

Convex Optimization Algorithms

Dimitri P. Bertsekas

Massachusetts Institute of Technology

WWW site for book information and orders

<http://www.athenasc.com>



Athena Scientific, Belmont, Massachusetts

Athena Scientific
Post Office Box 805
Nashua, NH 03061-0805
U.S.A.

Email: info@athenasc.com
WWW: <http://www.athenasc.com>

© 2015 Dimitri P. Bertsekas

All rights reserved. No part of this book may be reproduced in any form by any electronic or mechanical means (including photocopying, recording, or information storage and retrieval) without permission in writing from the publisher.

Publisher's Cataloging-in-Publication Data

Bertsekas, Dimitri P.

Convex Optimization Algorithms

Includes bibliographical references and index

1. Nonlinear Programming 2. Mathematical Optimization. I. Title.

T57.8.B475 2015 519.703

Library of Congress Control Number: 2002092168

ISBN-10: 1-886529-28-0, ISBN-13: 978-1-886529-28-1

Contents

1. Convex Optimization Models: An Overview	p. 1
1.1. Lagrange Duality	p. 2
1.1.1. Separable Problems – Decomposition	p. 7
1.1.2. Partitioning	p. 9
1.2. Fenchel Duality and Conic Programming	p. 10
1.2.1. Linear Conic Problems	p. 15
1.2.2. Second Order Cone Programming	p. 17
1.2.3. Semidefinite Programming	p. 22
1.3. Additive Cost Problems	p. 25
1.4. Large Number of Constraints	p. 34
1.5. Exact Penalty Functions	p. 39
1.6. Notes, Sources, and Exercises	p. 47
2. Optimization Algorithms: An Overview	p. 53
2.1. Iterative Descent Algorithms	p. 55
2.1.1. Differentiable Cost Function Descent – Unconstrained Problems	p. 58
2.1.2. Constrained Problems – Feasible Direction Methods	p. 71
2.1.3. Nondifferentiable Problems – Subgradient Methods	p. 78
2.1.4. Alternative Descent Methods	p. 80
2.1.5. Incremental Algorithms	p. 83
2.1.6. Distributed Asynchronous Iterative Algorithms	p. 104
2.2. Approximation Methods	p. 106
2.2.1. Polyhedral Approximation	p. 107
2.2.2. Penalty, Augmented Lagrangian, and Interior Point Methods	p. 108
2.2.3. Proximal Algorithm, Bundle Methods, and Tikhonov Regularization	p. 110
2.2.4. Alternating Direction Method of Multipliers	p. 111
2.2.5. Smoothing of Nondifferentiable Problems	p. 113
2.3. Notes, Sources, and Exercises	p. 119
3. Subgradient Methods	p. 135
3.1. Subgradients of Convex Real-Valued Functions	p. 136

3.1.1. Characterization of the Subdifferential	p. 146
3.2. Convergence Analysis of Subgradient Methods	p. 148
3.3. ϵ -Subgradient Methods	p. 162
3.3.1. Connection with Incremental Subgradient Methods	p. 166
3.4. Notes, Sources, and Exercises	p. 167
4. Polyhedral Approximation Methods	p. 181
4.1. Outer Linearization – Cutting Plane Methods	p. 182
4.2. Inner Linearization – Simplicial Decomposition	p. 188
4.3. Duality of Outer and Inner Linearization	p. 194
4.4. Generalized Polyhedral Approximation	p. 196
4.5. Generalized Simplicial Decomposition	p. 209
4.5.1. Differentiable Cost Case	p. 213
4.5.2. Nondifferentiable Cost and Side Constraints	p. 213
4.6. Polyhedral Approximation for Conic Programming	p. 217
4.7. Notes, Sources, and Exercises	p. 228
5. Proximal Algorithms	p. 233
5.1. Basic Theory of Proximal Algorithms	p. 234
5.1.1. Convergence	p. 235
5.1.2. Rate of Convergence	p. 239
5.1.3. Gradient Interpretation	p. 246
5.1.4. Fixed Point Interpretation, Overrelaxation,	
and Generalization	p. 248
5.2. Dual Proximal Algorithms	p. 256
5.2.1. Augmented Lagrangian Methods	p. 259
5.3. Proximal Algorithms with Linearization	p. 268
5.3.1. Proximal Cutting Plane Methods	p. 270
5.3.2. Bundle Methods	p. 272
5.3.3. Proximal Inner Linearization Methods	p. 276
5.4. Alternating Direction Methods of Multipliers	p. 280
5.4.1. Applications in Machine Learning	p. 286
5.4.2. ADMM Applied to Separable Problems	p. 289
5.5. Notes, Sources, and Exercises	p. 293
6. Additional Algorithmic Topics	p. 301
6.1. Gradient Projection Methods	p. 302
6.2. Gradient Projection with Extrapolation	p. 322
6.2.1. An Algorithm with Optimal Iteration Complexity	p. 323
6.2.2. Nondifferentiable Cost – Smoothing	p. 326
6.3. Proximal Gradient Methods	p. 330
6.4. Incremental Subgradient Proximal Methods	p. 340
6.4.1. Convergence for Methods with Cyclic Order	p. 344

6.4.2. Convergence for Methods with Randomized Order . . .	p. 353
6.4.3. Application in Specially Structured Problems	p. 361
6.4.4. Incremental Constraint Projection Methods	p. 365
6.5. Coordinate Descent Methods	p. 369
6.5.1. Variants of Coordinate Descent	p. 373
6.5.2. Distributed Asynchronous Coordinate Descent	p. 376
6.6. Generalized Proximal Methods	p. 382
6.7. ϵ -Descent and Extended Monotropic Programming	p. 396
6.7.1. ϵ -Subgradients	p. 397
6.7.2. ϵ -Descent Method	p. 400
6.7.3. Extended Monotropic Programming Duality	p. 406
6.7.4. Special Cases of Strong Duality	p. 408
6.8. Interior Point Methods	p. 412
6.8.1. Primal-Dual Methods for Linear Programming	p. 416
6.8.2. Interior Point Methods for Conic Programming	p. 423
6.8.3. Central Cutting Plane Methods	p. 425
6.9. Notes, Sources, and Exercises	p. 426
Appendix A: Mathematical Background	p. 443
A.1. Linear Algebra	p. 445
A.2. Topological Properties	p. 450
A.3. Derivatives	p. 456
A.4. Convergence Theorems	p. 458
Appendix B: Convex Optimization Theory: A Summary . . .	p. 467
B.1. Basic Concepts of Convex Analysis	p. 467
B.2. Basic Concepts of Polyhedral Convexity	p. 489
B.3. Basic Concepts of Convex Optimization	p. 494
B.4. Geometric Duality Framework	p. 498
B.5. Duality and Optimization	p. 505
References	p. 519
Index	p. 557

ATHENA SCIENTIFIC
OPTIMIZATION AND COMPUTATION SERIES

1. Convex Optimization Algorithms, by Dimitri P. Bertsekas, 2015, ISBN 978-1-886529-28-1, 576 pages
2. Abstract Dynamic Programming, by Dimitri P. Bertsekas, 2013, ISBN 978-1-886529-42-7, 256 pages
3. Dynamic Programming and Optimal Control, Two-Volume Set, by Dimitri P. Bertsekas, 2012, ISBN 1-886529-08-6, 1020 pages
4. Convex Optimization Theory, by Dimitri P. Bertsekas, 2009, ISBN 978-1-886529-31-1, 256 pages
5. Introduction to Probability, 2nd Edition, by Dimitri P. Bertsekas and John N. Tsitsiklis, 2008, ISBN 978-1-886529-23-6, 544 pages
6. Convex Analysis and Optimization, by Dimitri P. Bertsekas, Angelia Nedić, and Asuman E. Ozdaglar, 2003, ISBN 1-886529-45-0, 560 pages
7. Nonlinear Programming, 2nd Edition, by Dimitri P. Bertsekas, 1999, ISBN 1-886529-00-0, 791 pages
8. Network Optimization: Continuous and Discrete Models, by Dimitri P. Bertsekas, 1998, ISBN 1-886529-02-7, 608 pages
9. Network Flows and Monotropic Optimization, by R. Tyrrell Rockafellar, 1998, ISBN 1-886529-06-X, 634 pages
10. Introduction to Linear Optimization, by Dimitris Bertsimas and John N. Tsitsiklis, 1997, ISBN 1-886529-19-1, 608 pages
11. Parallel and Distributed Computation: Numerical Methods, by Dimitri P. Bertsekas and John N. Tsitsiklis, 1997, ISBN 1-886529-01-9, 718 pages
12. Neuro-Dynamic Programming, by Dimitri P. Bertsekas and John N. Tsitsiklis, 1996, ISBN 1-886529-10-8, 512 pages
13. Constrained Optimization and Lagrange Multiplier Methods, by Dimitri P. Bertsekas, 1996, ISBN 1-886529-04-3, 410 pages
14. Stochastic Optimal Control: The Discrete-Time Case, by Dimitri P. Bertsekas and Steven E. Shreve, 1996, ISBN 1-886529-03-5, 330 pages

ABOUT THE AUTHOR

Dimitri Bertsekas studied Mechanical and Electrical Engineering at the National Technical University of Athens, Greece, and obtained his Ph.D. in system science from the Massachusetts Institute of Technology. He has held faculty positions with the Engineering-Economic Systems Department, Stanford University, and the Electrical Engineering Department of the University of Illinois, Urbana. Since 1979 he has been teaching at the Electrical Engineering and Computer Science Department of the Massachusetts Institute of Technology (M.I.T.), where he is currently the McAfee Professor of Engineering.

His teaching and research spans several fields, including deterministic optimization, dynamic programming and stochastic control, large-scale and distributed computation, and data communication networks. He has authored or coauthored numerous research papers and sixteen books, several of which are currently used as textbooks in MIT classes, including “Nonlinear Programming,” “Dynamic Programming and Optimal Control,” “Data Networks,” “Introduction to Probability,” “Convex Optimization Theory,” as well as the present book. He often consults with private industry and has held editorial positions in several journals.

Professor Bertsekas was awarded the INFORMS 1997 Prize for Research Excellence in the Interface Between Operations Research and Computer Science for his book “Neuro-Dynamic Programming” (co-authored with John Tsitsiklis), the 2001 AACC John R. Ragazzini Education Award, the 2009 INFORMS Expository Writing Award, the 2014 AACC Richard Bellman Heritage Award for “contributions to the foundations of deterministic and stochastic optimization-based methods in systems and control,” the 2014 Khachiyan Prize for “life-time accomplishments in optimization,” and the SIAM/MOS 2015 George B. Dantzig Prize for “original research, which by its originality, breadth, and scope, is having a major impact on the field of mathematical programming.” In 2001, he was elected to the United States National Academy of Engineering for “pioneering contributions to fundamental research, practice and education of optimization/control theory, and especially its application to data communication networks.”

Preface

**There is no royal way to geometry
(Euclid to king Ptolemy of Alexandria)**

Interest in convex optimization has become intense due to widespread applications in fields such as large-scale resource allocation, signal processing, and machine learning. This book aims at an up-to-date and accessible development of algorithms for solving convex optimization problems.

The book complements the author's 2009 "Convex Optimization Theory" book, but can be read independently. The latter book focuses on convexity theory and optimization duality, while the present book focuses on algorithmic issues. The two books share mathematical prerequisites, notation, and style, and together cover the entire finite-dimensional convex optimization field. Both books rely on rigorous mathematical analysis, but also aim at an intuitive exposition that makes use of visualization where possible. This is facilitated by the extensive use of analytical and algorithmic concepts of duality, which by nature lend themselves to geometrical interpretation.

To enhance readability, the statements of definitions and results of the "theory book" are reproduced without proofs in Appendix B. Moreover, some of the theory needed for the present book, has been replicated and/or adapted to its algorithmic nature. For example the theory of subgradients for real-valued convex functions is fully developed in Chapter 3. Thus the reader who is already familiar with the analytical foundations of convex optimization need not consult the "theory book" except for the purpose of studying the proofs of some specific results.

The book covers almost all the major classes of convex optimization algorithms. Principal among these are gradient, subgradient, polyhedral approximation, proximal, and interior point methods. Most of these methods rely on convexity (but not necessarily differentiability) in the cost and constraint functions, and are often connected in various ways to duality. I have provided numerous examples describing in detail applications to specially structured problems. The reader may also find a wealth of analysis and discussion of applications in books on large-scale convex optimization, network optimization, parallel and distributed computation, signal processing, and machine learning.

The chapter-by-chapter description of the book follows:

Chapter 1: Here we provide a broad overview of some important classes of convex optimization problems, and their principal characteristics. Several

problem structures are discussed, often arising from Lagrange duality theory and Fenchel duality theory, together with its special case, conic duality. Some additional structures involving a large number of additive terms in the cost, or a large number of constraints are also discussed, together with their applications in machine learning and large-scale resource allocation.

Chapter 2: Here we provide an overview of algorithmic approaches, focusing primarily on algorithms for differentiable optimization, and we discuss their differences from their nondifferentiable convex optimization counterparts. We also highlight the main ideas of the two principal algorithmic approaches of this book, iterative descent and approximation, and we illustrate their application with specific algorithms, reserving detailed analysis for subsequent chapters.

Chapter 3: Here we discuss subgradient methods for minimizing a convex cost function over a convex constraint set. The cost function may be nondifferentiable, as is often the case in the context of duality and machine learning applications. These methods are based on the idea of reduction of distance to the optimal set, and include variations aimed at algorithmic efficiency, such as ϵ -subgradient and incremental subgradient methods.

Chapter 4: Here we discuss polyhedral approximation methods for minimizing a convex function over a convex constraint set. The two main approaches here are outer linearization (also called the cutting plane approach) and inner linearization (also called the simplicial decomposition approach). We show how these two approaches are intimately connected by conjugacy and duality, and we generalize our framework for polyhedral approximation to the case where the cost function is a sum of two or more convex component functions.

Chapter 5: Here we focus on proximal algorithms for minimizing a convex function over a convex constraint set. At each iteration of the basic proximal method, we solve an approximation to the original problem. However, unlike the preceding chapter, the approximation is not polyhedral, but rather it is based on quadratic regularization, i.e., adding a quadratic term to the cost function, which is appropriately adjusted at each iteration. We discuss several variations of the basic algorithm. Some of these include combinations with the polyhedral approximation methods of the preceding chapter, yielding the class of bundle methods. Others are obtained via duality from the basic proximal algorithm, including the augmented Lagrangian method (also called method of multipliers) for constrained optimization. Finally, we discuss extensions of the proximal algorithm for finding a zero of a maximal monotone operator, and a major special case: the alternating direction method of multipliers, which is well suited for taking advantage of the structure of several types of large-scale problems.

Chapter 6: Here we discuss a variety of algorithmic topics that supplement our discussion of the descent and approximation methods of the

preceding chapters. We first discuss gradient projection methods and variations with extrapolation that have good complexity properties, including Nesterov's optimal complexity algorithm. These were developed for differentiable problems, and can be extended to the nondifferentiable case by means of a smoothing scheme. Then we discuss a number of combinations of gradient, subgradient, and proximal methods that are well suited for specially structured problems. We pay special attention to incremental versions for the case where the cost function consists of the sum of a large number of component terms. We also describe additional methods, such as the classical block coordinate descent approach, the proximal algorithm with a nonquadratic regularization term, and the ϵ -descent method. We close the chapter with a discussion of interior point methods.

Our lines of analysis are largely based on differential calculus-type ideas, which are central in nonlinear programming, and on concepts of hyperplane separation, conjugacy, and duality, which are central in convex analysis. A traditional use of duality is to establish the equivalence and the connections between a pair of primal and dual problems, which may in turn enhance insight and enlarge the set of options for analysis and computation. The book makes heavy use of this type of problem duality, but also emphasizes a qualitatively different, algorithm-oriented type of duality that is largely based on conjugacy. In particular, some fundamental algorithmic operations turn out to be dual to each other, and whenever they arise in various algorithms they admit dual implementations, often with significant gains in insight and computational convenience. Some important examples are the duality between the subdifferentials of a convex function and its conjugate, the duality of a proximal operation using a convex function and an augmented Lagrangian minimization using its conjugate, and the duality between outer linearization of a convex function and inner linearization of its conjugate. Several interesting algorithms in Chapters 4-6 admit dual implementations based on these pairs of operations.

The book contains a fair number of exercises, many of them supplementing the algorithmic development and analysis. In addition a large number of theoretical exercises (with carefully written solutions) for the "theory book," together with other related material, can be obtained from the book's web page <http://www.athenasc.com/convexalgorithms.html>, and the author's web page <http://web.mit.edu/dimitrib/www/home.html>. The MIT OpenCourseWare site <http://ocw.mit.edu/index.htm>, also provides lecture slides and other relevant material.

The mathematical prerequisites for the book are a first course in linear algebra and a first course in real analysis. A summary of the relevant material is provided in Appendix A. Prior exposure to linear and nonlinear optimization algorithms is not assumed, although it will undoubtedly be helpful in providing context and perspective. Other than this background, the development is self-contained, with proofs provided throughout.

The present book, in conjunction with its “theory” counterpart may be used as a text for a one-semester or two-quarter convex optimization course; I have taught several variants of such a course at MIT and elsewhere over the last fifteen years. Still the book may not provide all of the convex optimization material an instructor may wish for, and it may need to be supplemented by works that aim primarily at specific types of convex optimization models, or address more comprehensively computational complexity issues. I have added representative citations for such works, which, however, are far from complete in view of the explosive growth of the literature on the subject.

The book may also be used as a supplementary source for nonlinear programming classes that are primarily focused on classical differentiable nonconvex optimization material (Kuhn-Tucker theory, Newton-like and conjugate direction methods, interior point, penalty, and augmented Lagrangian methods). For such courses, it may provide a nondifferentiable convex optimization component.

I was fortunate to have several outstanding collaborators in my research on various aspects of convex optimization: Vivek Borkar, Jon Eckstein, Eli Gafni, Xavier Luque, Angelia Nedić, Asuman Ozdaglar, John Tsitsiklis, Mengdi Wang, and Huizhen (Janey) Yu. Substantial portions of our joint research have found their way into the book. In addition, I am grateful for interactions and suggestions I received from several colleagues, including Leon Bottou, Steve Boyd, Tom Luo, Steve Wright, and particularly Mark Schmidt and Lin Xiao who read with care major portions of the book. I am also very thankful for the valuable proofreading of parts of the book by Mengdi Wang and Huizhen (Janey) Yu, and particularly by Ivan Pejcic who went through most of the book with a keen eye. I developed the book through convex optimization classes at MIT over a fifteen-year period, and I want to express appreciation for my students who provided continuing motivation and inspiration.

