARGER, *Conditional gradient algorithms with open loop step size rules*, Journal of Mathematical Analysis and Applications, 62 (1978), pp. 432–444.
- [Eav71] B. C. EAVES, *On quadratic programming*, Management Science, 17 (1971), pp. 698–711.
- [EHL01] T. F. EDGAR, D. M. HIMMELBLAU, AND L. S. LASDON, *Optimization of Chemical Processes*, McGraw-Hill, New York, NY, second ed., 2001.

- [Eke74] I. EKELAND, *On the variational principle*, Journal of Mathematical Analysis and Applications, 47 (1974), pp. 324–353.
- [Erm66] YU. M. ERMOL'EV, *Methods for solving nonlinear extremal problems*, Kibernetika, 2 (1966), pp. 1–17. In Russian, translated into English in Cybernetics, 2 (1966), pp. 1–14.
- [Eva70] J. P. EVANS, *On constraint qualifications in nonlinear programming*, Naval Research Logistics Quarterly, 17 (1970), pp. 281–286.
- [Eve63] H. EVERETT, III, *Generalized Lagrange multiplier method for solving problems of optimum allocation of resources*, Operations Research, 11 (1963), pp. 399–417.
- [Fac95] F. FACCHINEI, *Minimization of SC^1 functions and the Maratos effect*, Operations Research Letters, 17 (1995), pp. 131–137.
- [Fal67] J. E. FALK, *Lagrange multipliers and nonlinear programming*, Journal of Mathematical Analysis and Applications, 19 (1967), pp. 141–159.
- [Far1902] J. FARKAS, *Über die Theorie der einfachen Ungleichungen*, Journal für die Reine und Angewandte Mathematik, 124 (1902), pp. 1–24.
- [Fen51] W. FENCHEL, *Convex cones, sets and functions*, mimeographed lecture notes, Princeton University, Princeton, NY, 1951.
- [Fia83] A. V. FIACCO, *Introduction to sensitivity and stability analysis in nonlinear programming*, vol. 165 of Mathematics in Science and Engineering, Academic Press Inc., Orlando, FL, 1983.
- [FiM68] A. V. FIACCO AND G. P. MCCORMICK, *Nonlinear Programming: Sequential Unconstrained Minimization Techniques*, John Wiley & Sons, New York, NY, 1968. Also published as volume 4 in the Classics in Applied Mathematics Series, SIAM, Philadelphia, PA, 1990.
- [Fis81] M. L. FISHER, *The Lagrangian relaxation method for solving integer programming problems*, Management Science, 27 (1981), pp. 1–18.
- [Fis85] ———, *An applications oriented guide to Lagrangian relaxation*, Interfaces, 15 (1985), pp. 10–21.
- [Fle70] R. FLETCHER, *A new approach to variable metric algorithms*, Computer Journal, 13 (1970), pp. 317–322.
- [Fle87] ———, *Practical Methods of Optimization*, John Wiley & Sons, Chichester, U.K., second ed., 1987.
- [FLT02] R. FLETCHER, S. LEYFFER, AND PH. L. TOINT, *On the global convergence of a filter-SQP algorithm*, SIAM Journal on Optimization, 13 (2002), pp. 44–59.
- [FIP63] R. FLETCHER AND M. J. D. POWELL, *A rapidly convergent descent method for minimization*, Computer Journal, 6 (1963), pp. 163–168.