Finally, I would like to mention Paul Tseng, a major contributor to numerous topics in this book, who was my close friend and research collaborator on optimization algorithms for many years, and whom we unfortunately lost while he was still at his prime. I am dedicating the book to his memory.

Dimitri P. Bertsekas
dimitrib@mit.edu
January 2015

1

Convex Optimization Models: An Overview

Contents

1.1. Lagrange Duality	p. 2
1.1.1. Separable Problems – Decomposition	p. 7
1.1.2. Partitioning	p. 9
1.2. Fenchel Duality and Conic Programming	p. 10
1.2.1. Linear Conic Problems	p. 15
1.2.2. Second Order Cone Programming	p. 17
1.2.3. Semidefinite Programming	p. 22
1.3. Additive Cost Problems	p. 25
1.4. Large Number of Constraints	p. 34
1.5. Exact Penalty Functions	p. 39
1.6. Notes, Sources, and Exercises	p. 47

In this chapter we provide an overview of some broad classes of convex optimization models. Our primary focus will be on large challenging problems, often connected in some way to duality. We will consider two types of duality. The first is *Lagrange duality* for constrained optimization, which is obtained by assigning dual variables to the constraints. The second is *Fenchel duality* together with its special case, conic duality, which involves a cost function that is the sum of two convex function components. Both of these duality structures arise often in applications, and in Sections 1.1 and 1.2 we provide an overview, and discuss some examples.[†]

In Sections 1.3 and 1.4, we discuss additional model structures involving a large number of additive terms in the cost, or a large number of constraints. These types of problems also arise often in the context of duality, as well as in other contexts such as machine learning and signal processing with large amounts of data. In Section 1.5, we discuss the exact penalty function technique, whereby we can transform a convex constrained optimization problem to an equivalent unconstrained problem.

1.1 LAGRANGE DUALITY

We start our overview of Lagrange duality with the basic case of nonlinear inequality constraints, and then consider extensions involving linear inequality and equality constraints. Consider the problem[‡]

$$\begin{aligned} & \text{minimize} && f(x) \\ & \text{subject to} && x \in X, \quad g(x) \leq 0, \end{aligned} \tag{1.1}$$

where X is a nonempty set,

$$g(x) = (g_1(x), \dots, g_r(x))',$$

and $f : X \mapsto \Re$ and $g_j : X \mapsto \Re$, $j = 1, \dots, r$, are given functions. We refer to this as the *primal problem*, and we denote its optimal value by f^* . A vector x satisfying the constraints of the problem is referred to as *feasible*. The *dual* of problem (1.1) is given by

$$\begin{aligned} & \text{maximize} && q(\mu) \\ & \text{subject to} && \mu \in \Re^r, \end{aligned} \tag{1.2}$$

[†] Consistent with its overview character, this chapter contains few proofs, and refers frequently to the literature, and to Appendix B, which contains a full list of definitions and propositions (without proofs) relating to nonalgorithmic aspects of convex optimization. This list reflects and summarizes the content of the author's "Convex Optimization Theory" book [Ber09]. The proposition numbers of [Ber09] have been preserved, so all omitted proofs of propositions in Appendix B can be readily accessed from [Ber09].

[‡] Appendix A contains an overview of the mathematical notation, terminology, and results from linear algebra and real analysis that we will be using.

where the dual function q is

$$q(\mu) = \begin{cases} \inf_{x \in X} L(x, \mu) & \text{if } \mu \geq 0, \\ -\infty & \text{otherwise,} \end{cases}$$

and L is the Lagrangian function defined by

$$L(x, \mu) = f(x) + \mu'g(x), \quad x \in X, \mu \in \mathbb{R}^r;$$

(cf. Section 5.3 of Appendix B).

Note that the dual function is extended real-valued, and that the effective constraint set of the dual problem is

$$\left\{ \mu \geq 0 \mid \inf_{x \in X} L(x, \mu) > -\infty \right\}.$$

The optimal value of the dual problem is denoted by q^* .

The *weak duality* relation, $q^* \leq f^*$, always holds. It is easily shown by writing for all $\mu \geq 0$, and $x \in X$ with $g(x) \leq 0$,

$$q(\mu) = \inf_{z \in X} L(z, \mu) \leq L(x, \mu) = f(x) + \sum_{j=1}^r \mu_j g_j(x) \leq f(x),$$

so that

$$q^* = \sup_{\mu \in \mathbb{R}^r} q(\mu) = \sup_{\mu \geq 0} q(\mu) \leq \inf_{x \in X, g(x) \leq 0} f(x) = f^*.$$

We state this formally as follows (cf. Prop. 4.1.2 in Appendix B).

Proposition 1.1.1: (Weak Duality Theorem) Consider problem (1.1). For any feasible solution x and any $\mu \in \mathbb{R}^r$, we have $q(\mu) \leq f(x)$. Moreover, $q^* \leq f^*$.

When $q^* = f^*$, we say that *strong duality* holds. The following proposition gives necessary and sufficient conditions for strong duality, and primal and dual optimality (see Prop. 5.3.2 in Appendix B).

Proposition 1.1.2: (Optimality Conditions) Consider problem (1.1). There holds $q^* = f^*$, and (x^*, μ^*) are a primal and dual optimal solution pair if and only if x^* is feasible, $\mu^* \geq 0$, and

$$x^* \in \arg \min_{x \in X} L(x, \mu^*), \quad \mu_j^* g_j(x^*) = 0, \quad j = 1, \dots, r.$$

Both of the preceding propositions do not require any convexity assumptions on f , g , and X . However, generally the analytical and algorithmic solution process is simplified when strong duality ($q^* = f^*$) holds. This typically requires convexity assumptions, and in some cases conditions on $\text{ri}(X)$, the relative interior of X , as exemplified by the following result, given in Prop. 5.3.1 in Appendix B. The result delineates the two principal cases where there is no duality gap in an inequality-constrained problem.

Proposition 1.1.3: (Strong Duality – Existence of Dual Optimal Solutions) Consider problem (1.1) under the assumption that the set X is convex, and the functions f , and g_1, \dots, g_r are convex. Assume further that f^* is finite, and that one of the following two conditions holds:

- (1) There exists $\bar{x} \in X$ such that $g_j(\bar{x}) < 0$ for all $j = 1, \dots, r$.
- (2) The functions g_j , $j = 1, \dots, r$, are affine, and there exists $\bar{x} \in \text{ri}(X)$ such that $g(\bar{x}) \leq 0$.

Then $q^* = f^*$ and there exists at least one dual optimal solution. Under condition (1) the set of dual optimal solutions is also compact.

Convex Programming with Inequality and Equality Constraints

Let us consider an extension of problem (1.1), with additional linear equality constraints. It is our principal constrained optimization model under convexity assumptions, and it will be referred to as the *convex programming problem*. It is given by

$$\begin{aligned} & \text{minimize} && f(x) \\ & \text{subject to} && x \in X, \quad g(x) \leq 0, \quad Ax = b, \end{aligned} \tag{1.3}$$

where X is a convex set, $g(x) = (g_1(x), \dots, g_r(x))'$, $f : X \mapsto \mathbb{R}$ and $g_j : X \mapsto \mathbb{R}$, $j = 1, \dots, r$, are given convex functions, A is an $m \times n$ matrix, and $b \in \mathbb{R}^m$.

The preceding duality framework may be applied to this problem by converting the constraint $Ax = b$ to the equivalent set of linear inequality constraints

$$Ax \leq b, \quad -Ax \leq -b,$$

with corresponding dual variables $\lambda^+ \geq 0$ and $\lambda^- \geq 0$. The Lagrangian function is

$$f(x) + \mu'g(x) + (\lambda^+ - \lambda^-)'(Ax - b),$$

and by introducing a dual variable

$$\lambda = \lambda^+ - \lambda^-$$

with no sign restriction, it can be written as

$$L(x, \mu, \lambda) = f(x) + \mu'g(x) + \lambda'(Ax - b).$$

The dual problem is

$$\begin{aligned} & \text{maximize} && \inf_{x \in X} L(x, \mu, \lambda) \\ & \text{subject to} && \mu \geq 0, \lambda \in \mathbb{R}^m. \end{aligned}$$

In this manner, Prop. 1.1.3 under condition (2), together with Prop. 1.1.2, yield the following for the case where all constraint functions are linear.

Proposition 1.1.4: (Convex Programming – Linear Equality and Inequality Constraints) Consider problem (1.3).

- (a) Assume that f^* is finite, that the functions g_j are affine, and that there exists $\bar{x} \in \text{ri}(X)$ such that $A\bar{x} = b$ and $g(\bar{x}) \leq 0$. Then $q^* = f^*$ and there exists at least one dual optimal solution.
- (b) There holds $f^* = q^*$, and (x^*, μ^*, λ^*) are a primal and dual optimal solution pair if and only if x^* is feasible, $\mu^* \geq 0$, and

$$x^* \in \arg \min_{x \in X} L(x, \mu^*, \lambda^*), \quad \mu_j^* g_j(x^*) = 0, \quad j = 1, \dots, r.$$

In the special case where there are no inequality constraints:

$$\begin{aligned} & \text{minimize} && f(x) \\ & \text{subject to} && x \in X, \quad Ax = b, \end{aligned} \tag{1.4}$$

the Lagrangian function is

$$L(x, \lambda) = f(x) + \lambda'(Ax - b),$$

and the dual problem is

$$\begin{aligned} & \text{maximize} && \inf_{x \in X} L(x, \lambda) \\ & \text{subject to} && \lambda \in \mathbb{R}^m. \end{aligned}$$

The corresponding result, a simpler special case of Prop. 1.1.4, is given in the following proposition.

Proposition 1.1.5: (Convex Programming – Linear Equality Constraints) Consider problem (1.4).

- (a) Assume that f^* is finite and that there exists $\bar{x} \in \text{ri}(X)$ such that $A\bar{x} = b$. Then $f^* = q^*$ and there exists at least one dual optimal solution.
- (b) There holds $f^* = q^*$, and (x^*, λ^*) are a primal and dual optimal solution pair if and only if x^* is feasible and

$$x^* \in \arg \min_{x \in X} L(x, \lambda^*).$$

The following is an extension of Prop. 1.1.4(a) to the case where the inequality constraints may be nonlinear. It is the most general convex programming result relating to duality in this section (see Prop. 5.3.5 in Appendix B).

Proposition 1.1.6: (Convex Programming – Linear Equality and Nonlinear Inequality Constraints) Consider problem (1.3).

Assume that f^* is finite, that there exists $\bar{x} \in X$ such that $A\bar{x} = b$ and $g(\bar{x}) < 0$, and that there exists $\tilde{x} \in \text{ri}(X)$ such that $A\tilde{x} = b$. Then $q^* = f^*$ and there exists at least one dual optimal solution.

Aside from the preceding results, there are alternative optimality conditions for convex and nonconvex optimization problems, which are based on extended versions of the Fritz John theorem; see [Be002] and [BOT06], and the textbooks [Ber99] and [BNO03]. These conditions are derived using a somewhat different line of analysis and supplement the ones given here, but we will not have occasion to use them in this book.

Discrete Optimization and Lower Bounds

The preceding propositions deal mostly with situations where strong duality holds ($q^* = f^*$). However, duality can be useful even when there is duality gap, as often occurs in problems that have a finite constraint set X . An example is *integer programming*, where the components of x must be integers from a bounded range (usually 0 or 1). An important special case is the linear 0-1 integer programming problem

$$\begin{aligned} & \text{minimize} && c'x \\ & \text{subject to} && Ax \leq b, \quad x_i = 0 \text{ or } 1, \quad i = 1, \dots, n, \end{aligned}$$

where $x = (x_1, \dots, x_n)$.

A principal approach for solving discrete optimization problems with a finite constraint set is the *branch-and-bound method*, which is described in many sources; see e.g., one of the original works [LaD60], the survey [BaT85], and the book [NeW88]. The general idea of the method is that bounds on the cost function can be used to exclude from consideration portions of the feasible set. To illustrate, consider minimizing $F(x)$ over $x \in X$, and let Y_1, Y_2 be two subsets of X . Suppose that we have bounds

$$\underline{F}_1 \leq \min_{x \in Y_1} f(x), \quad \overline{F}_2 \geq \min_{x \in Y_2} f(x).$$

Then, if $\overline{F}_2 \leq \underline{F}_1$, the solutions in Y_1 may be disregarded since their cost cannot be smaller than the cost of the best solution in Y_2 . The lower bound \underline{F}_1 can often be conveniently obtained by minimizing f over a suitably enlarged version of Y_1 , while for the upper bound \overline{F}_2 , a value $f(x)$, where $x \in Y_2$, may be used.

Branch-and-bound is often based on weak duality (cf. Prop. 1.1.1) to obtain lower bounds to the optimal cost of restricted problems of the form

$$\begin{aligned} & \text{minimize} && f(x) \\ & \text{subject to} && x \in \tilde{X}, \quad g(x) \leq 0, \end{aligned} \tag{1.5}$$

where \tilde{X} is a subset of X ; for example in the 0-1 integer case where X specifies that all x_i should be 0 or 1, \tilde{X} may be the set of all 0-1 vectors x such that one or more components x_i are fixed at either 0 or 1 (i.e., are restricted to satisfy $x_i = 0$ for all $x \in \tilde{X}$ or $x_i = 1$ for all $x \in \tilde{X}$). These lower bounds can often be obtained by finding a dual-feasible (possibly dual-optimal) solution $\mu \geq 0$ of this problem and the corresponding dual value

$$q(\mu) = \inf_{x \in \tilde{X}} \{f(x) + \mu'g(x)\}, \tag{1.6}$$

which by weak duality, is a lower bound to the optimal value of the restricted problem (1.5). In a strengthened version of this approach, the given inequality constraints $g(x) \leq 0$ may be augmented by additional inequalities that are known to be satisfied by optimal solutions of the original problem.

An important point here is that when \tilde{X} is finite, the dual function q of Eq. (1.6) is concave and polyhedral. Thus solving the dual problem amounts to minimizing the polyhedral function $-q$ over the nonnegative orthant. This is a major context within which polyhedral functions arise in convex optimization.

1.1.1 Separable Problems – Decomposition

Let us now discuss an important problem structure that involves Lagrange duality and arises frequently in applications. Here x has m components,

$x = (x_1, \dots, x_m)$, with each x_i being a vector of dimension n_i (often $n_i = 1$). The problem has the form

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^m f_i(x_i) \\ & \text{subject to} && \sum_{i=1}^m g_{ij}(x_i) \leq 0, \quad x_i \in X_i, \quad i = 1, \dots, m, \quad j = 1, \dots, r, \end{aligned} \tag{1.7}$$

where $f_i : \mathbb{R}^{n_i} \mapsto \mathbb{R}$ and $g_{ij} : \mathbb{R}^{n_i} \mapsto \mathbb{R}^r$ are given functions, and X_i are given subsets of \mathbb{R}^{n_i} . By assigning a dual variable μ_j to the j th constraint, we obtain the dual problem [cf. Eq. (1.2)]

$$\begin{aligned} & \text{maximize} && \sum_{i=1}^m q_i(\mu) \\ & \text{subject to} && \mu \geq 0, \end{aligned} \tag{1.8}$$

where

$$q_i(\mu) = \inf_{x_i \in X_i} \left\{ f_i(x_i) + \sum_{j=1}^r \mu_j g_{ij}(x_i) \right\},$$

and $\mu = (\mu_1, \dots, \mu_r)$.

Note that the minimization involved in the calculation of the dual function has been decomposed into m simpler minimizations. These minimizations are often conveniently done either analytically or computationally, in which case the dual function can be easily evaluated. This is the key advantageous structure of separable problems: it facilitates computation of dual function values (as well as subgradients as we will see in Section 3.1), and it is amenable to decomposition and distributed computation.

Let us also note that in the special case where the components x_i are one-dimensional, and the functions f_i and sets X_i are convex, there is a particularly favorable duality result for the separable problem (1.7): essentially, strong duality holds without any qualifications such as the linearity of the constraint functions, or the Slater condition of Prop. 1.1.3; see [Tse09].

Duality Gap Estimates for Nonconvex Separable Problems

The separable structure is additionally helpful when the cost and/or the constraints are not convex, and there is a duality gap. In particular, in this case *the duality gap turns out to be relatively small and can often be shown to diminish to zero relative to the optimal primal value as the number m of separable terms increases*. As a result, one can often obtain a near-optimal primal solution, starting from a dual-optimal solution, without resorting to costly branch-and-bound procedures.

The small duality gap size is a consequence of the structure of the set S of constraint-cost pairs of problem (1.7), which in the case of a separable problem, can be written as a vector sum of m sets, one for each separable term, i.e.,

$$S = S_1 + \cdots + S_m,$$

where

$$S_i = \{(g_i(x_i), f_i(x_i)) \mid x_i \in X_i\},$$

and $g_i : \mathbb{R}^{n_i} \mapsto \mathbb{R}^r$ is the function $g_i(x_i) = (g_{i1}(x_i), \dots, g_{im}(x_i))$. It can be shown that the duality gap is related to how much S “differs” from its convex hull (a geometric explanation is given in [Ber99], Section 5.1.6, and [Ber09], Section 5.7). Generally, a set that is the vector sum of a large number of possibly nonconvex but roughly similar sets “tends to be convex” in the sense that any vector in its convex hull can be closely approximated by a vector in the set. As a result, the duality gap tends to be relatively small. The analytical substantiation is based on a theorem by Shapley and Folkman (see [Ber99], Section 5.1, or [Ber09], Prop. 5.7.1, for a statement and proof of this theorem). In particular, it is shown in [AuE76], and also [BeS82], [Ber82a], Section 5.6.1, under various reasonable assumptions, that the duality gap satisfies

$$f^* - q^* \leq (r + 1) \max_{i=1, \dots, m} \rho_i,$$

where for each i , ρ_i is a nonnegative scalar that depends on the structure of the functions $f_i, g_{ij}, j = 1, \dots, r$, and the set X_i (the paper [AuE76] focuses on the case where the problem is nonconvex but continuous, while [BeS82] and [Ber82a] focus on an important class of mixed integer programming problems). This estimate suggests that as $m \rightarrow \infty$ and $|f^*| \rightarrow \infty$, the duality gap is bounded, while the “relative” duality gap $(f^* - q^*)/|f^*|$ diminishes to 0 as $m \rightarrow \infty$.

The duality gap has also been investigated in the author’s book [Ber09] within the more general min common-max crossing framework (Section 4.1 of Appendix B). This framework includes as special cases minimax and zero-sum game problems. In particular, consider a function $\phi : X \times Z \mapsto \mathbb{R}$ defined over nonempty subsets $X \subset \mathbb{R}^n$ and $Z \subset \mathbb{R}^m$. Then it can be shown that the gap between “infsup” and “supinf” of ϕ can be decomposed into the sum of two terms that can be computed separately: one term can be attributed to the lack of convexity and/or closure of ϕ with respect to x , and the other can be attributed to the lack of concavity and/or upper semicontinuity of ϕ with respect to z . We refer to [Ber09], Section 5.7.2, for the analysis.

1.1.2 Partitioning

It is important to note that there are several different ways to introduce duality in the solution of large-scale optimization problems. For example a

strategy, often called *partitioning*, is to divide the variables in two subsets, and minimize first with respect to one subset while taking advantage of whatever simplification may arise by fixing the variables in the other subset.

As an example, the problem

$$\begin{aligned} & \text{minimize} && F(x) + G(y) \\ & \text{subject to} && Ax + By = c, \quad x \in X, \quad y \in Y, \end{aligned}$$

can be written as

$$\begin{aligned} & \text{minimize} && F(x) + \inf_{By=c-Ax, y \in Y} G(y) \\ & \text{subject to} && x \in X, \end{aligned}$$

or

$$\begin{aligned} & \text{minimize} && F(x) + p(c - Ax) \\ & \text{subject to} && x \in X, \end{aligned}$$

where p is given by

$$p(u) = \inf_{By=u, y \in Y} G(y).$$

In favorable cases, p can be dealt with conveniently (see e.g., the book [Las70] and the paper [Geo72]).

Strategies of splitting or transforming the variables to facilitate algorithmic solution will be frequently encountered in what follows, and in a variety of contexts, including duality. The next section describes some significant contexts of this type.

1.2 FENCHEL DUALITY AND CONIC PROGRAMMING

Let us consider the Fenchel duality framework (see Section 5.3.5 of Appendix B). It involves the problem

$$\begin{aligned} & \text{minimize} && f_1(x) + f_2(Ax) \\ & \text{subject to} && x \in \mathbb{R}^n, \end{aligned} \tag{1.9}$$

where A is an $m \times n$ matrix, $f_1 : \mathbb{R}^n \mapsto (-\infty, \infty]$ and $f_2 : \mathbb{R}^m \mapsto (-\infty, \infty]$ are closed proper convex functions, and we assume that there exists a feasible solution, i.e., an $x \in \mathbb{R}^n$ such that $x \in \text{dom}(f_1)$ and $Ax \in \text{dom}(f_2)$.[†]

The problem is equivalent to the following constrained optimization problem in the variables $x_1 \in \mathbb{R}^n$ and $x_2 \in \mathbb{R}^m$:

$$\begin{aligned} & \text{minimize} && f_1(x_1) + f_2(x_2) \\ & \text{subject to} && x_1 \in \text{dom}(f_1), \quad x_2 \in \text{dom}(f_2), \quad x_2 = Ax_1. \end{aligned} \tag{1.10}$$

[†] We remind the reader that our convex analysis notation, terminology, and nonalgorithmic theory are summarized in Appendix B.

Viewing this as a convex programming problem with the linear equality constraint $x_2 = Ax_1$, we obtain the dual function as

$$\begin{aligned} q(\lambda) &= \inf_{x_1 \in \text{dom}(f_1), x_2 \in \text{dom}(f_2)} \{f_1(x_1) + f_2(x_2) + \lambda'(x_2 - Ax_1)\} \\ &= \inf_{x_1 \in \mathbb{R}^n} \{f_1(x_1) - \lambda'Ax_1\} + \inf_{x_2 \in \mathbb{R}^m} \{f_2(x_2) + \lambda'x_2\}. \end{aligned}$$

The dual problem of maximizing q over $\lambda \in \mathbb{R}^m$, after a sign change to convert it to a minimization problem, takes the form

$$\begin{aligned} &\text{minimize} && f_1^*(A'\lambda) + f_2^*(-\lambda) \\ &\text{subject to} && \lambda \in \mathbb{R}^m, \end{aligned} \tag{1.11}$$

where f_1^* and f_2^* are the conjugate functions of f_1 and f_2 . We denote by f^* and q^* the corresponding optimal primal and dual values [q^* is the negative of the optimal value of problem (1.11)].

The following Fenchel duality result is given as Prop. 5.3.8 in Appendix B. Parts (a) and (b) are obtained by applying Prop. 1.1.5(a) to problem (1.10), viewed as a problem with $x_2 = Ax_1$ as the only linear equality constraint. The first equation of part (c) is a consequence of Prop. 1.1.5(b). Its equivalence with the last two equations is a consequence of the Conjugate Subgradient Theorem (Prop. 5.4.3, App. B), which states that for a closed proper convex function f , its conjugate f^* , and any pair of vectors (x, y) , we have

$$x \in \arg \min_{z \in \mathbb{R}^n} \{f(z) - z'y\} \quad \text{iff} \quad y \in \partial f(x) \quad \text{iff} \quad x \in \partial f^*(y),$$

with all of these three relations being equivalent to $x'y = f(x) + f^*(y)$. Here $\partial f(x)$ denotes the subdifferential of f at x (the set of all subgradients of f at x); see Section 5.4 of Appendix B.

Proposition 1.2.1: (Fenchel Duality) Consider problem (1.9).

- (a) If f^* is finite and $(A \cdot \text{ri}(\text{dom}(f_1))) \cap \text{ri}(\text{dom}(f_2)) \neq \emptyset$, then $f^* = q^*$ and there exists at least one dual optimal solution.
- (b) If q^* is finite and $\text{ri}(\text{dom}(f_1^*)) \cap (A' \cdot \text{ri}(-\text{dom}(f_2^*))) \neq \emptyset$, then $f^* = q^*$ and there exists at least one primal optimal solution.
- (c) There holds $f^* = q^*$, and (x^*, λ^*) is a primal and dual optimal solution pair if and only if any one of the following three equivalent conditions hold:

$$x^* \in \arg \min_{x \in \mathbb{R}^n} \{f_1(x) - x'A'\lambda^*\} \quad \text{and} \quad Ax^* \in \arg \min_{z \in \mathbb{R}^m} \{f_2(z) + z'\lambda^*\}, \tag{1.12}$$

$$A'\lambda^* \in \partial f_1(x^*) \quad \text{and} \quad -\lambda^* \in \partial f_2(Ax^*), \tag{1.13}$$

$$x^* \in \partial f_1^*(A'\lambda^*) \quad \text{and} \quad Ax^* \in \partial f_2^*(-\lambda^*). \tag{1.14}$$

Minimax Problems

Minimax problems involve minimization over a set X of a function \overline{F} of the form

$$\overline{F}(x) = \sup_{z \in Z} \phi(x, z),$$

where X and Z are subsets of \Re^n and \Re^m , respectively, and $\phi : \Re^n \times \Re^m \mapsto \Re$ is a given function. Some (but not all) problems of this type are related to constrained optimization and Fenchel duality.

Example 1.2.1: (Connection with Constrained Optimization)

Let ϕ and Z have the form

$$\phi(x, z) = f(x) + z'g(x), \quad Z = \{z \mid z \geq 0\},$$

where $f : \Re^n \mapsto \Re$ and $g : \Re^n \mapsto \Re^m$ are given functions. Then it is seen that

$$\overline{F}(x) = \sup_{z \in Z} \phi(x, z) = \begin{cases} f(x) & \text{if } g(x) \leq 0, \\ \infty & \text{otherwise.} \end{cases}$$

Thus minimization of \overline{F} over $x \in X$ is equivalent to solving the constrained optimization problem

$$\begin{aligned} & \text{minimize} && f(x) \\ & \text{subject to} && x \in X, \quad g(x) \leq 0. \end{aligned} \tag{1.15}$$

The dual problem is to maximize over $z \geq 0$ the function

$$\underline{F}(z) = \inf_{x \in X} \{f(x) + z'g(x)\} = \inf_{x \in X} \phi(x, z),$$

and the minimax equality

$$\inf_{x \in X} \sup_{z \in Z} \phi(x, z) = \sup_{z \in Z} \inf_{x \in X} \phi(x, z) \tag{1.16}$$

is equivalent to problem (1.15) having no duality gap.

Example 1.2.2: (Connection with Fenchel Duality)

Let ϕ have the special form

$$\phi(x, z) = f(x) + z'Ax - g(z),$$

where $f : \Re^n \mapsto \Re$ and $g : \Re^m \mapsto \Re$ are given functions, and A is a given $m \times n$ matrix. Then we have

$$\overline{F}(x) = \sup_{z \in Z} \phi(x, z) = f(x) + \sup_{z \in Z} \{(Ax)'z - g(z)\} = f(x) + \hat{g}^*(Ax),$$

where \hat{g}^* is the conjugate of the function

$$\hat{g}(z) = \begin{cases} g(z) & \text{if } z \in Z, \\ \infty & \text{otherwise.} \end{cases}$$

Thus the minimax problem of minimizing \overline{F} over $x \in X$ comes under the Fenchel framework (1.9) with $f_2 = \hat{g}^*$ and f_1 given by

$$f_1(x) = \begin{cases} f(x) & \text{if } x \in X, \\ \infty & \text{if } x \notin X. \end{cases}$$

It can also be verified that the Fenchel dual problem (1.11) is equivalent to maximizing over $z \in Z$ the function $\underline{F}(z) = \inf_{x \in X} \phi(x, z)$. Again having no duality gap is equivalent to the minimax equality (1.16) holding.

Finally note that strong duality theory is connected with minimax problems primarily when X and Z are convex sets, and ϕ is convex in x and concave in z . When Z is a finite set, there is a different connection with constrained optimization that does not involve Fenchel duality and applies without any convexity conditions. In particular, the problem

$$\begin{aligned} & \text{minimize} && \max \{g_1(x), \dots, g_r(x)\} \\ & \text{subject to} && x \in X, \end{aligned}$$

where $g_j : \Re^n \mapsto \Re$ are any real-valued functions, is equivalent to the constrained optimization problem

$$\begin{aligned} & \text{minimize} && y \\ & \text{subject to} && x \in X, \quad g_j(x) \leq y, \quad j = 1, \dots, r, \end{aligned}$$

where y is an additional scalar optimization variable. Minimax problems will be discussed further later, in Section 1.4, as an example of problems that may involve a large number of constraints.

Conic Programming

An important problem structure, which can be analyzed as a special case of the Fenchel duality framework is *conic programming*. This is the problem

$$\begin{aligned} & \text{minimize} && f(x) \\ & \text{subject to} && x \in C, \end{aligned} \tag{1.17}$$

where $f : \Re^n \mapsto (-\infty, \infty]$ is a closed proper convex function and C is a closed convex cone in \Re^n .

Indeed, let us apply Fenchel duality with A equal to the identity and the definitions

$$f_1(x) = f(x), \quad f_2(x) = \begin{cases} 0 & \text{if } x \in C, \\ \infty & \text{if } x \notin C. \end{cases}$$

The corresponding conjugates are

$$f_1^*(\lambda) = \sup_{x \in \mathbb{R}^n} \{\lambda'x - f(x)\}, \quad f_2^*(\lambda) = \sup_{x \in C} \lambda'x = \begin{cases} 0 & \text{if } \lambda \in C^*, \\ \infty & \text{if } \lambda \notin C^*, \end{cases}$$

where

$$C^* = \{\lambda \mid \lambda'x \leq 0, \forall x \in C\}$$

is the polar cone of C (note that f_2^* is the support function of C ; cf. Section 1.6 of Appendix B). The dual problem is

$$\begin{aligned} & \text{minimize} && f^*(\lambda) \\ & \text{subject to} && \lambda \in \hat{C}, \end{aligned} \tag{1.18}$$

where f^* is the conjugate of f and \hat{C} is the negative polar cone (also called the *dual cone* of C):

$$\hat{C} = -C^* = \{\lambda \mid \lambda'x \geq 0, \forall x \in C\}.$$

Note the symmetry between primal and dual problems. The strong duality relation $f^* = q^*$ can be written as

$$\inf_{x \in C} f(x) = - \inf_{\lambda \in \hat{C}} f^*(\lambda).$$

The following proposition translates the conditions of Prop. 1.2.1(a), which guarantees that there is no duality gap and that the dual problem has an optimal solution.

Proposition 1.2.2: (Conic Duality Theorem) Assume that the primal conic problem (1.17) has finite optimal value, and moreover $\text{ri}(\text{dom}(f)) \cap \text{ri}(C) \neq \emptyset$. Then, there is no duality gap and the dual problem (1.18) has an optimal solution.

Using the symmetry of the primal and dual problems, we also obtain that there is no duality gap and the primal problem (1.17) has an optimal solution if the optimal value of the dual conic problem (1.18) is finite and $\text{ri}(\text{dom}(f^*)) \cap \text{ri}(\hat{C}) \neq \emptyset$. It is also possible to derive primal and dual optimality conditions by translating the optimality conditions of the Fenchel duality framework [Prop. 1.2.1(c)].

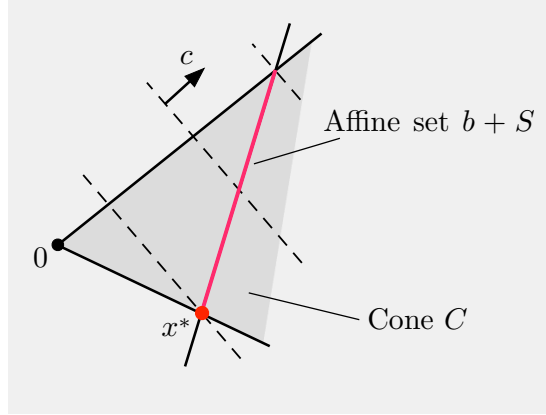


Figure 1.2.1. Illustration of a linear-conic problem: minimizing a linear function $c'x$ over the intersection of an affine set $b + S$ and a convex cone C .

1.2.1 Linear-Conic Problems

An important special case of conic programming, called *linear-conic problem*, arises when $\text{dom}(f)$ is an affine set and f is linear over $\text{dom}(f)$, i.e.,

$$f(x) = \begin{cases} c'x & \text{if } x \in b + S, \\ \infty & \text{if } x \notin b + S, \end{cases}$$

where b and c are given vectors, and S is a subspace. Then the primal problem can be written as

$$\begin{aligned} & \text{minimize} && c'x \\ & \text{subject to} && x - b \in S, \quad x \in C; \end{aligned} \tag{1.19}$$

see Fig. 1.2.1.

To derive the dual problem, we note that

$$\begin{aligned} f^*(\lambda) &= \sup_{x-b \in S} (\lambda - c)'x \\ &= \sup_{y \in S} (\lambda - c)'(y + b) \\ &= \begin{cases} (\lambda - c)'b & \text{if } \lambda - c \in S^\perp, \\ \infty & \text{if } \lambda - c \notin S^\perp. \end{cases} \end{aligned}$$

It can be seen that the dual problem $\min_{\lambda \in \hat{C}} f^*(\lambda)$ [cf. Eq. (1.18)], after discarding the superfluous term $c'b$ from the cost, can be written as

$$\begin{aligned} & \text{minimize} && b'\lambda \\ & \text{subject to} && \lambda - c \in S^\perp, \quad \lambda \in \hat{C}, \end{aligned} \tag{1.20}$$

where \hat{C} is the dual cone:

$$\hat{C} = \{\lambda \mid \lambda'x \geq 0, \forall x \in C\}.$$

By specializing the conditions of the Conic Duality Theorem (Prop. 1.2.2) to the linear-conic duality context, we obtain the following.

Proposition 1.2.3: (Linear-Conic Duality Theorem) Assume that the primal problem (1.19) has finite optimal value, and moreover $(b+S) \cap \text{ri}(C) \neq \emptyset$. Then, there is no duality gap and the dual problem has an optimal solution.

Special Forms of Linear-Conic Problems

The primal and dual linear-conic problems (1.19) and (1.20) have been placed in an elegant symmetric form. There are also other useful formats that parallel and generalize similar formats in linear programming. For example, we have the following dual problem pairs:

$$\min_{Ax=b, x \in C} c'x \quad \Longleftrightarrow \quad \max_{c-A'\lambda \in \hat{C}} b'\lambda, \quad (1.21)$$

$$\min_{Ax-b \in C} c'x \quad \Longleftrightarrow \quad \max_{A'\lambda=c, \lambda \in \hat{C}} b'\lambda, \quad (1.22)$$

where A is an $m \times n$ matrix, and $x \in \mathbb{R}^n$, $\lambda \in \mathbb{R}^m$, $c \in \mathbb{R}^n$, $b \in \mathbb{R}^m$.

To verify the duality relation (1.21), let \bar{x} be any vector such that $A\bar{x} = b$, and let us write the primal problem on the left in the primal conic form (1.19) as

$$\begin{aligned} &\text{minimize} && c'x \\ &\text{subject to} && x - \bar{x} \in N(A), \quad x \in C, \end{aligned}$$

where $N(A)$ is the nullspace of A . The corresponding dual conic problem (1.20) is to solve for μ the problem

$$\begin{aligned} &\text{minimize} && \bar{x}'\mu \\ &\text{subject to} && \mu - c \in N(A)^\perp, \quad \mu \in \hat{C}. \end{aligned} \quad (1.23)$$

Since $N(A)^\perp$ is equal to $\text{Ra}(A')$, the range of A' , the constraints of problem (1.23) can be equivalently written as $c - \mu \in -\text{Ra}(A') = \text{Ra}(A')$, $\mu \in \hat{C}$, or

$$c - \mu = A'\lambda, \quad \mu \in \hat{C},$$

for some $\lambda \in \mathbb{R}^m$. Making the change of variables $\mu = c - A'\lambda$, the dual problem (1.23) can be written as

$$\begin{aligned} & \text{minimize} \quad \bar{x}'(c - A'\lambda) \\ & \text{subject to} \quad c - A'\lambda \in \hat{C}. \end{aligned}$$

By discarding the constant $\bar{x}'c$ from the cost function, using the fact $A\bar{x} = b$, and changing from minimization to maximization, we see that this dual problem is equivalent to the one in the right-hand side of the duality pair (1.21). The duality relation (1.22) is proved similarly.