References

- [FlR64] R. FLETCHER AND C. M. REEVES, *Function minimization by conjugate gradients*, Computer Journal, 7 (1964), pp. 149–154.
- [FrW56] M. FRANK AND P. WOLFE, *An algorithm for quadratic programming*, Naval Research Logistics Quarterly, 3 (1956), pp. 95–110.
- [GKT51] D. H. GALE, H. W. KUHN, AND A. W. TUCKER, *Linear programming and the theory of games*, in Activity Analysis of Production and Allocation, Tj. C. Koopmans, ed., New York, NY, 1951, Wiley, pp. 317–329.
- [GaJ79] M. R. GAREY AND D. S. JOHNSON, *Computers and Intractability: A Guide to the Theory of NP-Completeness*, Freeman, New York, NY, 1979.
- [GaA05] S. I. GASS AND A. A. ASSAD, *An Annotated Timeline of Operations Research. An Informal History*, vol. 75 of International Series in Operations Research & Management Science, Kluwer Academic Publishers, New York, NY, 2005.
- [Geo74] A. M. GEOFFRION, *Lagrangian relaxation for integer programming: Approaches to integer programming*, Mathematical Programming Study, 2 (1974), pp. 82–114.
- [Gil66] E. G. GILBERT, *An iterative procedure for computing the minimum of a quadratic form on a convex set*, SIAM Journal on Control, 4 (1966), pp. 61–80.
- [GiM73] P. E. GILL AND W. MURRAY, *A numerically stable form of the simplex algorithm*, Linear Algebra and Its Applications, 7 (1973), pp. 99–138.
- [GMS05] P. E. GILL, W. MURRAY, AND M. A. SAUNDERS, *SNOPT: An SQP algorithm for large-scale constrained optimization*, SIAM Review, 47 (2005), pp. 99–131.
- [GMSW89] P. E. GILL, W. MURRAY, M. A. SAUNDERS, AND M. H. WRIGHT, *A practical anti-cycling procedure for linearly constrained optimization*, Mathematical Programming, 45 (1989), pp. 437–474.
- [Gol70] D. GOLDFARB, *A family of variable-metric methods derived by variational means*, Mathematics of Computation, 24 (1970), pp. 23–26.
- [Gol64] A. A. GOLDSTEIN, *Convex programming in Hilbert space*, Bulletin of the American Mathematical Society, 70 (1964), pp. 709–710.
- [GrD03] A. GRANAS AND J. DUGUNDJI, *Fixed Point Theory*, Springer Monographs in Mathematics, Springer-Verlag, New York, NY, 1969.
- [Gri00] A. GRIEWANK, *Evaluating Derivatives: Principles and Techniques of Algorithmic Differentiation*, vol. 19 of Frontiers in Applied Mathematics, SIAM, Philadelphia, PA, 2000.
- [Gui69] M. GUIGNARD, *Generalized Kuhn–Tucker conditions for mathematical programming problems in a Banach space*, SIAM Journal on Control, 7 (1969), pp. 232–241.

- [Had10] J. HADAMARD, *Sur quelques applications de l'indice de Kronecker*, in Introduction à la théorie des fonctions d'une variable, J. Tannary, ed., vol. 2, Hermann, Paris, 1910, pp. 875–915.
- [Han75] S. P. HAN, *Penalty Lagrangian methods in a quasi-Newton approach*, Report TR 75-252, Computer Science, Cornell University, Ithaca, NY, 1975.
- [HaH96] G. K. HAUER AND H. M. HOGANSON, *Tailoring a decomposition method to a large forest management scheduling problem in northern Ontario*, INFOR, 34 (1996), pp. 209–231.
- [HLV87] D. W. HEARN, S. LAWPHONGPANICH, AND J. A. VENTURA, *Restricted simplicial decomposition: Computation and extensions*, Mathematical Programming Study, 31 (1987), pp. 99–118.
- [HWC74] M. HELD, P. WOLFE, AND H. P. CROWDER, *Validation of subgradient optimization*, Mathematical Programming, 6 (1974), pp. 62–88.
- [HiL93] J.-B. HIRIART-URRUTY AND C. LEMARÉCHAL, *Convex Analysis and Minimization Algorithms*, vol. 305 and 306 of Grundlehren der mathematischen Wissenschaften, Springer-Verlag, Berlin, 1993.
- [Hof53] A. HOFFMAN, *Cycling in the simplex algorithm*, Report 2974, National Bureau of Standards, Gaithersburg, MD, 1953.
- [Ius03] A. N. IUSEM, *On the convergence properties of the projected gradient method for convex optimization*, Computational and Applied Mathematics, 22 (2003), pp. 37–52.
- [Joh48] F. JOHN, *Extremum problems with inequalities as subsidiary conditions*, in Studies and Essays Presented to R. Courant on his 60th Birthday, January 8, 1948, Interscience Publishers, Inc., New York, NY, 1948, pp. 187–204.
- [JoM74] L. A. JOHNSON AND D. S. MONTGOMERY, *Operations Research in Production Planning, Scheduling and Inventory Control*, John Wiley & Sons, New York, NY, 1974.
- [Jos03] M. JOSEFSSON, *Sensitivity analysis of traffic equilibria*, Master's thesis, Department of Mathematics, Chalmers University of Technology, Gothenburg, Sweden, 2003.
- [Kar84a] N. KARMARKAR, *A new polynomial-time algorithm for linear programming*, in Proceedings of the 16th Annual ACM Symposium on Theory of Computing, STOC'84 (Washington, DC, April 30–May 2, 1984), New York, NY, 1984, ACM Press, pp. 302–311.
- [Kar84b] ———, *A new polynomial-time algorithm for linear programming*, Combinatorica, 4 (1984), pp. 373–395.
- [Kha79] L. G. KHACHIYAN, *A polynomial algorithm in linear programming*, Doklady Akademii Nauk SSSR, 244 (1979), pp. 1093–1096.
- [Kha80] ———, *Polynomial algorithms in linear programming*, Akademiya Nauk SSSR. Zhurnal Vychislitel'noy Matematiki i Matematicheskoy Fiziki, 20 (1980), pp. 51–68.