We next discuss two important special cases of conic programming: *second order cone programming* and *semidefinite programming*. These problems involve two different special cones, and an explicit definition of the affine set constraint. They arise in a variety of applications, and their computational difficulty in practice tends to lie between that of linear and quadratic programming on one hand, and general convex programming on the other hand.

1.2.2 Second Order Cone Programming

In this section we consider the linear-conic problem (1.22), with the cone

$$C = \left\{ (x_1, \dots, x_n) \mid x_n \geq \sqrt{x_1^2 + \dots + x_{n-1}^2} \right\},$$

which is known as the *second order cone* (see Fig. 1.2.2). The dual cone is

$$\hat{C} = \{y \mid 0 \leq y'x, \forall x \in C\} = \left\{ y \mid 0 \leq \inf_{\|(x_1, \dots, x_{n-1})\| \leq x_n} y'x \right\},$$

and it can be shown that $\hat{C} = C$. This property is referred to as *self-duality* of the second order cone, and is fairly evident from Fig. 1.2.2. For a proof, we write

$$\begin{aligned} \inf_{\|(x_1, \dots, x_{n-1})\| \leq x_n} y'x &= \inf_{x_n \geq 0} \left\{ y_n x_n + \inf_{\|(x_1, \dots, x_{n-1})\| \leq x_n} \sum_{i=1}^{n-1} y_i x_i \right\} \\ &= \inf_{x_n \geq 0} \{ y_n x_n - \|(y_1, \dots, y_{n-1})\| x_n \} \\ &= \begin{cases} 0 & \text{if } \|(y_1, \dots, y_{n-1})\| \leq y_n, \\ -\infty & \text{otherwise,} \end{cases} \end{aligned}$$

where the second equality follows because the minimum of the inner product of a vector $z \in \mathbb{R}^{n-1}$ with vectors in the unit ball of \mathbb{R}^{n-1} is $-\|z\|$. Combining the preceding two relations, we have

$$y \in \hat{C} \quad \text{if and only if} \quad 0 \leq y_n - \|(y_1, \dots, y_{n-1})\|,$$

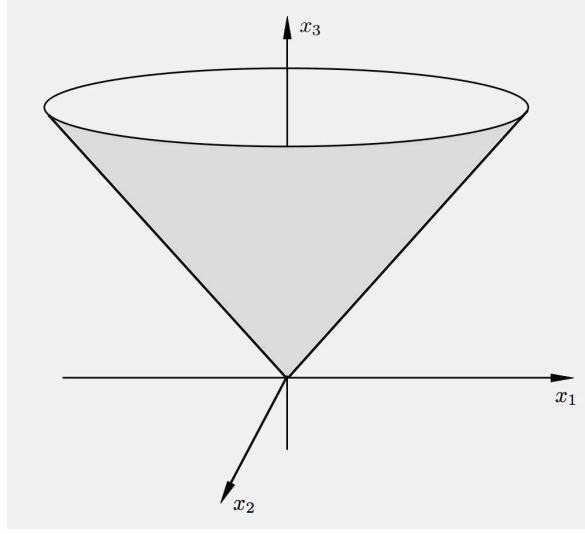


Figure 1.2.2. The second order cone

$$C = \left\{ (x_1, \dots, x_n) \mid x_n \geq \sqrt{x_1^2 + \dots + x_{n-1}^2} \right\},$$

in \mathbb{R}^3 .

so $\hat{C} = C$.

The second order cone programming problem (SOCP for short) is

$$\begin{aligned} & \text{minimize} && c'x \\ & \text{subject to} && A_i x - b_i \in C_i, \quad i = 1, \dots, m, \end{aligned} \tag{1.24}$$

where $x \in \mathbb{R}^n$, c is a vector in \mathbb{R}^n , and for $i = 1, \dots, m$, A_i is an $n_i \times n$ matrix, b_i is a vector in \mathbb{R}^{n_i} , and C_i is the second order cone of \mathbb{R}^{n_i} . It is seen to be a special case of the primal problem in the left-hand side of the duality relation (1.22), where

$$A = \begin{pmatrix} A_1 \\ \vdots \\ A_m \end{pmatrix}, \quad b = \begin{pmatrix} b_1 \\ \vdots \\ b_m \end{pmatrix}, \quad C = C_1 \times \dots \times C_m.$$

Note that linear inequality constraints of the form $a'_i x - b_i \geq 0$ can be written as

$$\begin{pmatrix} 0 \\ a'_i \end{pmatrix} x - \begin{pmatrix} 0 \\ b_i \end{pmatrix} \in C_i,$$

where C_i is the second order cone of \mathbb{R}^2 . As a result, linear-conic problems involving second order cones contain as special cases linear programming problems.

We now observe that from the right-hand side of the duality relation (1.22), and the self-duality relation $C = \hat{C}$, the corresponding dual linear-conic problem has the form

$$\begin{aligned} & \text{maximize} && \sum_{i=1}^m b'_i \lambda_i \\ & \text{subject to} && \sum_{i=1}^m A'_i \lambda_i = c, \quad \lambda_i \in C_i, \quad i = 1, \dots, m, \end{aligned} \tag{1.25}$$

where $\lambda = (\lambda_1, \dots, \lambda_m)$. By applying the Linear-Conic Duality Theorem (Prop. 1.2.3), we have the following.

Proposition 1.2.4: (Second Order Cone Duality Theorem)

Consider the primal SOCP (1.24), and its dual problem (1.25).

- (a) If the optimal value of the primal problem is finite and there exists a feasible solution \bar{x} such that

$$A_i \bar{x} - b_i \in \text{int}(C_i), \quad i = 1, \dots, m,$$

then there is no duality gap, and the dual problem has an optimal solution.

- (b) If the optimal value of the dual problem is finite and there exists a feasible solution $\bar{\lambda} = (\bar{\lambda}_1, \dots, \bar{\lambda}_m)$ such that

$$\bar{\lambda}_i \in \text{int}(C_i), \quad i = 1, \dots, m,$$

then there is no duality gap, and the primal problem has an optimal solution.

Note that while the Linear-Conic Duality Theorem requires a relative interior point condition, the preceding proposition requires an interior point condition. The reason is that the second order cone has nonempty interior, so its relative interior coincides with its interior.

The SOCP arises in many application contexts, and significantly, it can be solved numerically with powerful specialized algorithms that belong to the class of interior point methods, which will be discussed in Section 6.8. We refer to the literature for a more detailed description and analysis (see e.g., the books [BeN01], [BoV04]).

Generally, SOCPs can be recognized from the presence of convex quadratic functions in the cost or the constraint functions. The following are illustrative examples. The first example relates to the field of robust optimization, which involves optimization under uncertainty described by set membership.

Example 1.2.3: (Robust Linear Programming)

Frequently, there is uncertainty about the data of an optimization problem, so one would like to have a solution that is adequate for a whole range of the uncertainty. A popular formulation of this type, is to assume that the constraints contain parameters that take values in a given set, and require that the constraints are satisfied for all values in that set. This approach is also known as a set membership description of the uncertainty and has been used in fields other than optimization, such as set membership estimation, and minimax control (see the textbook [Ber07], which also surveys earlier work).

As an example, consider the problem

$$\begin{aligned} & \text{minimize} && c'x \\ & \text{subject to} && a_j'x \leq b_j, \quad \forall (a_j, b_j) \in T_j, \quad j = 1, \dots, r, \end{aligned} \quad (1.26)$$

where $c \in \mathbb{R}^n$ is a given vector, and T_j is a given subset of \mathbb{R}^{n+1} to which the constraint parameter vectors (a_j, b_j) must belong. The vector x must be chosen so that the constraint $a_j'x \leq b_j$ is satisfied for all $(a_j, b_j) \in T_j$, $j = 1, \dots, r$.

Generally, when T_j contains an infinite number of elements, this problem involves a correspondingly infinite number of constraints. To convert the problem to one involving a finite number of constraints, we note that

$$a_j'x \leq b_j, \quad \forall (a_j, b_j) \in T_j \quad \text{if and only if} \quad g_j(x) \leq 0,$$

where

$$g_j(x) = \sup_{(a_j, b_j) \in T_j} \{a_j'x - b_j\}. \quad (1.27)$$

Thus, the robust linear programming problem (1.26) is equivalent to

$$\begin{aligned} & \text{minimize} && c'x \\ & \text{subject to} && g_j(x) \leq 0, \quad j = 1, \dots, r. \end{aligned}$$

For special choices of the set T_j , the function g_j can be expressed in closed form, and in the case where T_j is an ellipsoid, it turns out that the constraint $g_j(x) \leq 0$ can be expressed in terms of a second order cone. To see this, let

$$T_j = \{(\bar{a}_j + P_j u_j, \bar{b}_j + q_j' u_j) \mid \|u_j\| \leq 1, u_j \in \mathbb{R}^{n_j}\}, \quad (1.28)$$

where P_j is a given $n \times n_j$ matrix, $\bar{a}_j \in \mathbb{R}^n$ and $q_j \in \mathbb{R}^{n_j}$ are given vectors, and \bar{b}_j is a given scalar. Then, from Eqs. (1.27) and (1.28),

$$\begin{aligned} g_j(x) &= \sup_{\|u_j\| \leq 1} \{(\bar{a}_j + P_j u_j)'x - (\bar{b}_j + q_j' u_j)\} \\ &= \sup_{\|u_j\| \leq 1} (P_j'x - q_j)'u_j + \bar{a}_j'x - \bar{b}_j, \end{aligned}$$

and finally

$$g_j(x) = \|P'_j x - q_j\| + \bar{a}'_j x - \bar{b}_j.$$

Thus,

$$g_j(x) \leq 0 \quad \text{if and only if} \quad (P'_j x - q_j, \bar{b}_j - \bar{a}'_j x) \in C_j,$$

where C_j is the second order cone of \mathbb{R}^{n_j+1} ; i.e., the “robust” constraint $g_j(x) \leq 0$ is equivalent to a second order cone constraint. It follows that in the case of ellipsoidal uncertainty, the robust linear programming problem (1.26) is a SOCP of the form (1.24).

Example 1.2.4: (Quadratically Constrained Quadratic Problems)

Consider the quadratically constrained quadratic problem

$$\begin{aligned} & \text{minimize} && x' Q_0 x + 2q'_0 x + p_0 \\ & \text{subject to} && x' Q_j x + 2q'_j x + p_j \leq 0, \quad j = 1, \dots, r, \end{aligned}$$

where Q_0, \dots, Q_r are symmetric $n \times n$ positive definite matrices, q_0, \dots, q_r are vectors in \mathbb{R}^n , and p_0, \dots, p_r are scalars. We show that the problem can be converted to the second order cone format. A similar conversion is also possible for the quadratic programming problem where Q_0 is positive definite and $Q_j = 0$, $j = 1, \dots, r$.

Indeed, since each Q_j is symmetric and positive definite, we have

$$\begin{aligned} x' Q_j x + 2q'_j x + p_j &= \left(Q_j^{1/2} x \right)' Q_j^{1/2} x + 2 \left(Q_j^{-1/2} q_j \right)' Q_j^{1/2} x + p_j \\ &= \|Q_j^{1/2} x + Q_j^{-1/2} q_j\|^2 + p_j - q'_j Q_j^{-1} q_j, \end{aligned}$$

for $j = 0, 1, \dots, r$. Thus, the problem can be written as

$$\begin{aligned} & \text{minimize} && \|Q_0^{1/2} x + Q_0^{-1/2} q_0\|^2 + p_0 - q'_0 Q_0^{-1} q_0 \\ & \text{subject to} && \|Q_j^{1/2} x + Q_j^{-1/2} q_j\|^2 + p_j - q'_j Q_j^{-1} q_j \leq 0, \quad j = 1, \dots, r, \end{aligned}$$

or, by neglecting the constant $p_0 - q'_0 Q_0^{-1} q_0$,

$$\begin{aligned} & \text{minimize} && \|Q_0^{1/2} x + Q_0^{-1/2} q_0\| \\ & \text{subject to} && \|Q_j^{1/2} x + Q_j^{-1/2} q_j\| \leq (q'_j Q_j^{-1} q_j - p_j)^{1/2}, \quad j = 1, \dots, r. \end{aligned}$$

By introducing an auxiliary variable x_{n+1} , the problem can be written as

$$\begin{aligned} & \text{minimize} && x_{n+1} \\ & \text{subject to} && \|Q_0^{1/2} x + Q_0^{-1/2} q_0\| \leq x_{n+1} \\ & && \|Q_j^{1/2} x + Q_j^{-1/2} q_j\| \leq (q'_j Q_j^{-1} q_j - p_j)^{1/2}, \quad j = 1, \dots, r. \end{aligned}$$

It can be seen that this problem has the second order cone form (1.24). In particular, the first constraint is of the form $A_0x - b_0 \in C$, where C is the second order cone of \mathbb{R}^{n+1} and the $(n+1)$ st component of $A_0x - b_0$ is x_{n+1} . The remaining r constraints are of the form $A_jx - b_j \in C$, where the $(n+1)$ st component of $A_jx - b_j$ is the scalar $(q'_j Q_j^{-1} q_j - p_j)^{1/2}$.

We finally note that the problem of this example is special in that it has no duality gap, assuming its optimal value is finite, i.e., there is no need for the interior point conditions of Prop. 1.2.4. This can be traced to the fact that linear transformations preserve the closure of sets defined by quadratic constraints (see e.g., BNO03], Section 1.5.2).

1.2.3 Semidefinite Programming

In this section we consider the linear-conic problem (1.21) with C being the cone of matrices that are positive semidefinite.[†] This is called the *positive semidefinite cone*. To define the problem, we view the space of symmetric $n \times n$ matrices as the space \mathbb{R}^{n^2} with the inner product

$$\langle X, Y \rangle = \text{trace}(XY) = \sum_{i=1}^n \sum_{j=1}^n x_{ij} y_{ij}.$$

The interior of C is the set of positive definite matrices.

The dual cone is

$$\hat{C} = \{Y \mid \text{trace}(XY) \geq 0, \forall X \in C\},$$

and it can be shown that $\hat{C} = C$, i.e., C is self-dual. Indeed, if $Y \notin C$, there exists a vector $v \in \mathbb{R}^n$ such that

$$0 > v'Yv = \text{trace}(vv'Y).$$

Hence the positive semidefinite matrix $X = vv'$ satisfies $0 > \text{trace}(XY)$, so $Y \notin \hat{C}$ and it follows that $C \supset \hat{C}$. Conversely, let $Y \in C$, and let X be any positive semidefinite matrix. We can express X as

$$X = \sum_{i=1}^n \lambda_i e_i e_i',$$

where λ_i are the nonnegative eigenvalues of X , and e_i are corresponding orthonormal eigenvectors. Then,

$$\text{trace}(XY) = \text{trace}\left(Y \sum_{i=1}^n \lambda_i e_i e_i'\right) = \sum_{i=1}^n \lambda_i e_i' Y e_i \geq 0.$$

[†] As noted in Appendix A, throughout this book a positive semidefinite matrix is implicitly assumed to be symmetric.

It follows that $Y \in \hat{C}$ and $C \subset \hat{C}$. Thus C is self-dual, $C = \hat{C}$.

The semidefinite programming problem (SDP for short) is to minimize a linear function of a symmetric matrix over the intersection of an affine set with the positive semidefinite cone. It has the form

$$\begin{aligned} & \text{minimize} && \langle D, X \rangle \\ & \text{subject to} && \langle A_i, X \rangle = b_i, \quad i = 1, \dots, m, \quad X \in C, \end{aligned} \quad (1.29)$$

where D, A_1, \dots, A_m , are given $n \times n$ symmetric matrices, and b_1, \dots, b_m , are given scalars. It is seen to be a special case of the primal problem in the left-hand side of the duality relation (1.21).

We can view the SDP as a problem with linear cost, linear constraints, and a convex set constraint. Then, similar to the case of SOCP, it can be verified that the dual problem (1.20), as given by the right-hand side of the duality relation (1.21), takes the form

$$\begin{aligned} & \text{maximize} && b' \lambda \\ & \text{subject to} && D - (\lambda_1 A_1 + \dots + \lambda_m A_m) \in C, \end{aligned} \quad (1.30)$$

where $b = (b_1, \dots, b_m)$ and the maximization is over the vector $\lambda = (\lambda_1, \dots, \lambda_m)$. By applying the Linear-Conic Duality Theorem (Prop. 1.2.3), we have the following proposition.

Proposition 1.2.5: (Semidefinite Duality Theorem) Consider the primal SDP (1.29), and its dual problem (1.30).

- (a) If the optimal value of the primal problem is finite and there exists a primal-feasible solution, which is positive definite, then there is no duality gap, and the dual problem has an optimal solution.
- (b) If the optimal value of the dual problem is finite and there exist scalars $\bar{\lambda}_1, \dots, \bar{\lambda}_m$ such that $D - (\bar{\lambda}_1 A_1 + \dots + \bar{\lambda}_m A_m)$ is positive definite, then there is no duality gap, and the primal problem has an optimal solution.

The SDP is a fairly general problem. In particular, it can be shown that a SOCP can be cast as a SDP. Thus SDP involves a more general structure than SOCP. This is consistent with the practical observation that the latter problem is generally more amenable to computational solution. We provide some examples of problem formulation as an SDP.

Example 1.2.5: (Minimizing the Maximum Eigenvalue)

Given a symmetric $n \times n$ matrix $M(\lambda)$, which depends on a parameter vector $\lambda = (\lambda_1, \dots, \lambda_m)$, we want to choose λ so as to minimize the maximum

eigenvalue of $M(\lambda)$. We pose this problem as

$$\begin{aligned} & \text{minimize} && z \\ & \text{subject to} && \text{maximum eigenvalue of } M(\lambda) \leq z, \end{aligned}$$

or equivalently

$$\begin{aligned} & \text{minimize} && z \\ & \text{subject to} && zI - M(\lambda) \in C, \end{aligned}$$

where I is the $n \times n$ identity matrix, and C is the semidefinite cone. If $M(\lambda)$ is an affine function of λ ,

$$M(\lambda) = M_0 + \lambda_1 M_1 + \cdots + \lambda_m M_m,$$

the problem has the form of the dual problem (1.30), with the optimization variables being $(z, \lambda_1, \dots, \lambda_m)$.

Example 1.2.6: (Semidefinite Relaxation – Lower Bounds for Discrete Optimization Problems)

Semidefinite programming provides a means for deriving lower bounds to the optimal value of several types of discrete optimization problems. As an example, consider the following quadratic problem with quadratic equality constraints

$$\begin{aligned} & \text{minimize} && x'Q_0x + a_0'x + b_0 \\ & \text{subject to} && x'Q_i x + a_i'x + b_i = 0, \quad i = 1, \dots, m, \end{aligned} \tag{1.31}$$

where Q_0, \dots, Q_m are symmetric $n \times n$ matrices, a_0, \dots, a_m are vectors in \mathbb{R}^n , and b_0, \dots, b_m are scalars.

This problem can be used to model broad classes of discrete optimization problems. To see this, consider an integer constraint that a variable x_i must be either 0 or 1. Such a constraint can be expressed by the quadratic equality $x_i^2 - x_i = 0$. Furthermore, a linear inequality constraint $a_j'x \leq b_j$ can be expressed as the quadratic equality constraint $y_j^2 + a_j'x - b_j = 0$, where y_j is an additional variable.

Introducing a multiplier vector $\lambda = (\lambda_1, \dots, \lambda_m)$, the dual function is given by

$$q(\lambda) = \inf_{x \in \mathbb{R}^n} \{x'Q(\lambda)x + a(\lambda)'x + b(\lambda)\},$$

where

$$Q(\lambda) = Q_0 + \sum_{i=1}^m \lambda_i Q_i, \quad a(\lambda) = a_0 + \sum_{i=1}^m \lambda_i a_i, \quad b(\lambda) = b_0 + \sum_{i=1}^m \lambda_i b_i.$$

Let f^* and q^* be the optimal values of problem (1.31) and its dual, and note that by weak duality, we have $f^* \geq q^*$. By introducing an auxiliary

scalar variable ξ , we see that the dual problem is to find a pair (ξ, λ) that solves the problem

$$\begin{aligned} & \text{maximize} && \xi \\ & \text{subject to} && q(\lambda) \geq \xi. \end{aligned}$$

The constraint $q(\lambda) \geq \xi$ of this problem can be written as

$$\inf_{x \in \mathbb{R}^n} \{x'Q(\lambda)x + a(\lambda)'x + b(\lambda) - \xi\} \geq 0,$$

or equivalently, introducing a scalar variable t and multiplying with t^2 ,

$$\inf_{x \in \mathbb{R}^n, t \in \mathbb{R}} \{(tx)'Q(\lambda)(tx) + a(\lambda)'(tx)t + (b(\lambda) - \xi)t^2\} \geq 0.$$

Writing $y = tx$, this relation takes the form of a quadratic in (y, t) ,

$$\inf_{y \in \mathbb{R}^n, t \in \mathbb{R}} \{y'Q(\lambda)y + a(\lambda)'yt + (b(\lambda) - \xi)t^2\} \geq 0,$$

or

$$\begin{pmatrix} Q(\lambda) & \frac{1}{2}a(\lambda) \\ \frac{1}{2}a(\lambda)' & b(\lambda) - \xi \end{pmatrix} \in C, \quad (1.32)$$

where C is the positive semidefinite cone. Thus the dual problem is equivalent to the SDP of maximizing ξ over all (ξ, λ) satisfying the constraint (1.32), and its optimal value q^* is a lower bound to f^* .

1.3 ADDITIVE COST PROBLEMS

In this section we focus on a structural characteristic that arises in several important contexts: a cost function f that is the sum of a large number of components $f_i : \mathbb{R}^n \mapsto \mathbb{R}$,

$$f(x) = \sum_{i=1}^m f_i(x). \quad (1.33)$$

Such cost functions can be minimized with specialized methods, called *incremental*, which exploit their additive structure, by updating x using one component function f_i at a time (see Section 2.1.5). Problems with additive cost functions can also be treated with specialized outer and inner linearization methods that approximate the component functions f_i individually (rather than approximating f); see Section 4.4.

An important special case is the cost function of the dual of a separable problem

$$\begin{aligned} & \text{maximize} && \sum_{i=1}^m q_i(\mu) \\ & \text{subject to} && \mu \geq 0, \end{aligned}$$

where

$$q_i(\mu) = \inf_{x_i \in X_i} \left\{ f_i(x_i) + \sum_{j=1}^r \mu_j g_{ij}(x_i) \right\},$$

and $\mu = (\mu_1, \dots, \mu_r)$ [cf. Eq. (1.8)]. After a sign change to convert to minimization it takes the form (1.33) with $f_i(\mu) = -q_i(\mu)$. This is a major class of additive cost problems.

We will next describe some applications from a variety of fields. The following five examples arise in many machine learning contexts.

Example 1.3.1: (Regularized Regression)

This is a broad class of applications that relate to parameter estimation. The cost function involves a sum of terms $f_i(x)$, each corresponding to the error between some data and the output of a parametric model, with x being the vector of parameters. An example is linear least squares problems, also referred to as *linear regression* problems, where f_i has quadratic structure. Often a convex regularization function $R(x)$ is added to the least squares objective, to induce desirable properties of the solution and/or the corresponding algorithms. This gives rise to problems of the form

$$\begin{aligned} & \text{minimize} && R(x) + \frac{1}{2} \sum_{i=1}^m (c_i'x - b_i)^2 \\ & \text{subject to} && x \in \mathbb{R}^n, \end{aligned}$$

where c_i and b_i are given vectors and scalars, respectively. The regularization function R is often taken to be differentiable, and particularly quadratic. However, there are practically important examples of nondifferentiable choices (see the next example).

In statistical applications, such a problem arises when constructing a linear model for an unknown input-output relation. The model involves a vector of parameters x , to be determined, which weigh input data (the components of the vectors c_i). The inner products $c_i'x$ produced by the model are matched against the scalars b_i , which are observed output data, corresponding to inputs c_i from the true input-output relation that we try to represent. The optimal vector of parameters x^* provides the model that (in the absence of a regularization function) minimizes the sum of the squared errors $(c_i'x^* - b_i)^2$.

In a more general version of the problem, a nonlinear parametric model is constructed, giving rise to a nonlinear least squares problem of the form

$$\begin{aligned} & \text{minimize} && R(x) + \sum_{i=1}^m |g_i(x)|^2 \\ & \text{subject to} && x \in \mathbb{R}^n, \end{aligned}$$

where $g_i : \mathbb{R}^n \mapsto \mathbb{R}$ are given nonlinear functions that depend on the data. This is also a common problem, referred to as *nonlinear regression*, which, however, is often nonconvex [it is convex if the functions g_i are convex and also nonnegative, i.e., $g_i(x) \geq 0$ for all $x \in \mathbb{R}^n$].

It is also possible to use a nonquadratic function of the error between some data and the output of a linear parametric model. Thus in place of the squared error $(1/2)(c'_i x - b_i)^2$, we may use $h_i(c'_i x - b_i)$, where $h_i : \Re \mapsto \Re$ is a convex function, leading to the problem

$$\begin{aligned} & \text{minimize} && R(x) + \sum_{i=1}^m h_i(c'_i x - b_i) \\ & \text{subject to} && x \in \Re^n. \end{aligned}$$

Generally the choice of the function h_i is dictated by statistical modeling considerations, for which the reader may consult the relevant literature. An example is

$$h_i(c'_i x - b_i) = |c'_i x - b_i|,$$

which tends to result in a more robust estimate than least squares in the presence of large outliers in the data. This is known as the *least absolute deviations* method.

There are also constrained variants of the problems just discussed, where the parameter vector x is required to belong to some subset of \Re^n , such as the nonnegative orthant or a “box” formed by given upper and lower bounds on the components of x . Such constraints may be used to encode into the model some prior knowledge about the nature of the solution.

Example 1.3.2: (ℓ_1 -Regularization)

A popular approach to regularized regression involves ℓ_1 -regularization, where

$$R(x) = \gamma \|x\|_1 = \gamma \sum_{j=1}^n |x^j|,$$

γ is a positive scalar and x^j is the j th coordinate of x . The reason for the popularity of the ℓ_1 norm $\|x\|_1$ is that it tends to produce optimal solutions where a greater number of components x^j are zero, relative to the case of quadratic regularization (see Fig. 1.3.1). This is considered desirable in many statistical applications, where the number of parameters to include in a model may not be known a priori; see e.g., [Tib96], [DoE03], [BJM12]. The special case where a linear least squares model is used,

$$\begin{aligned} & \text{minimize} && \gamma \|x\|_1 + \frac{1}{2} \sum_{i=1}^m (c'_i x - b_i)^2 \\ & \text{subject to} && x \in \Re^n, \end{aligned}$$

is known as the *lasso problem*.

In a generalization of the lasso problem, the ℓ_1 regularization function $\|x\|_1$ is replaced by a scaled version $\|Sx\|_1$, where S is some scaling matrix.

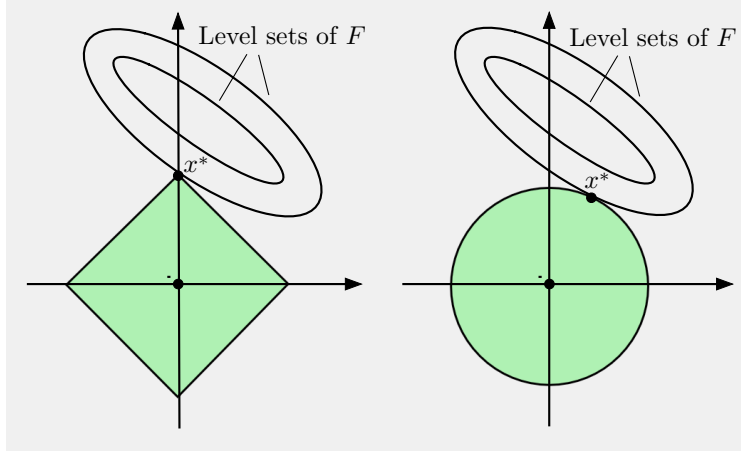


Figure 1.3.1. Illustration of the effect of ℓ_1 -regularization for cost functions of the form $\gamma\|x\|_1 + F(x)$, where $\gamma > 0$ and $F : \mathbb{R}^n \mapsto \mathbb{R}$ is differentiable (figure in the left-hand side). The optimal solution x^* tends to have more zero components than in the corresponding quadratic regularization case, illustrated in the right-hand side.

The term $\|Sx\|_1$ then induces a penalty on some undesirable characteristic of the solution. For example the problem

$$\begin{aligned} & \text{minimize} && \gamma \sum_{i=1}^{n-1} |x_{i+1} - x_i| + \frac{1}{2} \sum_{i=1}^m (c'_i x - b_i)^2 \\ & \text{subject to} && x \in \mathbb{R}^n, \end{aligned}$$

is known as the *total variation denoising problem*; see e.g., [ROF92], [Cha04], [BeT09a]. The regularization term here encourages consecutive variables to take similar values, and tends to produce more smoothly varying solutions.

Another related example is *matrix completion with nuclear norm regularization*; see e.g., [CaR09], [CaT10], [RFP10], [Rec11], [ReR13]. Here the minimization is over all $m \times n$ matrices X , with components denoted X_{ij} . We have a set of entries M_{ij} , $(i, j) \in \Omega$, where Ω is a subset of index pairs, and we want to find X whose entries X_{ij} are close to M_{ij} for $(i, j) \in \Omega$, and has as small rank as possible, a property that is desirable on the basis of statistical considerations. The following more tractable version of the problem is solved instead:

$$\begin{aligned} & \text{minimize} && \gamma \|X\|_* + \frac{1}{2} \sum_{(i,j) \in \Omega} (X_{ij} - M_{ij})^2 \\ & \text{subject to} && X \in \mathbb{R}^{m \times n}, \end{aligned}$$

where $\|X\|_*$ is the *nuclear norm* of X , defined as the sum of the singular values of X . There is substantial theory that justifies this approximation, for which we refer to the literature. It turns out that the nuclear norm is a convex function with some nice properties. In particular, its subdifferential at any X can be conveniently characterized for use in algorithms.

Let us finally note that sometimes additional regularization functions are used in conjunction with ℓ_1 -type terms. An example is the sum of a quadratic and an ℓ_1 -type term.

Example 1.3.3: (Classification)

In the regression problems of the preceding examples we aim to construct a parametric model that matches well an input-output relationship based on given data. Similar problems arise in a classification context, where we try to construct a parametric model for predicting whether an object with certain characteristics (also called features) belongs to a given category or not.

We assume that each object is characterized by a *feature vector* c that belongs to \mathbb{R}^n and a *label* b that takes the values $+1$ or -1 , if the object belongs to the category or not, respectively. As illustration consider a credit card company that wishes to classify applicants as “low risk” ($+1$) or “high risk” (-1), with each customer characterized by n scalar features of financial and personal type.

We are given data, which is a set of feature-label pairs (c_i, b_i) , $i = 1, \dots, m$. Based on this data, we want to find a parameter vector $x \in \mathbb{R}^n$ and a scalar $y \in \mathbb{R}$ such that the sign of $c'_i x + y$ is a good predictor of the label of an object with feature vector c . Thus, loosely speaking, x and y should be such that for “most” of the given feature-label data (c_i, b_i) we have

$$c'_i x + y > 0, \quad \text{if } b_i = +1,$$

$$c'_i x + y < 0, \quad \text{if } b_i = -1.$$

In the statistical literature, $c'_i x + y$ is often called the *discriminant function*, and the value of

$$b_i(c'_i x + y),$$

for a given object i provides a measure of “margin” to misclassification of the object. In particular, a classification error is made for object i when $b_i(c'_i x + y) < 0$.

Thus it makes sense to formulate classification as an optimization problem where negative values of $b_i(c'_i x + y)$ are penalized. This leads to the problem

$$\text{minimize} \quad R(x) + \sum_{i=1}^m h(b_i(c'_i x + y))$$

$$\text{subject to} \quad x \in \mathbb{R}^n, \quad y \in \mathbb{R},$$

where R is a suitable regularization function, and $h : \mathbb{R} \mapsto \mathbb{R}$ is a convex function that penalizes negative values of its argument. It would make some sense to use a penalty of one unit for misclassification, i.e.,

$$h(z) = \begin{cases} 0 & \text{if } z \geq 0, \\ 1 & \text{if } z < 0, \end{cases}$$

but such a penalty function is discontinuous. To obtain a continuous cost function, we allow a continuous transition of h from negative to positive

values, leading to a variety of nonincreasing functions h . The choice of h depends on the given application and other theoretical considerations for which we refer to the literature. Some common examples are

$$\begin{aligned} h(z) &= e^{-z}, & (\text{exponential loss}), \\ h(z) &= \log(1 + e^{-z}), & (\text{logistic loss}), \\ h(z) &= \max\{0, 1 - z\}, & (\text{hinge loss}). \end{aligned}$$

For the case of logistic loss the method comes under the methodology of *logistic regression*, and for the case of hinge loss the method comes under the methodology of *support vector machines*. As in the case of regression, the regularization function R could be quadratic, the ℓ_1 norm, or some scaled version or combination thereof. There is extensive literature on these methodologies and their applications, to which we refer for further discussion.

Example 1.3.4: (Nonnegative Matrix Factorization)

The nonnegative matrix factorization problem is to approximately factor a given nonnegative matrix B as CX , where C and X are nonnegative matrices to be determined via the optimization

$$\begin{aligned} &\text{minimize} && \|CX - B\|_F^2 \\ &\text{subject to} && C \geq 0, X \geq 0. \end{aligned}$$

Here $\|\cdot\|_F$ denotes the Frobenius norm of a matrix ($\|M\|_F^2$ is the sum of the squares of the scalar components of M). The matrices B , C , and X must have compatible dimensions, with the column dimension of C usually being much smaller than its row dimension, so that CX is a low-rank approximation of B . In some versions of the problem some of the nonnegativity constraints on the components of C and X may be relaxed. Moreover, regularization terms may be added to the cost function to induce sparsity or some other effect, similar to earlier examples in this section.

This problem, formulated in the 90s, [PaT94], [Paa97], [LeS99], has become a popular model for regression-type applications such as the ones of Example 1.3.1, but with the vectors c_i in the least squares objective $\sum_{i=1}^m (c_i'x - b_i)^2$ being unknown and subject to optimization. In the regression context of Example 1.3.1, we aim to (approximately) represent the data in the range space of the matrix C whose rows are the vectors c_i' , and we may view C as a matrix of known basis functions. In the matrix factorization context of the present example, we aim to discover a “good” matrix C of basis functions that represents well the given data, i.e., the matrix B .

An important characteristic of the problem is that its cost function is not convex jointly in (C, X) . However, it is convex in each of the matrices C and X individually, when the other matrix is held fixed. This facilitates the application of algorithms that involve alternate minimizations with respect to C and with respect to X ; see Section 6.5. We refer to the literature, e.g., the papers [BBL07], [Lin07], [GoZ12], for a discussion of related algorithmic issues.

Example 1.3.5: (Maximum Likelihood Estimation)

The maximum likelihood approach is a major statistical inference methodology for parameter estimation, which is described in many sources (see e.g., the textbooks [Was04], [HTF09]). In fact in many cases, a maximum likelihood formulation is used to provide a probabilistic justification of the regression and classification models of the preceding examples.

Here we observe a sample of a random vector Z whose distribution $P_Z(\cdot; x)$ depends on an unknown parameter vector $x \in \mathbb{R}^n$. For simplicity we assume that Z can take only a finite set of values, so that $P_Z(z; x)$ is the probability that Z takes the value z when the parameter vector has the value x . We estimate x based on the given sample value z , by solving the problem

$$\begin{aligned} & \text{maximize} && P_Z(z; x) \\ & \text{subject to} && x \in \mathbb{R}^n. \end{aligned} \tag{1.34}$$

The cost function $P_Z(z; \cdot)$ of this problem may either have an additive structure or may be equivalent to a problem that has an additive structure. For example the event that $Z = z$ may be the union of a large number of disjoint events, so $P_Z(z; x)$ is the sum of the probabilities of these events. For another important context, suppose that the data z consists of m independent samples z_1, \dots, z_m drawn from a distribution $P(\cdot; x)$, in which case

$$P_Z(z; x) = P(z_1; x) \cdots P(z_m; x).$$

Then the maximization (1.34) is equivalent to the additive cost minimization

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^m f_i(x) \\ & \text{subject to} && x \in \mathbb{R}^n, \end{aligned}$$

where

$$f_i(x) = -\log P(z_i; x).$$

In many applications the number of samples m is very large, in which case special methods that exploit the additive structure of the cost are recommended. Often a suitable regularization term is added to the cost function, similar to the preceding examples.