References

- [Kir1847] G. KIRCHHOFF, *Über die Ausflösung der Gleichungen auf welche man bei der Untersuchungen der Linearen Vertheilung Galvanischer Ströme geführt wird*, Pogendorff Annalen Der Physik, 72 (1847), pp. 497–508. English translation, IRE Transactions on Circuit Theory, CT-5 (1958), pp. 4–8.
- [KIM72] V. KLEE AND G. J. MINTY, *How good is the simplex algorithm?*, in Inequalities, III. Proceedings of the Third Symposium on Inequalities held at the University of California, Los Angeles, CA, September 1–9, 1969; dedicated to the memory of Theodore S. Motzkin, O. Shisha, ed., New York, NY, 1972, Academic Press, pp. 159–175.
- [KLT03] T. G. KOLDA, R. M. LEWIS, AND V. TORCZON, *Optimization by direct search: New perspectives on some classical and modern methods*, SIAM Review, 45 (2003), pp. 385–482.
- [Kre78] E. KREYSZIG, *Introductory Functional Analysis with Applications*, John Wiley & Sons, New York, NY, 1978.
- [KuT51] H. W. KUHN AND A. W. TUCKER, *Nonlinear programming*, in Proceedings of the Second Berkeley Symposium on Mathematical Statistics and Probability, 1950, Berkeley and Los Angeles, CA, 1951, University of California Press, pp. 481–492.
- [LaP05] T. LARSSON AND M. PATRIKSSON, *Global optimality conditions for discrete and nonconvex optimization—with applications to Lagrangian heuristics and column generation*, technical report, Department of Mathematics, Chalmers University of Technology, Gothenburg, Sweden, 2005. To appear in *Operations Research*.
- [LPS96] T. LARSSON, M. PATRIKSSON, AND A.-B. STRÖMBERG, *Conditional subgradient optimization—theory and applications*, European Journal of Operational Research, 88 (1996), pp. 382–403.
- [LPS99] ———, *Ergodic, primal convergence in dual subgradient schemes for convex programming*, Mathematical Programming, 86 (1999), pp. 283–312.
- [Las70] L. S. LASDON, *Optimization Theory for Large Systems*, Macmillan, New York, NY, 1970.
- [Law76] E. LAWLER, *Combinatorial Optimization: Networks and Matroids*, Holt, Rinehart and Winston, New York, NY, 1976.
- [LRS91] J. K. LENSTRA, A. H. G. RINNOOY KAN, AND A. SCHRIJVER, eds., *History of Mathematical Programming. A Collection of Personal Reminiscences*, North-Holland, Amsterdam, 1991.
- [LeP66] E. S. LEVITIN AND B. T. POLYAK, *Constrained minimization methods*, USSR Computational Mathematics and Mathematical Physics, 6 (1966), pp. 1–50.
- [Lip1877] R. LIPSCHITZ, *Lehrbuch der Analysis*, Cohn & Sohn, Leipzig, 1877.
- [Lue84] D. G. LUENBERGER, *Linear and Nonlinear Programming*, Addison Wesley, Reading, MA, second ed., 1984. Reprinted by Kluwer Academic Publishers, Boston, MA, 2003.