Example 1.3.6: (Minimization of an Expected Value - Stochastic Programming)

An important context where additive cost functions arise is the minimization of an expected value

$$\begin{aligned} & \text{minimize} && E\{F(x, w)\} \\ & \text{subject to} && x \in X, \end{aligned}$$

where w is a random variable taking a finite but very large number of values w_i , $i = 1, \dots, m$, with corresponding probabilities π_i . Then the cost function consists of the sum of the m functions $\pi_i F(x, w_i)$.

For example, in *stochastic programming*, a classical model of two-stage optimization under uncertainty, a vector $x \in X$ is selected, a random event occurs that has m possible outcomes w_1, \dots, w_m , and another vector $y \in Y$ is selected with knowledge of the outcome that occurred (see e.g., the books [BiL97], [KaW94], [Pre95], [SDR09]). Then for optimization purposes, we need to specify a different vector $y_i \in Y$ for each outcome w_i . The problem is to minimize the expected cost

$$F(x) + \sum_{i=1}^m \pi_i G_i(y_i),$$

where $G_i(y_i)$ is the cost associated with the choice y_i and the occurrence of w_i , and π_i is the corresponding probability. This is a problem with an additive cost function.

Additive cost functions also arise when the expected value cost function $E\{F(x, w)\}$ is approximated by an m -sample average

$$f(x) = \frac{1}{m} \sum_{i=1}^m F(x, w_i),$$

where w_i are independent samples of the random variable w . The minimum of the sample average $f(x)$ is then taken as an approximation of the minimum of $E\{F(x, w)\}$.

Generally additive cost problems arise when we want to strike a balance between several types of costs by lumping them into a single cost function. The following is an example of a different character than the preceding ones.

Example 1.3.7: (Weber Problem in Location Theory)

A basic problem in location theory is to find a point x in the plane whose sum of weighted distances from a given set of points y_1, \dots, y_m is minimized. Mathematically, the problem is

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^m w_i \|x - y_i\| \\ & \text{subject to} && x \in \mathbb{R}^n, \end{aligned}$$

where w_1, \dots, w_m are given positive scalars. This problem has many variations, including constrained versions, and descends from the famous Fermat-Torricelli-Viviani problem (see [BMS99] for an account of the history of this problem). We refer to the book [DrH04] for a survey of recent research, and to the paper [BeT10] for a discussion that is relevant to our context.

The structure of the additive cost function (1.33) often facilitates the use of a distributed computing system that is well-suited for the incremental approach. The following is an illustrative example.

Example 1.3.8: (Distributed Incremental Optimization – Sensor Networks)

Consider a network of m sensors where data are collected and are used to solve some inference problem involving a parameter vector x . If $f_i(x)$ represents an error penalty for the data collected by the i th sensor, the inference problem involves an additive cost function $\sum_{i=1}^m f_i$. While it is possible to collect all the data at a fusion center where the problem will be solved in centralized manner, it may be preferable to adopt a distributed approach in order to save in data communication overhead and/or take advantage of parallelism in computation. In such an approach the current iterate x_k is passed on from one sensor to another, with each sensor i performing an incremental iteration involving just its local component f_i . The entire cost function need not be known at any one location. For further discussion we refer to representative sources such as [RaN04], [RaN05], [BHG08], [MRS10], [GSW12], and [Say14].

The approach of computing incrementally the values and subgradients of the components f_i in a distributed manner can be substantially extended to apply to general systems of asynchronous distributed computation, where the components are processed at the nodes of a computing network, and the results are suitably combined [NBB01] (see our discussion in Sections 2.1.5 and 2.1.6).

Let us finally note a constrained version of additive cost problems where the functions f_i are extended real-valued. This is essentially equivalent to constraining x to lie in the intersection of the domains

$$X_i = \text{dom}(f_i),$$

resulting in a problem of the form

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^m f_i(x) \\ & \text{subject to} && x \in \bigcap_{i=1}^m X_i, \end{aligned}$$

where each f_i is real-valued over the set X_i . Methods that are well-suited for the unconstrained version of the problem where $X_i \equiv \mathbb{R}^n$ can often be modified to apply to the constrained version, as we will see in Chapter 6, where we will discuss incremental constraint projection methods. However, the case of constraint sets with many components arises independently of whether the cost function is additive or not, and has its own character, as we discuss in the next section.

1.4 LARGE NUMBER OF CONSTRAINTS

In this section we consider problems of the form

$$\begin{aligned} & \text{minimize} && f(x) \\ & \text{subject to} && x \in X, \quad g_j(x) \leq 0, \quad j = 1, \dots, r, \end{aligned} \tag{1.35}$$

where the number r of constraints is very large. Problems of this type occur often in practice, either directly or via reformulation from other problems. A similar type of problem arises when the abstract constraint set X consists of the intersection of many simpler sets:

$$X = \cap_{\ell \in L} X_\ell,$$

where L is a finite or infinite index set. There may or may not be additional inequality constraints $g_j(x) \leq 0$ like the ones in problem (1.35). We provide a few examples.

Example 1.4.1: (Feasibility and Minimum Distance Problems)

A simple but important problem, which arises in many contexts and embodies important algorithmic ideas, is a classical *feasibility problem*, where the objective is to find a common point within a collection of sets X_ℓ , $\ell \in L$, where each X_ℓ is a closed convex set. In the feasibility problem the cost function is zero. A somewhat more complex problem with a similar structure arises when there is a cost function, i.e., a problem of the form

$$\begin{aligned} & \text{minimize} && f(x) \\ & \text{subject to} && x \in \cap_{\ell \in L} X_\ell, \end{aligned}$$

where $f : \Re^n \mapsto \Re$. An important example is the minimum distance problem, where

$$f(x) = \|x - z\|,$$

for a given vector z and some norm $\|\cdot\|$. The following example is a special case.

Example 1.4.2: (Basis Pursuit)

Consider the problem

$$\begin{aligned} & \text{minimize} && \|x\|_1 \\ & \text{subject to} && Ax = b, \end{aligned} \tag{1.36}$$

where $\|\cdot\|_1$ is the ℓ_1 norm in \Re^n , A is a given $m \times n$ matrix, and b is a vector in \Re^m that consists of m given measurements. We are trying to construct a linear model of the form $Ax = b$, where x is a vector of n scalar

weights for a large number n of basis functions ($m < n$). We want to satisfy exactly the measurement equations $Ax = b$, while using only a few of the basis functions in our model. Consequently, we introduce the ℓ_1 norm in the cost function of problem (1.36), aiming to delineate a small subset of basis functions, corresponding to nonzero coordinates of x at the optimal solution. This is called the *basis pursuit* problem (see, e.g., [CDS01], [VaF08]), and its underlying idea is similar to the one of ℓ_1 -regularization (cf. Example 1.3.2).

It is also possible to consider a norm other than ℓ_1 in Eq. (1.36). An example is the *atomic norm* $\|\cdot\|_{\mathcal{A}}$ induced by a subset \mathcal{A} that is centrally symmetric around the origin ($a \in \mathcal{A}$ if and only if $-a \in \mathcal{A}$):

$$\|x\|_{\mathcal{A}} = \inf \{t > 0 \mid x \in t \cdot \text{conv}(\mathcal{A})\}.$$

This problem, and other related problems involving atomic norms, have many applications; see for example [CRP12], [SBT12], [RSW13].

A related problem is

$$\begin{aligned} & \text{minimize} && \|X\|_* \\ & \text{subject to} && AX = B, \end{aligned}$$

where the optimization is over all $m \times n$ matrices X . The matrices A , B are given and have dimensions $\ell \times m$ and $\ell \times n$, respectively, and $\|X\|_*$ is the nuclear norm of X . This problem aims to produce a low-rank matrix X that satisfies an underdetermined set of linear equations $AX = B$ (see e.g., [CaR09], [RFP10], [RXB11]). When these equations specify that a subset of entries X_{ij} , $(i, j) \in \Omega$, are fixed at given values M_{ij} ,

$$X_{ij} = M_{ij}, \quad (i, j) \in \Omega,$$

we obtain an alternative formulation of the matrix completion problem discussed in Example 1.3.2.

Example 1.4.3: (Minimax Problems)

In a minimax problem the cost function has the form

$$f(x) = \sup_{z \in Z} \phi(x, z),$$

where Z is a subset of some space and $\phi(\cdot, z)$ is a real-valued function for each $z \in Z$. We want to minimize f subject to $x \in X$, where X is a given constraint set. By introducing an artificial scalar variable y , we may transform such a problem to the general form

$$\begin{aligned} & \text{minimize} && y \\ & \text{subject to} && x \in X, \quad \phi(x, z) \leq y, \quad \forall z \in Z, \end{aligned}$$

which involves a large number of constraints (one constraint for each z in the set Z , which could be infinite). Of course in this problem the set X may also be of the form $X = \bigcap_{\ell \in L} X_{\ell}$ as in earlier examples.

Example 1.4.4: (Basis Function Approximation for Separable Problems – Approximate Dynamic Programming)

Let us consider a large-scale separable problem of the form

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^m f_i(y_i) \\ & \text{subject to} && \sum_{i=1}^m g_{ij}(y_i) \leq 0, \quad \forall j = 1, \dots, r, \quad y \geq 0, \end{aligned} \tag{1.37}$$

where $f_i : \mathbb{R} \mapsto \mathbb{R}$ are scalar functions, and the dimension m of the vector $y = (y_1, \dots, y_m)$ is very large. One possible way to address this problem is to approximate y with a vector of the form Φx , where Φ is an $m \times n$ matrix. The columns of Φ may be relatively few, and may be viewed as basis functions for a low-dimensional approximation subspace $\{\Phi x \mid x \in \mathbb{R}^n\}$. We replace problem (1.37) with the approximate version

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^m f_i(\phi'_i x) \\ & \text{subject to} && \sum_{i=1}^m g_{ij}(\phi'_i x) \leq 0, \quad \forall j = 1, \dots, r, \\ & && \phi'_i x \geq 0, \quad i = 1, \dots, m, \end{aligned} \tag{1.38}$$

where ϕ'_i denotes the i th row of Φ , and $\phi'_i x$ is viewed as an approximation of y_i . Thus the dimension of the problem is reduced from m to n . However, the constraint set of the problem became more complicated, because the simple constraints $y_i \geq 0$ take the more complex form $\phi'_i x \geq 0$. Moreover the number m of additive components in the cost function, as well as the number of its constraints is still large. Thus the problem has the additive cost structure of the preceding section, as well as a large number of constraints.

An important application of this approach is in approximate dynamic programming (see e.g., [BeT96], [SuB98], [Pow11], [Ber12]), where the functions f_i and g_{ij} are linear. The corresponding problem (1.37) relates to the solution of the optimality condition (Bellman equation) of an infinite horizon Markovian decision problem (the constraint $y \geq 0$ may not be present in this context). Here the numbers m and r are often astronomical (in fact r can be much larger than m), in which case an exact solution cannot be obtained. For such problems, approximation based on problem (1.38) has been one of the major algorithmic approaches (see [Ber12] for a textbook presentation and references). For very large m , it may be impossible to calculate the cost function value $\sum_{i=1}^m f_i(\phi'_i x)$ for a given x , and one may at most be able to sample individual cost components f_i . For this reason optimization by stochastic simulation is one of the most prominent approaches in large scale dynamic programming.

Let us also mention that related approaches based on randomization and simulation have been proposed for the solution of large scale instances of classical linear algebra problems; see [BeY09], [Ber12] (Section 7.3), [DMM06], [StV09], [HMT10], [Nee10], [DMM11], [WaB13a], [WaB13b].

A large number of constraints also arises often in problems involving a graph, and may be handled with algorithms that take into account the graph structure. The following example is typical.

Example 1.4.5: (Optimal Routing in a Network – Multicommodity Flows)

Consider a directed graph that is used to transfer “commodities” from given supply points to given demand points. We are given a set W of ordered node pairs $w = (i, j)$. The nodes i and j are referred to as the *origin* and the *destination* of w , respectively, and w is referred to as an OD pair. For each w , we are given a scalar r_w referred to as the *input* of w . For example, in the context of routing of data in a communication network, r_w (measured in data units/second) is the arrival rate of traffic entering and exiting the network at the origin and the destination of w , respectively. The objective is to divide each r_w among the many paths from origin to destination in a way that the resulting total arc flow pattern minimizes a suitable cost function.

We denote:

P_w : A given set of paths that start at the origin and end at the destination of w . All arcs on each of these paths are oriented in the direction from the origin to the destination.

x_p : The portion of r_w assigned to path p , also called the *flow of path p* .

The collection of all path flows $\{x_p \mid p \in P_w, w \in W\}$ must satisfy the constraints

$$\sum_{p \in P_w} x_p = r_w, \quad \forall w \in W, \quad (1.39)$$

$$x_p \geq 0, \quad \forall p \in P_w, w \in W. \quad (1.40)$$

The total flow F_{ij} of arc (i, j) is the sum of all path flows traversing the arc:

$$F_{ij} = \sum_{\substack{\text{all paths } p \\ \text{containing } (i, j)}} x_p. \quad (1.41)$$

Consider a cost function of the form

$$\sum_{(i, j)} D_{ij}(F_{ij}). \quad (1.42)$$

The problem is to find a set of path flows $\{x_p\}$ that minimize this cost function subject to the constraints of Eqs. (1.39)-(1.41). It is typically assumed that D_{ij} is a convex function of F_{ij} . In data routing applications, the form of D_{ij} is often based on a queueing model of average delay, in which case D_{ij} is continuously differentiable within its domain (see e.g., [BeG92]). In a related context, arising in optical networks, the problem involves additional integer constraints on x_p , but may be addressed as a problem with continuous flow variables (see [OzB03]).

The preceding problem is known as a *multicommodity network flow problem*. The terminology reflects the fact that the arc flows consist of several different commodities; in the present example the different commodities are the data of the distinct OD pairs. This problem also arises in essentially identical form in traffic network equilibrium problems (see e.g., [FIH95], [Ber98], [Ber99], [Pat99], [Pat04]). The special case where all OD pairs have the same end node, or all OD pairs have the same start node, is known as the *single commodity network flow problem*, a much easier type of problem, for which there are efficient specialized algorithms that tend to be much faster than their multicommodity counterparts (see textbooks such as [Ber91], [Ber98]).

By expressing the total flows F_{ij} in terms of the path flows in the cost function (1.42) [using Eq. (1.41)], the problem can be formulated in terms of the path flow variables $\{x_p \mid p \in P_w, w \in W\}$ as

$$\begin{aligned} & \text{minimize} && D(x) \\ & \text{subject to} && \sum_{p \in P_w} x_p = r_w, \quad \forall w \in W, \\ & && x_p \geq 0, \quad \forall p \in P_w, w \in W, \end{aligned}$$

where

$$D(x) = \sum_{(i,j)} D_{ij} \left(\sum_{\substack{\text{all paths } p \\ \text{containing } (i,j)}} x_p \right)$$

and x is the vector of path flows x_p . There is a potentially huge number of variables as well as constraints in this problem. However, by judiciously taking into account the special structure of the problem, the constraint set can be simplified and approximated by the convex hull of a small number of vectors x , and the number of variables and constraints can be reduced to a manageable size (see e.g., [BeG83], [FIH95], [OMV00], and our discussion in Section 4.2).

There are several approaches to handle a large number of constraints. One possibility, which points the way to some major classes of algorithms, is to initially discard some of the constraints, solve the corresponding less constrained problem, and later selectively reintroduce constraints that seem to be violated at the optimum. In Chapters 4-6, we will discuss methods of this type in some detail.

Another possibility is to replace constraints with penalties that assign high cost for their violation. In particular, we may replace problem (1.35) with

$$\begin{aligned} & \text{minimize} && f(x) + c \sum_{j=1}^r P(g_j(x)) \\ & \text{subject to} && x \in X, \end{aligned}$$

where $P(\cdot)$ is a scalar penalty function satisfying $P(u) = 0$ if $u \leq 0$, and $P(u) > 0$ if $u > 0$, and c is a positive penalty parameter. We discuss this possibility in the next section.

1.5 EXACT PENALTY FUNCTIONS

In this section we discuss a transformation that is often useful in the context of constrained optimization algorithms. We will derive a form of equivalence between a constrained convex optimization problem, and a penalized problem that is less constrained or is entirely unconstrained. The motivation is that some convex optimization algorithms do not have constrained counterparts, but can be applied to a penalized unconstrained problem. Furthermore, in some analytical contexts, it is useful to be able to work with an equivalent problem that is less constrained.

We consider the convex programming problem

$$\begin{aligned} & \text{minimize} && f(x) \\ & \text{subject to} && x \in X, \quad g_j(x) \leq 0, \quad j = 1, \dots, r, \end{aligned} \quad (1.43)$$

where X is a convex subset of \Re^n , and $f : X \rightarrow \Re$ and $g_j : X \rightarrow \Re$ are given convex functions. We denote by f^* the primal optimal value, and by q^* the dual optimal value, i.e.,

$$q^* = \sup_{\mu \geq 0} q(\mu),$$

where

$$q(\mu) = \inf_{x \in X} \{f(x) + \mu'g(x)\}, \quad \forall \mu \geq 0,$$

with $g(x) = (g_1(x), \dots, g_r(x))'$. We assume that $-\infty < q^* = f^* < \infty$.