- [Man65] O. L. MANGASARIAN, *Pseudo-convex functions*, SIAM Journal on Control, 3 (1965), pp. 281–290.
- [Man69] ———, *Nonlinear Programming*, McGraw-Hill, New York, NY, 1969. Also published as volume 10 in the Classics in Applied Mathematics Series, SIAM, Philadelphia, PA, 1994.
- [Man88] ———, *A simple characterization of solution sets of convex programs*, Operations Research Letters, 7 (1988), pp. 21–26.
- [MaF67] O. L. MANGASARIAN AND S. FROMOVITZ, *The Fritz John necessary optimality conditions in the presence of equality and inequality constraints*, Journal of Mathematical Analysis and Applications, 17 (1967), pp. 37–47.
- [Mar78] N. MARATOS, *Exact penalty function algorithms for finite dimensional and control optimization problems*, PhD thesis, Imperial College of Science and Technology, University of London, London, U.K., 1978.
- [Max1865] J. C. MAXWELL, *A dynamical theory of the electromagnetic field*, Philosophical Transactions of the Royal Society of London, 155 (1865), pp. 459–512.
- [Min10] H. MINKOWSKI, *Geometrie der Zahlen*, Teubner, Leipzig, 1910.
- [Min11] ———, *Gesammelte Abhandlungen*, vol. II, Teubner, Leipzig, 1911, ch. Theorie der konvexen Körper, Insbesondere Begründung ihres Oberflächenbegriffs.
- [Mot36] T. MOTZKIN, *Beiträge zur Theorie der linearen Ungleichungen*, Azriel, Israel, 1936.
- [Mur83] K. G. MURTY, *Linear Programming*, John Wiley & Sons, New York, NY, 1983.
- [Mur95] ———, *Operations Research: Deterministic Optimization Models*, Prentice Hall, Englewood Cliffs, NJ, 1995.
- [Nas50] J. F. NASH, JR., *Equilibrium points in n-person games*, Proceedings of the National Academy of Sciences of the United States of America, 36 (1950), pp. 48–49.
- [Nas51] ———, *Non-cooperative games*, Annals of Mathematics, 54 (1951), pp. 286–295.
- [NaS96] S. G. NASH AND A. SOFER, *Linear and Nonlinear Programming*, MacGraw-Hill, Singapore, 1996.
- [NeW88] G. L. NEMHAUSER AND L. WOLSEY, *Integer and Combinatorial Optimization*, Wiley-Interscience Series in Discrete Mathematics and Optimization, John Wiley & Sons, New York, NY, 1988.
- [New1687] I. S. NEWTON, *Philosophiae Naturalis Principia Mathematica*, London, U.K., 1687.
- [NoW99] J. NOCEDAL AND S. J. WRIGHT, *Numerical Optimization*, Springer Series in Operations Research, Springer-Verlag, New York, NY, 1999.

References

- [Orc54] W. ORCHARD-HAYS, *Background, development and extensions of the revised simplex method*, Report RM 1433, The Rand Corporation, Santa Monica, CA, 1954.
- [OrR70] J. M. ORTEGA AND W. C. RHEINBOLDT, *Iterative Solution of Nonlinear Equations in Several Variables*, Academic Press, New York, NY, 1970. Also published as volume 30 in the Classics in Applied Mathematics Series, SIAM, Philadelphia, PA, 2000.
- [Pad99] M. PADBERG, *Linear Optimization and Extensions*, vol. 12 of Algorithms and Combinatorics, Springer-Verlag, Berlin, second ed., 1999.
- [PaT91] E. R. PANIER AND A. L. TITS, *Avoiding the Maratos effect by means of a nonmonotone line search, I. General constrained problems*, SIAM Journal on Numerical Analysis, 28 (1991), pp. 1183–1195.
- [PaS82] C. H. PAPADIMITRIOU AND K. STEIGLITZ, *Combinatorial Optimization: Algorithms and Complexity*, Prentice Hall, Englewood Cliffs, 1982.
- [PaS88] P. M. PARDALOS AND G. SCHNITGER, *Checking local optimality in constrained quadratic programming is NP-hard*, Operations Research Letters, 7 (1988), pp. 33–35.
- [PaV91] P. M. PARDALOS AND S. VAVASIS, *Quadratic programming with one negative eigenvalue is NP-hard*, Journal of Global Optimization, 1 (1991), pp. 15–22.
- [Pat94] M. PATRIKSSON, *The Traffic Assignment Problem—Models and Methods*, Topics in Transportation, VSP BV, Utrecht, The Netherlands, 1994.
- [Pat98] ———, *Nonlinear Programming and Variational Inequalities: A Unified Approach*, vol. 23 of Applied Optimization, Kluwer Academic Publishers, Dordrecht, The Netherlands, 1998.
- [PoR69] E. POLAK AND G. RIBIÈRE, *Note sur la convergence de méthodes de directions conjuguées*, Revue Française d'Information et Recherche Opérationnelle, 3 (1969), pp. 35–43.
- [Pol69] B. T. POLYAK, *Minimization of unsmooth functionals*, USSR Computational Mathematics and Mathematical Physics, 9 (1969), pp. 14–29.
- [Pow78] M. J. D. POWELL, *A fast algorithm for nonlinearly constrained optimization calculations*, in Numerical Analysis, Proceedings of the Seventh Biennial Conference held at the University of Dundee, Dundee, June 28–July 1, 1977, G. A. Watson, ed., vol. 630 of Lecture Notes in Mathematics, Berlin, 1978, Springer-Verlag, pp. 144–157.
- [PsD78] B. N. PSHENICHNYJ AND YU. M. DANILIN, *Numerical Methods in Extremal Problems*, MIR Publishers, Moscow, 1978.
- [Rad19] H. RADEMACHER, *Über partielle und totale Differenzierbarkeit von*