We introduce a convex penalty function $P : \Re^r \mapsto \Re$, which satisfies

$$P(u) = 0, \quad \forall u \leq 0, \quad (1.44)$$

$$P(u) > 0, \quad \text{if } u_j > 0 \text{ for some } j = 1, \dots, r. \quad (1.45)$$

We consider solving in place of the original problem (1.43), the “penalized” problem

$$\begin{aligned} & \text{minimize} && f(x) + P(g(x)) \\ & \text{subject to} && x \in X, \end{aligned} \quad (1.46)$$

where the inequality constraints have been replaced by the extra cost $P(g(x))$ for their violation. Some interesting examples of penalty functions are based on the squared or the absolute value of constraint violation:

$$P(u) = \frac{c}{2} \sum_{j=1}^r (\max\{0, u_j\})^2,$$

and

$$P(u) = c \sum_{j=1}^r \max\{0, u_j\},$$

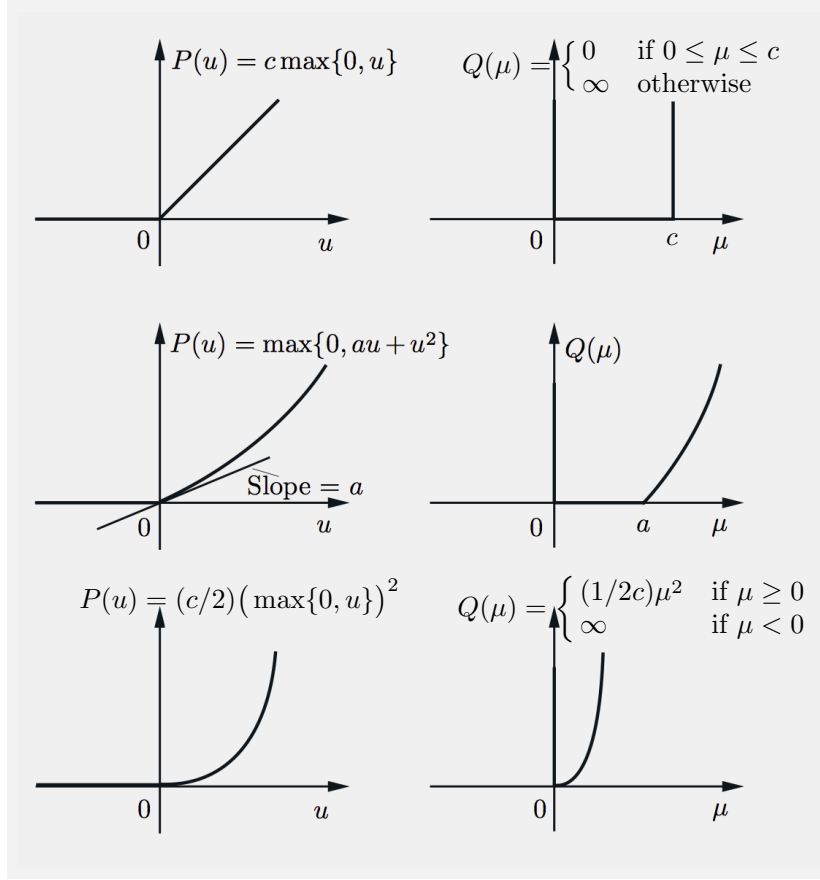


Figure 1.5.1. Illustration of various penalty functions P and their conjugate functions, denoted by Q . Because $P(u) = 0$ for $u \leq 0$, we have $Q(\mu) = \infty$ for μ outside the nonnegative orthant.

where c is a positive penalty parameter. However, there are other possibilities that may be well-matched with the problem at hand.

The conjugate function of P is given by

$$Q(\mu) = \sup_{u \in \mathbb{R}^r} \{u' \mu - P(u)\},$$

and it can be seen that

$$Q(\mu) \geq 0, \quad \forall \mu \in \mathbb{R}^r,$$

$$Q(\mu) = \infty, \quad \text{if } \mu_j < 0 \text{ for some } j = 1, \dots, r.$$

Figure 1.5.1 shows some examples of one-dimensional penalty functions P , together with their conjugates.

Consider the primal function of the original constrained problem,

$$p(u) = \inf_{x \in X, g(x) \leq u} f(x), \quad u \in \mathbb{R}^r.$$

We have,

$$\begin{aligned} \inf_{x \in X} \{f(x) + P(g(x))\} &= \inf_{x \in X} \inf_{u \in \mathbb{R}^r, g(x) \leq u} \{f(x) + P(g(x))\} \\ &= \inf_{x \in X} \inf_{u \in \mathbb{R}^r, g(x) \leq u} \{f(x) + P(u)\} \\ &= \inf_{x \in X, u \in \mathbb{R}^r, g(x) \leq u} \{f(x) + P(u)\} \\ &= \inf_{u \in \mathbb{R}^r} \inf_{x \in X, g(x) \leq u} \{f(x) + P(u)\} \\ &= \inf_{u \in \mathbb{R}^r} \{p(u) + P(u)\}, \end{aligned}$$

where for the second equality, we use the monotonicity relation[†]

$$u \leq v \quad \Rightarrow \quad P(u) \leq P(v).$$

Moreover, $-\infty < q^*$ and $f^* < \infty$ by assumption, and since for any μ with $q(\mu) > -\infty$, we have

$$p(u) \geq q(\mu) - \mu'u > -\infty, \quad \forall u \in \mathbb{R}^r,$$

it follows that $p(0) < \infty$ and $p(u) > -\infty$ for all $u \in \mathbb{R}^r$, so p is proper.

We can now apply the Fenchel Duality Theorem (Prop. 1.2.1) with the identifications $f_1 = p$, $f_2 = P$, and $A = I$. We use the conjugacy relation between the primal function p and the dual function q to write

$$\inf_{u \in \mathbb{R}^r} \{p(u) + P(u)\} = \sup_{\mu \geq 0} \{q(\mu) - Q(\mu)\}, \quad (1.47)$$

so that

$$\inf_{x \in X} \{f(x) + P(g(x))\} = \sup_{\mu \geq 0} \{q(\mu) - Q(\mu)\}; \quad (1.48)$$

see Fig. 1.5.2. Note that the conditions for application of the theorem are satisfied since the penalty function P is real-valued, so that the relative

[†] To show this relation, we argue by contradiction. If there exist u and v with $u \leq v$ and $P(u) > P(v)$, then by continuity of P , there must exist \bar{u} close enough to u such that $\bar{u} < v$ and $P(\bar{u}) > P(v)$. Since P is convex, it is monotonically increasing along the halfline $\{\bar{u} + \alpha(\bar{u} - v) \mid \alpha \geq 0\}$, and since $P(\bar{u}) > P(v) \geq 0$, P takes positive values along this halfline. However, since $\bar{u} < v$, this halfline eventually enters the negative orthant, where P takes the value 0 by Eq. (1.44), a contradiction.

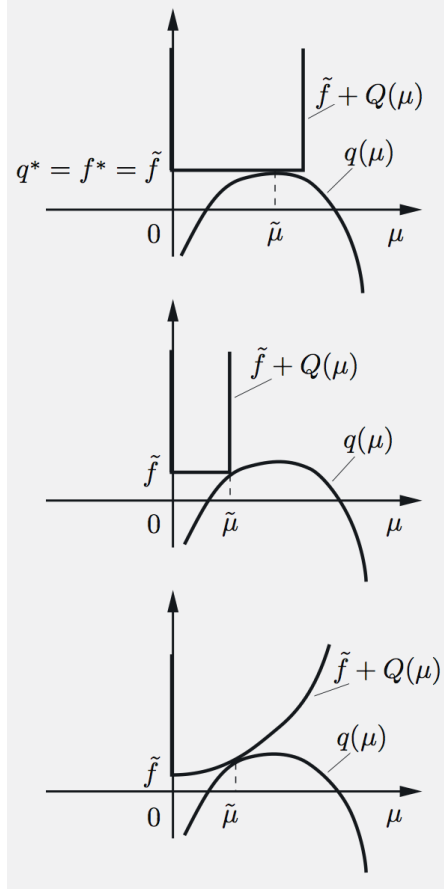


Figure 1.5.2. Illustration of the duality relation (1.48), and the optimal values of the penalized and the dual problem. Here f^* is the optimal value of the original problem, which is assumed to be equal to the optimal dual value q^* , while \tilde{f} is the optimal value of the penalized problem,

$$\tilde{f} = \inf_{x \in X} \{f(x) + P(g(x))\}.$$

The point of contact of the graphs of the functions $\tilde{f} + Q(\mu)$ and $q(\mu)$ corresponds to the vector $\tilde{\mu}$ that attains the maximum in the relation

$$\tilde{f} = \max_{\mu \geq 0} \{q(\mu) - Q(\mu)\}.$$

interiors of $\text{dom}(p)$ and $\text{dom}(P)$ have nonempty intersection. Furthermore, as part of the conclusions of part (a) of the Fenchel Duality Theorem, it follows that the supremum over $\mu \geq 0$ in Eq. (1.48) is attained.

Figure 1.5.2 suggests that in order for the penalized problem (1.46) to have the same optimal value as the original constrained problem (1.43), the conjugate Q must be “sufficiently flat” so that it is minimized by some dual optimal solution μ^* . This can be interpreted in terms of properties of subgradients, which are stated in Appendix B, Section 5.4: we must have $0 \in \partial Q(\mu^*)$ for some dual optimal solution μ^* , which by Prop. 5.4.3 in Appendix B, is equivalent to $\mu^* \in \partial P(0)$. This is part (a) of the following proposition, which was given in [Ber75a]. Parts (b) and (c) of the proposition deal with issues of equality of corresponding optimal solutions. The proposition assumes the convexity and other assumptions made in the early part in this section regarding problem (1.43) and the penalty function P .

Proposition 1.5.1: Consider problem (1.43), where we assume that $-\infty < q^* = f^* < \infty$.

- (a) The penalized problem (1.46) and the original constrained problem (1.43) have equal optimal values if and only if there exists a dual optimal solution μ^* such that $\mu^* \in \partial P(0)$.
- (b) In order for some optimal solution of the penalized problem (1.46) to be an optimal solution of the constrained problem (1.43), it is necessary that there exists a dual optimal solution μ^* such that

$$u' \mu^* \leq P(u), \quad \forall u \in \mathbb{R}^r. \quad (1.49)$$

- (c) In order for the penalized problem (1.46) and the constrained problem (1.43) to have the same set of optimal solutions, it is sufficient that there exists a dual optimal solution μ^* such that

$$u' \mu^* < P(u), \quad \forall u \in \mathbb{R}^r \text{ with } u_j > 0 \text{ for some } j. \quad (1.50)$$

Proof: (a) We have using Eqs. (1.47) and (1.48),

$$p(0) \geq \inf_{u \in \mathbb{R}^r} \{p(u) + P(u)\} = \sup_{\mu \geq 0} \{q(\mu) - Q(\mu)\} = \inf_{x \in X} \{f(x) + P(g(x))\}. \quad (1.51)$$

Since $f^* = p(0)$, we have

$$f^* = \inf_{x \in X} \{f(x) + P(g(x))\}$$

if and only if equality holds in Eq. (1.51). This is true if and only if

$$0 \in \arg \min_{u \in \mathbb{R}^r} \{p(u) + P(u)\},$$

which by Prop. 5.4.7 in Appendix B, is true if and only if there exists some $\mu^* \in -\partial p(0)$ with $\mu^* \in \partial P(0)$ (in view of the fact that P is real-valued). Since the set of dual optimal solutions is $-\partial p(0)$ (under our assumption $-\infty < q^* = f^* < \infty$; see Example 5.4.2, [Ber09]), the result follows.

(b) If x^* is an optimal solution of both problems (1.43) and (1.46), then by feasibility of x^* , we have $P(g(x^*)) = 0$, so these two problems have equal optimal values. From part (a), there must exist a dual optimal solution $\mu^* \in \partial P(0)$, which is equivalent to Eq. (1.49), by the subgradient inequality.

(c) If x^* is an optimal solution of the constrained problem (1.43), then $P(g(x^*)) = 0$, so we have

$$f^* = f(x^*) = f(x^*) + P(g(x^*)) \geq \inf_{x \in X} \{f(x) + P(g(x))\}.$$

The condition (1.50) implies the condition (1.49), so that by part (a), equality holds throughout in the above relation, showing that x^* is also an optimal solution of the penalized problem (1.46).

Conversely, let $x^* \in X$ be an optimal solution of the penalized problem (1.46). If x^* is feasible [i.e., satisfies in addition $g(x^*) \leq 0$], then it is an optimal solution of the constrained problem (1.43) [since $P(g(x)) = 0$ for all feasible vectors x], and we are done. Otherwise x^* is infeasible in which case $g_j(x^*) > 0$ for some j . Then, by using the given condition (1.50), it follows that there exists a dual optimal solution μ^* and an $\epsilon > 0$ such that

$$\mu^{*\prime} g(x^*) + \epsilon < P(g(x^*)).$$

Let \tilde{x} be a feasible vector such that $f(\tilde{x}) \leq f^* + \epsilon$. Since $P(g(\tilde{x})) = 0$ and $f^* = \min_{x \in X} \{f(x) + \mu^{*\prime} g(x)\}$, we obtain

$$f(\tilde{x}) + P(g(\tilde{x})) = f(\tilde{x}) \leq f^* + \epsilon \leq f(x^*) + \mu^{*\prime} g(x^*) + \epsilon.$$

By combining the last two relations, it follows that

$$f(\tilde{x}) + P(g(\tilde{x})) < f(x^*) + P(g(x^*)),$$

which contradicts the hypothesis that x^* is an optimal solution of the penalized problem (1.46). This completes the proof. **Q.E.D.**

As an illustration, consider the minimization of $f(x) = -x$ over all $x \in X = \{x \mid x \geq 0\}$ with $g(x) = x \leq 0$. The dual function is

$$q(\mu) = \inf_{x \geq 0} (\mu - 1)x, \quad \mu \geq 0,$$

so $q(\mu) = 0$ for $\mu \in [1, \infty)$ and $q(\mu) = -\infty$ otherwise. Let $P(u) = c \max\{0, u\}$, so the penalized problem is $\min_{x \geq 0} \{-x + c \max\{0, x\}\}$. Then parts (a) and (b) of the proposition apply if $c \geq 1$. However, part (c) applies only if $c > 1$. In terms of Fig. 1.5.2, the conjugate of P is $Q(\mu) = 0$ if $\mu \in [0, c]$ and $Q(\mu) = \infty$ otherwise, so when $c = 1$, Q is “flat” over an area not including an interior point of the dual optimal solution set $[1, \infty)$.

To elaborate on the idea of the preceding example, let

$$P(u) = c \sum_{j=1}^r \max\{0, u_j\},$$

where $c > 0$. The condition $\mu^* \in \partial P(0)$, or equivalently,

$$u' \mu^* \leq P(u), \quad \forall u \in \Re^r$$

[cf. Eq. (1.49)], is equivalent to

$$\mu_j^* \leq c, \quad \forall j = 1, \dots, r.$$

Similarly, the condition $u'\mu^* < P(u)$ for all $u \in \Re^r$ with $u_j > 0$ for some j [cf. Eq. (1.50)], is equivalent to

$$\mu_j^* < c, \quad \forall j = 1, \dots, r.$$

The reader may consult the literature for other results on exact penalty functions, starting with their first proposal in the book [Zan69]. The preceding development is based on [Ber75], and focuses on convex programming problems. For additional representative references, some of which also discuss nonconvex problems, see [HaM79], [Ber82a], [Bur91], [FeM91], [BNO03], [FrT07]. In what follows we develop an exact penalty function result for the case of an abstract constraint set, which will be used in the context of incremental constraint projection algorithms in Section 6.4.4.

A Distance-Based Exact Penalty Function

Let us discuss the case of a general Lipschitz continuous (not necessarily convex) cost function and an abstract constraint set $X \subset \Re^n$. The idea is to use a penalty that is proportional to the distance from X :

$$\text{dist}(x; X) = \inf_{y \in X} \|x - y\|.$$

The next proposition from [Ber11] provides the basic result (see Fig. 1.5.3).

Proposition 1.5.2: Let $f : \Re^n \mapsto \Re$ be a function that is Lipschitz continuous with constant L over a set $Y \subset \Re^n$, i.e.,

$$|f(x) - f(y)| \leq L\|x - y\|, \quad \forall x, y \in Y.$$

Let also X be a nonempty closed subset of Y , and let c be a scalar such that $c > L$. Then x^* minimizes f over X if and only if x^* minimizes

$$F_c(x) = f(x) + c \text{dist}(x; X)$$

over Y .

Proof: For any $x \in Y$, let \hat{x} denote a vector of X that is at minimum distance from x (such a vector exists by the closure of X and Weierstrass' Theorem). If $c > L$, we have for all $x \in Y$,

$$\begin{aligned} F(x) &= f(x) + c\|x - \hat{x}\| \\ &= f(\hat{x}) + (f(x) - f(\hat{x})) + c\|x - \hat{x}\| \\ &\geq f(\hat{x}) + (c - L)\|x - \hat{x}\| \\ &\geq F(\hat{x}), \end{aligned}$$

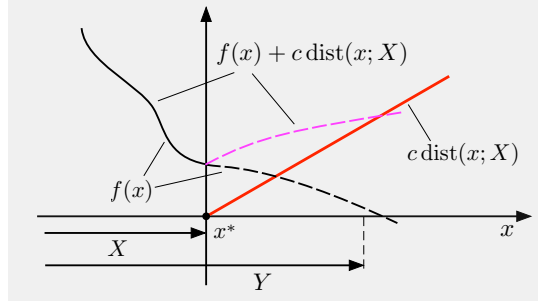


Figure 1.5.3. Illustration of Prop. 1.5.2. For c greater than the Lipschitz constant of f , the “slope” of the penalty function counteracts the “slope” of f at the optimal solution x^* .

with strict inequality if $x \neq \hat{x}$; here the first inequality follows using the Lipschitz property of f to write

$$f(x) - f(\hat{x}) \geq -L\|x - \hat{x}\|,$$

while the second inequality follows from the fact $f(\hat{x}) = F(\hat{x})$. In words, the value of $F(x)$ is strictly reduced when we project an $x \in Y$ with $x \notin X$ onto X . Hence the minima of F over Y can only lie within X , while $F = f$ within X . Thus all minima of F over Y must lie in X and also minimize f over X (since $F = f$ on X). Conversely, all minima of f over X are also minima of F over X (since $F = f$ on X), and by the preceding inequality, they are also minima of F over Y . **Q.E.D.**

The following proposition provides a generalization for constraints that involve the intersection of several sets.