- Funktionen mehrerer Variablen unter über die Transformation der Doppelintegrale*, Mathematische Annalen, 79 (1919), pp. 340–359.
- [Rar98] R. L. RARDIN, *Optimization in Operations Research*, Prentice Hall, Englewood Cliffs, NJ, 1998.
- [Roc70] R. T. ROCKAFELLAR, *Convex Analysis*, Princeton University Press, Princeton, NJ, 1970. Reprinted in the series *Princeton Landmarks in Mathematics* by Princeton University Press, Princeton, NJ, 1997.
- [RoW97] R. T. ROCKAFELLAR AND R. J.-B. WETS, *Variational Analysis*, vol. 317 of Grundlehren der mathematischen Wissenschaften, Springer-Verlag, Berlin, 1997.
- [Sau72] M. A. SAUNDERS, *Large-scale linear programming using the Cholesky factorization*, technical report Stan-CS-72-252, Computer Sciences Department, Stanford University, Stanford, CA, 1972.
- [Sch86] A. SCHRIJVER, *Theory of Linear and Integer Programming*, Wiley, Chichester, 1986.
- [Sch03] ———, *Combinatorial optimization*, vol. 24 of Algorithms and Combinatorics, Springer-Verlag, Berlin, 2003.
- [Sha70] D. F. SHANNO, *Conditioning of quasi-Newton methods for function minimization*, Mathematics of Computation, 24 (1970), pp. 647–656.
- [She85] Y. SHEFFI, *Urban Transportation Networks: Equilibrium Analysis with Mathematical Programming Methods*, Prentice-Hall, Englewood Cliffs, NJ, 1985.
- [She76] M. A. SHEPILOV, *Method of the generalized gradient for finding the absolute minimum of a convex function*, Cybernetics, 12 (1976), pp. 547–553.
- [Sho77] N. Z. SHOR, *Cut-off method with space extension in convex programming problems*, Cybernetics, 13 (1977), pp. 94–96.
- [Sho85] ———, *Minimization Methods for Non-Differentiable Functions*, Springer-Verlag, Berlin, 1985. Translated from the Russian by K. C. Kiwiel and A. Ruszczyński.
- [StW70] J. STOER AND C. WITZGALL, *Convexity and Optimization in Finite Dimensions I*, Springer-Verlag, Berlin, 1970.
- [Tah03] H. A. TAHA, *Operations Research: An Introduction*, Prentice Hall, Englewood Cliffs, NJ, seventh ed., 2003.
- [UUV04] M. ULBRICH, S. ULBRICH, AND L. N. VICENTE, *A globally convergent primal-dual interior-point filter method for nonlinear programming*, Mathematical Programming, 100 (2004), pp. 379–410.
- [Van01] R. J. VANDERBEI, *Linear Programming. Foundations and Extensions*, vol. 37 of International Series in Operations Research & Management Science, Kluwer Academic Publishers, Boston, MA, second ed., 2001.

References

- [vHo77] B. VON HOHENBALKEN, *Simplicial decomposition in nonlinear programming algorithms*, Mathematical Programming, 13 (1977), pp. 49–68.
- [vNe28] J. VON NEUMANN, *Zur Theorie der Gesellschaftsspiele*, Mathematische Annalen, 100 (1928), pp. 295–320.
- [vNe47] ———, *On a maximization problem*, unpublished manuscript, Institute for Advanced Study, Princeton, NJ, 1947.
- [vNM43] J. VON NEUMANN AND O. MORGENSTERN, *Theory of Games and Economic Behavior*, Princeton University Press, Princeton, NJ, 1943.
- [Wag75] H. M. WAGNER, *Principles of Operations Research: With Applications to Managerial Decisions*, Prentice Hall, Englewood Cliffs, NJ, second ed., 1975.
- [War52] J. G. WARDROP, *Some theoretical aspects of road traffic research*, Proceedings of the Institute of Civil Engineers, Part II, (1952), pp. 325–378.
- [Wil99] H. P. WILLIAMS, *Model Building in Mathematical Programming*, John Wiley & Sons, Chichester, UK, fourth ed., 1999.
- [Wil63] R. B. WILSON, *A simplicial algorithm for concave programming*, PhD thesis, Graduate School of Business Administration, Harvard University, Cambridge, MA, 1963.
- [Wol69] P. WOLFE, *Convergence conditions for ascent methods*, SIAM Review, 11 (1969), pp. 226–235.
- [Wol75] ———, *A method of conjugate subgradients for minimizing nondifferentiable functions*, Mathematical Programming Study, 3 (1975), pp. 145–173.
- [Wol98] L. A. WOLSEY, *Integer Programming*, Wiley-Interscience Series in Discrete Mathematics and Optimization, John Wiley & Sons, New York, NY, 1998.
- [YuN77] D. B. YUDIN AND A. S. NEMIROVSKII, *Informational complexity and efficient methods for the solution of convex extremal problems*, Matekon, 13 (1977), pp. 25–45.
- [Zan69] W. I. ZANGWILL, *Nonlinear Programming: A Unified Approach*, Prentice Hall, Englewood Cliffs, NJ, 1969.

Index

- Abadie's CQ, 125
- active constraint ($\mathcal{I}(\boldsymbol{x})$), 89
- adjacent extreme point, 220
- affine combination, 43
- affine function, 11, 59
- affine hull, 43
- affine independence, 34
- affine subspace, 34
- affine transformation, 300
- algebraic characterization of adjacency, 221
- approximate line search, 276
- Armijo step, 277, 298, 312
- artificial variables, 233
- augmented Lagrangian function, 350
- augmented Lagrangian method, 350
- automatic differentiation, 297

- Banach's Theorem, 101
- barrier function, 330
- barrier problem, 331
- basic feasible solution, 215
- basic solution, 215
- basic variables, 215
- basis, 35
- BFGS method, 274
- Bland's rule, 237
- boundary, 38
- bounded set, 37
- Brouwer's Theorem, 101
- bundle method, 166

- calculus rules, 39
- canonical form, 243
- Carathéodory's Theorem, 45

- Cartesian product set, 148
- Cauchy–Bunyakowski–Schwarz inequality, 34
- central difference formula, 297
- classification of optimization models, 11
- closed mapping, 159
- closed sets, 37
- closure, 37
- coercive function, 78
- column dropping, 308
- column generation, 6
- combinatorial optimization, 185
- complementarity, 148
- Complementary Slackness Theorem, 250
- composite function, 60, 104
- composite operator, 104
- concave function, 58
- cone, 50
- cone of feasible directions, 115
- conjugate direction, 286, 294
- conjugate gradient, 289
- conjugate gradient method, 289
- constrained optimization, 12, 88–96, 303–350
- constraint qualification (CQ), 17, 124, 125, 130, 131
- constraints, 4
- continuity, 96
- continuous function, 38
- continuous optimization, 12
- continuous relaxation, 185
- continuously differentiable function, 39

Index

- contractive operator, 101
- convergence rate, 296
 - geometric, 101, 170
 - linear, 296
 - quadratic, 296
 - superlinear, 296
- convex analysis, 41–72
- convex combination, 43
- convex function, 57, 96, 159
- convex hull, 43
- convex programming, 12
- convex set, 41
- coordinates, 35
- CQ, 124

- Danskin’s Theorem, 160
- Dantzig–Wolfe algorithm, 158
- decision science, 10
- decision variable, 6
- degenerate basic solution, 215
- descent direction, 85, 165
- descent lemma, 322
- DFP method, 293
- Diet problem, 10
- differentiability, 163
- differentiable function, 38
- differentiable optimization, 12
- Dijkstra’s algorithm, 320
- diode, 182
- direction of unboundedness, 226
- directional derivative, 38, 86, 159
- distance function, 67
- divergent series step length rule, 165, 307
- domination, 346
- dual feasible basis, 254
- dual infeasible basis, 254
- dual linear program, 242
- dual simplex algorithm, 255
- dual simplex method, 254
- duality gap, 145

- effective domain, 96, 144
- efficient frontier, 346
- eigenvalue, 36
- eigenvector, 36
- Ekeland’s variational principle, 110

- electrical circuit, 181
- electrical network, 181
- eligible entering variable, 238
- eligible leaving variable, 238
- epigraph, 60, 80
- ε -optimal solution, 96
- equality constraint, 11
- equivalent systems, 224
- Euclidean projection, 66
- Everett’s Theorem, 176
- exact penalty function, 340
- existence of optimal solution, 219
- exterior penalty method, 326–330
- extreme direction, 218
- extreme point, 46

- Farkas’ Lemma, 57, 249
- feasibility heuristic, 188
- feasible direction, 88, 89
- feasible solution, 5, 13
- feasible-direction methods, 303
- filter, 346
- filter-SQP methods, 346
- finite termination, 281
- finitely generated cone, 55
- fixed point, 100
- Fletcher–Reeves formula, 292
- forward difference formula, 297
- Frank–Wolfe algorithm, 305
- Frank–Wolfe Theorem, 82
- Fritz John conditions, 121

- Gauss–Seidel method, 105
- geometric convergence rate, 101, 170
- global minimum, 76
- global optimality conditions, 133, 147
- global optimum, 76
 - necessary and sufficient conditions, 88, 91
 - sufficient conditions, 133
- Golden section, 277
- gradient, 38
- gradient projection algorithm, 311
- gradient related, 271
- gradient related method, 279, 280
- Gram–Schmidt procedure, 289

- hard constraint, 18
- Hessian matrix, 39
- $\mathcal{I}(\mathbf{x})$, 89
- identity matrix \mathbf{I}^n , 36
- ill-conditioning, 346
- implicit function, 13, 40, 296
- Implicit Function Theorem, 164
- indicator function (χ_S), 167, 325
- inequality constraint, 11
- infimum, 14
- infinite-dimensional optimization, 13
- integer programming, 12, 13
- integrable function, 318
- integrality property, 13
- interior, 38
- interior penalty function, 107
- interior point algorithm, 238, 330–337
- interpolation, 277
- iso-cost line, 268
- iso-curve, 268
- Jacobi method, 105
- Jacobian, 39, 300, 318, 339
- Karmarkar's algorithm, 238
- Karush–Kuhn–Tucker (KKT) conditions, 125, 133
- Kirchhoff's laws, 182
- Lagrange function, 142
- Lagrange multiplier method, 158, 174
- Lagrange multiplier vector, 143, 179
- Lagrange multipliers, 122
- Lagrangian dual function, 143
- Lagrangian dual problem, 143
- Lagrangian duality, 141–194
- Lagrangian relaxation, 18, 142, 143, 185
- least-squares data fitting, 267
- level curve, 268
- level set ($\text{lev}_g(b)$), 65, 66, 78, 80, 151, 167, 279, 281, 313
- Levenberg–Marquardt, 273, 301
- LICQ, 132
- limit, 37
- limit points, 37
- line search, 275
 - approximate, 276
 - Armijo step length rule, 277, 298, 312
 - Golden section, 277
 - interpolation, 277
 - Newton's method, 277
- linear convergence rate, 296
- linear function, 40
- linear independence, 34
- linear programming, 11, 13, 154, 197–264, 336–337
- linear programming duality, 241–264
- linear space, 34
- linear-fractional programming, 223
- Lipschitz continuity, 280
- local convergence, 339
- local minimum, 76
- local optimum, 76
 - necessary conditions, 85, 86, 90, 117, 121, 125, 130
 - sufficient conditions, 87
- logarithmic barrier, 331
- logical constraint, 5
- lower semi-continuity, 79
- Maratos effect, 344
- mathematical model, 4
- mathematical programming, 9
- matrix, 35
- matrix game, 105
- matrix inverse, 36
- matrix norm, 35
- matrix product, 35
- matrix transpose, 35
- max function, 160
- mean-value theorem, 39
- merit function, 341
- method of successive averages (MSA), 321
- MFCQ, 131
- minimax theorem, 105
- minimum, 14

Index

- minimum distance (dist_S), 167
- multi-objective optimization, 13, 346
- near-optimality, 95
- negative curvature, 270
- neighbourhood, 38
- Newton's method, 272, 277, 299
- Newton–Raphson method, 105, 272
- Nobel laureates, 10
- non-basic variables, 215
- non-convex programming, 12
- non-coordinability, 176
- non-differentiable function, 283
- non-differentiable optimization, 12
- non-expansive operator, 99
- nonlinear programming, 11, 13
- nonsingular matrix, 36
- norm, 34
- normal cone (N_X), 94, 126
- NP-hard problem, 77, 186
- objective function, 4
- Ohm's law, 183
- open ball, 37
- open set, 37
- operations research, 9
- optimal BFS, 226
- optimal solution, 5
- optimal value, 5
- optimality, 9
- optimality conditions, 84–88, 90–94, 111–140, 147–149, 175, 227, 228, 252–253
- optimization under uncertainty, 13
- optimize, 3
- orthogonality, 34, 148
- orthonormal basis, 35
- parametric optimization, 137
- Pareto set, 346
- partial pricing, 230
- pattern search methods, 298
- penalty, 19
- penalty function, 19
- penalty parameter, 326
- perturbation function ($p(\mathbf{u})$), 178
- Phase I, 315
- phase I problem, 233
- phase II problem, 233
- physical constraint, 5
- piece-wise linear function, 283
- Polak–Ribière formula, 292
- Polyak step, 165
- polyhedral cone, 50
- polyhedron, 47
- polytope, 45
- positive (semi)definite matrix, 37
- potential, 181
- potential difference, 181
- pre-conditioning, 292
- primal infeasibility criterion, 254
- primal simplex method, 254
- projection, 66
- projection arc, 312
- projection operator, 66, 91, 99
- projection problem, 315
- proof by contradiction, 33
- proper function, 15, 167
- proximal point algorithm, 301
- pseudo-convex function, 109
- Q**-orthogonal, 286
- quadratic convergence rate, 296
- quadratic function, 40, 64, 77
- quadratic programming, 77, 156, 315
- quasi-convex function, 109
- quasi-Newton methods, 273, 293, 342
- Rademacher's Theorem, 283
- rank-two update, 293
- recession cone, 81, 82
- reduced cost, 227
- redundant constraint, 107
- relaxation, 19, 141, 142, 185, 328
- Relaxation Theorem, 141
- Representation Theorem, 50, 218, 308
- resistor, 182
- restricted master problem, 308
- restricted simplicial decomposition, 309
- restrification, 350

- revised simplex method, 238
- saddle point, 105, 147, 270
- scalar product, 34
- secant method, 274
- sensitivity analysis, 137, 178, 179
- sensitivity analysis for LP, 257
- separation of convex sets, 98
- Separation Theorem, 52, 98, 163
- sequential linear programming (SLP), 349
- sequential quadratic programming (SQP), 337–346
- set covering problem, 18
- shadow price, 249
- shortest route, 320
- simplex method, 10, 225–240
- simplicial decomposition algorithm, 308
- slack variable, 7
- Slater CQ, 131
- SLP algorithm, 349
- soft constraint, 18, 177
- spectral theorem, 136
- SQP algorithm, 337–350
- square matrix, 36
- stalling, 239
- standard basis, 35
- stationary point, 16, 85, 91
- steepest descent, 269
- steepest-edge rule, 230
- stochastic programming, 13
- strict inequality, 81
- strict local minimum, 76
- strictly convex function, 58
- strictly quasi-convex function, 277
- strong duality, 149
- Strong Duality Theorem, 149, 152–154, 156, 248
- subdifferentiability, 162
- subdifferential, 158
- subgradient, 158, 284
- subgradient optimization, 166
- subgradient projection method, 165
- superlinear convergence rate, 296
- symmetric matrix, 36
- tangent cone, 115
- traffic assignment problem, 318
- traffic equilibrium, 317
- triangle inequality, 36
- trust region methods, 284
- twice differentiable function, 38
- unconstrained optimization, 12, 84–88, 267–302
- unimodal function, 277
- unique optimum, 83
- upper level set, 151
- upper semi-continuity, 79
- user equilibrium, 317
- variable, 4
- variational inequality, 90, 104
- vector, 34
- vector-valued functions, 39
- voltage source, 182
- von Neumann’s Minimax Theorem, 105
- Wardrop’s principle, 317
- Weak Duality Theorem, 144, 247
- weak Wolfe condition, 278
- weakly coercive function, 78
- Weierstrass’ Theorem, 80, 98, 162
- Wolfe condition, 278

Errata and comments list for “An Introduction to Continuous Optimization”

Michael Patriksson
17 August, 2005

Page	Row	Reads	Should read
76	−2	has a lower value	has a lower function value
94	−13	than the other	than any of the other
98	11	simplicity	the readers’ convenience
165	17	means fast	means that fast
165	20	$\alpha_k = \gamma + \beta/(k + 1)$	$\alpha_k = \beta/(k + 1)$
165	21	where $\beta > 0$, $\gamma \geq 0$	where $\beta > 0$
171	Figure 6.4	A convex min-function	A concave min-function
175	14	$k \leq m + 1$ such that	$k \leq m + 1$, such that
276	Figure 11.2(b)	$\mathbf{x}_k + \alpha^* \mathbf{p}_k$	α^*
367	Exercise 10.5	$\mathbf{y} \leq \mathbf{0}^m$	$\mathbf{y} \geq \mathbf{0}^m$