Proposition 1.5.3: Let $f : Y \mapsto \Re$ be a function defined on a subset Y of \Re^n , and let X_i , $i = 1, \dots, m$, be closed subsets of Y with nonempty intersection. Assume that f is Lipschitz continuous over Y with constant L , and that for some scalar $\beta > 0$, we have

$$\text{dist}(x; X_1 \cap \dots \cap X_m) \leq \beta \sum_{i=1}^m \text{dist}(x; X_i), \quad \forall x \in Y. \quad (1.52)$$

Let c be a scalar such that $c > \beta L$. Then the set of minima of f over $\cap_{i=1}^m X_i$ coincides with the set of minima of

$$f(x) + c \sum_{i=1}^m \text{dist}(x; X_i)$$

over Y .

Proof: The proof is similar to the proof of Prop. 1.5.2, using Eq. (1.52) to modify the main inequality. Denote $F(x) = f(x) + c \sum_{i=1}^m \text{dist}(x; X_i)$ and $X = X_1 \cap \cdots \cap X_m$. For a vector $x \in Y$, let \hat{x}_i denote a vector of X_i that is at minimum distance from x , and let \hat{x} denote a vector of X that is at minimum distance from x . If $c > \beta L$, we have for all $x \in Y$,

$$\begin{aligned} F(x) &= f(x) + c \sum_{i=1}^m \|x - \hat{x}_i\| \\ &\geq f(\hat{x}) + (f(x) - f(\hat{x})) + \frac{c}{\beta} \|x - \hat{x}\| \\ &\geq f(\hat{x}) + \left(\frac{c}{\beta} - L \right) \|x - \hat{x}\| \\ &\geq F(\hat{x}), \end{aligned}$$

with strict inequality if $x \neq \hat{x}$. The proof now proceeds as in Prop. 1.5.2. **Q.E.D.**

It can be shown that the condition (1.52) is satisfied if all the sets X_1, \dots, X_m are polyhedral (this is a consequence of the well-known Hoffman's Lemma). We finally note that exact penalty functions, and particularly the distance function $\text{dist}(x; X_i)$, are often relatively convenient in various contexts where difficult constraints complicate the algorithmic solution. As an example, see Section 6.4.4, where incremental proximal methods for highly constrained problems are discussed.

1.6 NOTES, SOURCES, AND EXERCISES

There is a very extensive literature on convex optimization, and in this section we will restrict ourselves to noting some books, research monographs, and surveys. In subsequent chapters, we will discuss in greater detail the literature that relates to the specialized content of these chapters.

Books relating primarily to duality theory are Rockafellar [Roc70], Stoer and Witzgall [StW70], Ekeland and Temam [EkT76], Bonnans and Shapiro [BoS00], Zalinescu [Zal02], Auslender and Teboulle [AuT03], and Bertsekas [Ber09].

The books by Rockafellar and Wets [RoW98], Borwein and Lewis [BoL00], and Bertsekas, Nedić, and Ozdaglar [BNO03] straddle the boundary between convex and variational analysis, a broad spectrum of topics that integrate classical analysis, convexity, and optimization of both convex and nonconvex (possibly nonsmooth) functions.

The book by Hiriart-Urruty and Lemarechal [HiL93] focuses on convex optimization algorithms. The books by Rockafellar [Roc84] and Bertsekas [Ber98] have a more specialized focus on network optimization algorithms and monotropic programming problems, which will be discussed in

Chapters 4 and 6. The book by Ben-Tal and Nemirovski [BeN01] focuses on conic and semidefinite programming [see also the 2005 class notes by Nemirovski (on line), and the representative survey papers by Alizadeh and Goldfarb [AlG03], and Todd [Tod01]]. The book by Wolkowicz, Saigal, and Vanderberghe [WSV00] contains a collection of survey articles on semidefinite programming. The book by Boyd and Vanderberghe [BoV04] describes many applications, and contains a lot of related material and references. The book by Ben-Tal, El Ghaoui, and Nemirovski [BGN09] focuses on robust optimization; see also the survey by Bertsimas, Brown, and Caramanis [BBC11]. The book by Bauschke and Combettes [BaC11] develops the connection of convex analysis and monotone operator theory in infinite dimensional spaces. The book by Rockafellar and Wets [RoW98] also has a substantial finite-dimensional treatment of this subject. The books by Cottle, Pang, and Stone [CPS92], and Facchinei and Pang [FaP03] focus on complementarity and variational inequality problems. The books by Palomar and Eldar [PaE10], and Vetterli, Kovacevic, and Goyal [VKG14], and the surveys in the May 2010 issue of the IEEE Signal Processing Magazine describe applications of convex optimization in communications and signal processing. The books by Hastie, Tibshirani, and Friedman [HTF09], and Sra, Nowozin, and Wright [SNW12] describe applications of convex optimization in machine learning.

E X E R C I S E S

1.1 (Support Vector Machines and Duality)

Consider the classification problem associated with a support vector machine,

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2}\|x\|^2 + \beta \sum_{i=1}^m \max\{0, 1 - b_i(c'_i x + y)\} \\ & \text{subject to} \quad x \in \mathbb{R}^n, \quad y \in \mathbb{R}, \end{aligned}$$

with quadratic regularization, where β is a positive regularization parameter (cf. Example 1.3.3).

(a) Write the problem in the equivalent form

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2}\|x\|^2 + \beta \sum_{i=1}^m \xi_i \\ & \text{subject to} \quad x \in \mathbb{R}^n, \quad y \in \mathbb{R}, \\ & \quad \quad \quad 0 \leq \xi_i, \quad 1 - b_i(c'_i x + y) \leq \xi_i, \quad i = 1, \dots, m. \end{aligned}$$

Associate dual variables $\mu_i \geq 0$ with the constraints $1 - b_i(c'_i x + y) \leq \xi_i$, and show that the dual function is given by

$$q(\mu) = \begin{cases} \hat{q}(\mu) & \text{if } \sum_{j=1}^m \mu_j b_j = 0, \quad 0 \leq \mu_i \leq \beta, \quad i = 1, \dots, m, \\ -\infty & \text{otherwise,} \end{cases}$$

where

$$\hat{q}(\mu) = \sum_{i=1}^m \mu_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m b_i b_j c'_i c'_j \mu_i \mu_j.$$

Does the dual problem, viewed as the equivalent quadratic program

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m b_i b_j c'_i c'_j \mu_i \mu_j - \sum_{i=1}^m \mu_i \\ & \text{subject to} && \sum_{j=1}^m \mu_j b_j = 0, \quad 0 \leq \mu_i \leq \beta, \quad i = 1, \dots, m, \end{aligned}$$

always have a solution? Is the solution unique? *Note:* The dual problem may have high dimension, but it has a generally more favorable structure than the primal. The reason is the simplicity of its constraint set, which makes it suitable for special types of quadratic programming methods, and the two-metric projection and coordinate descent methods of Section 2.1.2.

- (b) Consider an alternative formulation where the variable y is set to 0, leading to the problem

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \|x\|^2 + \beta \sum_{i=1}^m \max\{0, 1 - b_i c'_i x\} \\ & \text{subject to} && x \in \mathbb{R}^n. \end{aligned}$$

Show that the dual problem should be modified so that the constraint $\sum_{j=1}^m \mu_j b_j = 0$ is not present, thus leading to a bound-constrained quadratic dual problem.

Note: The literature of the support vector machine field is extensive. Many of the nondifferentiable optimization methods to be discussed in subsequent chapters have been applied in connection to this field; see e.g., [MaM01], [FeM02], [SmS04], [Bot05], [Joa06], [JFY09], [JoY09], [SSS07], [LeW11].

1.2 (Minimizing the Sum or the Maximum of Norms [LVB98])

Consider the problems

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^p \|F_i x + g_i\| \\ & \text{subject to} && x \in \mathbb{R}^n, \end{aligned} \tag{1.53}$$

and

$$\begin{aligned} & \text{minimize} && \max_{i=1, \dots, p} \|F_i x + g_i\| \\ & \text{subject to} && x \in \mathbb{R}^n, \end{aligned}$$

where F_i and g_i are given matrices and vectors, respectively. Convert these problems to second order cone form and derive the corresponding dual problems.

1.3 (Complex l_1 and l_∞ Approximation [LVB98])

Consider the complex l_1 approximation problem

$$\begin{aligned} & \text{minimize} && \|Ax - b\|_1 \\ & \text{subject to} && x \in \mathcal{C}^n, \end{aligned}$$

where \mathcal{C}^n is the set of n -dimensional vectors whose components are complex numbers, and A and b are given matrix and vector with complex components. Show that it is a special case of problem (1.53) and derive the corresponding dual problem. Repeat for the complex l_∞ approximation problem

$$\begin{aligned} & \text{minimize} && \|Ax - b\|_\infty \\ & \text{subject to} && x \in \mathcal{C}^n. \end{aligned}$$

1.4

The purpose of this exercise is to show that the SOCP can be viewed as a special case of SDP.

- (a) Show that a vector $x \in \mathbb{R}^n$ belongs to the second order cone if and only if the matrix

$$x_n I + \begin{pmatrix} 0 & 0 & \cdots & 0 & x_1 \\ 0 & 0 & \cdots & 0 & x_2 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 0 & x_{n-1} \\ x_1 & x_2 & \cdots & x_{n-1} & 0 \end{pmatrix}$$

is positive semidefinite. *Hint:* We have that for any positive definite symmetric $n \times n$ matrix A , vector $b \in \mathbb{R}^n$, and scalar d , the matrix

$$\begin{pmatrix} A & b \\ b' & c \end{pmatrix}$$

is positive definite if and only if

$$c - b' A^{-1} b > 0.$$

- (b) Use part (a) to show that the primal SOCP can be written in the form of the dual SDP.

1.5 (Explicit Form of a Second Order Cone Problem)

Consider the SOCP (1.24).

- (a) Partition the $n_i \times (n+1)$ matrices $(A_i \ b_i)$ as

$$(A_i \ b_i) = \begin{pmatrix} D_i & d_i \\ p_i' & q_i \end{pmatrix}, \quad i = 1, \dots, m,$$

where D_i is an $(n_i - 1) \times n$ matrix, $d_i \in \mathbb{R}^{n_i - 1}$, $p_i \in \mathbb{R}^n$, and $q_i \in \mathbb{R}$. Show that

$$A_i x - b_i \in C_i \quad \text{if and only if} \quad \|D_i x - d_i\| \leq p_i' x - q_i,$$

so we can write the SOCP (1.24) as

$$\begin{aligned} & \text{minimize} && c'x \\ & \text{subject to} && \|D_i x - d_i\| \leq p_i' x - q_i, \quad i = 1, \dots, m. \end{aligned}$$

(b) Similarly partition λ_i as

$$\lambda_i = \begin{pmatrix} \mu_i \\ \nu_i \end{pmatrix}, \quad i = 1, \dots, m,$$

where $\mu_i \in \mathbb{R}^{n_i - 1}$ and $\nu_i \in \mathbb{R}$. Show that the dual problem (1.25) can be written in the form

$$\begin{aligned} & \text{maximize} && \sum_{i=1}^m (d_i' \mu_i + q_i \nu_i) \\ & \text{subject to} && \sum_{i=1}^m (D_i' \mu_i + \nu_i p_i) = c, \quad \|\mu_i\| \leq \nu_i, \quad i = 1, \dots, m. \end{aligned}$$

(c) Show that the primal and dual interior point conditions for strong duality (Prop. 1.2.4) hold if there exist primal and dual feasible solutions \bar{x} and $(\bar{\mu}_i, \bar{\nu}_i)$ such that

$$\|D_i \bar{x} - d_i\| < p_i' \bar{x} - q_i, \quad i = 1, \dots, m,$$

and

$$\|\bar{\mu}_i\| < \bar{\nu}_i, \quad i = 1, \dots, m,$$

respectively.

1.6 (Separable Conic Problems)

Consider the problem

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^m f_i(x_i) \\ & \text{subject to} && x \in S \cap C, \end{aligned}$$

where $x = (x_1, \dots, x_m)$ with $x_i \in \mathbb{R}^{n_i}$, $i = 1, \dots, m$, and $f_i : \mathbb{R}^{n_i} \mapsto (-\infty, \infty]$ is a proper convex function for each i , and S and C are a subspace and a cone of $\mathbb{R}^{n_1 + \dots + n_m}$, respectively. Show that a dual problem is

$$\begin{aligned} & \text{maximize} && \sum_{i=1}^m q_i(\lambda_i) \\ & \text{subject to} && \lambda \in \hat{C} + S^\perp, \end{aligned}$$

where $\lambda = (\lambda_1, \dots, \lambda_m)$, \hat{C} is the dual cone of C , and

$$q_i(\lambda_i) = \inf_{z_i \in \mathfrak{R}} \{f_i(z_i) - \lambda'_i z_i\}, \quad i = 1, \dots, m.$$

1.7 (Weber Points)

Consider the problem of finding a circle of minimum radius that contains r points y_1, \dots, y_r in the plane, i.e., find x and z that minimize z subject to $\|x - y_j\| \leq z$ for all $j = 1, \dots, r$, where x is the center of the circle under optimization.

- (a) Introduce multipliers μ_j , $j = 1, \dots, r$, for the constraints, and show that the dual problem has an optimal solution and there is no duality gap.
- (b) Show that calculating the dual function at some $\mu \geq 0$ involves the computation of a Weber point of y_1, \dots, y_r with weights μ_1, \dots, μ_r , i.e., the solution of the problem

$$\min_{x \in \mathfrak{R}^2} \sum_{j=1}^r \mu_j \|x - y_j\|$$

(see Example 1.3.7).

1.8 (Inconsistent Convex Systems of Inequalities)

Let $g_j : \mathfrak{R}^n \mapsto \mathfrak{R}$, $j = 1, \dots, r$, be convex functions over the nonempty convex set $X \subset \mathfrak{R}^n$. Show that the system

$$g_j(x) < 0, \quad j = 1, \dots, r,$$

has no solution within X if and only if there exists a vector $\mu \in \mathfrak{R}^r$ such that

$$\sum_{j=1}^r \mu_j = 1, \quad \mu \geq 0,$$

$$\mu' g(x) \geq 0, \quad \forall x \in X.$$

Note: This is an example of what is known as a *theorem of the alternative*. There are many results of this type, with a long history, such as the Farkas Lemma, and the theorems of Gordan, Motzkin, and Stiemke, which address the feasibility (possibly strict) of linear inequalities. They can be found in many sources, including Section 5.6 of [Ber09]. *Hint:* Consider the convex program

$$\begin{aligned} &\text{minimize } y \\ &\text{subject to } x \in X, \quad y \in \mathfrak{R}, \quad g_j(x) \leq y, \quad j = 1, \dots, r. \end{aligned}$$

References

- [ACH97] Auslender, A., Cominetti, R., and Haddou, M., 1997. “Asymptotic Analysis for Penalty and Barrier Methods in Convex and Linear Programming,” *Math. of Operations Research*, Vol. 22, pp. 43-62.
- [ALS14] Abernethy, J., Lee, C., Sinha, A., and Tewari, A., 2014. “Online Linear Optimization via Smoothing,” *arXiv preprint arXiv:1405.6076*.
- [AgB14] Agarwal, A., and Bottou, L., 2014. “A Lower Bound for the Optimization of Finite Sums,” *arXiv preprint arXiv:1410.0723*.
- [AgD11] Agarwal, A., and Duchi, J. C., 2011. “Distributed Delayed Stochastic Optimization,” In *Advances in Neural Information Processing Systems* (NIPS 2011), pp. 873-881.
- [AlG03] Alizadeh, F., and Goldfarb, D., 2003. “Second-Order Cone Programming,” *Math. Programming*, Vol. 95, pp. 3-51.
- [AnH13] Andersen, M. S., and Hansen, P. C., 2013. “Generalized Row-Action Methods for Tomographic Imaging,” *Numerical Algorithms*, Vol. 67, pp. 1-24.
- [Arm66] Armijo, L., 1966. “Minimization of Functions Having Continuous Partial Derivatives,” *Pacific J. Math.*, Vol. 16, pp. 1-3.
- [Ash72] Ash, R. B., 1972. *Real Analysis and Probability*, Academic Press, NY.
- [AtV95] Atkinson, D. S., and Vaidya, P. M., 1995. “A Cutting Plane Algorithm for Convex Programming that Uses Analytic Centers,” *Math. Programming*, Vol. 69, pp. 1-44.
- [AuE76] Aubin, J. P., and Ekeland, I., 1976. “Estimates of the Duality Gap in Nonconvex Optimization,” *Math. of Operations Research*, Vol. 1, pp. 255- 245.
- [AuT03] Auslender, A., and Teboulle, M., 2003. *Asymptotic Cones and Functions in Optimization and Variational Inequalities*, Springer, NY.
- [AuT04] Auslender, A., and Teboulle, M., 2004. “Interior Gradient and Epsilon-Subgradient Descent Methods for Constrained Convex Minimization,” *Math. of Operations Research*, Vol. 29, pp. 1-26.
- [Aus76] Auslender, A., 1976. *Optimization: Methodes Numeriques*, Mason, Paris.
- [Aus92] Auslender, A., 1992. “Asymptotic Properties of the Fenchel Dual

- Functional and Applications to Decomposition Problems,” *J. of Optimization Theory and Applications*, Vol. 73, pp. 427-449.
- [BBC11] Bertsimas, D., Brown, D. B., and Caramanis, C., 2011. “Theory and Applications of Robust Optimization,” *SIAM Review*, Vol. 53, pp. 464-501.
- [BBG09] Bordes, A., Bottou, L., and Gallinari, P., 2009. “SGD-QN: Careful Quasi-Newton Stochastic Gradient Descent,” *J. of Machine Learning Research*, Vol. 10, pp. 1737-1754.
- [BBL07] Berry, M. W., Browne, M., Langville, A. N., Pauca, V. P., and Plemmons, R. J., 2007. “Algorithms and Applications for Approximate Nonnegative Matrix Factorization,” *Computational Statistics and Data Analysis*, Vol. 52, pp. 155-173.
- [BBY12] Borwein, J. M., Burachik, R. S., and Yao, L., 2012. “Conditions for Zero Duality Gap in Convex Programming,” *arXiv preprint arXiv:1211.4953*.
- [BCK06] Bauschke, H. H., Combettes, P. L., and Kruk, S. G., 2006. “Extrapolation Algorithm for Affine-Convex Feasibility Problems,” *Numer. Algorithms*, Vol. 41, pp. 239-274.
- [BGI95] Burachik, R., Grana Drummond, L. M., Iusem, A. N., and Svaiter, B. F., 1995. “Full Convergence of the Steepest Descent Method with Inexact Line Searches,” *Optimization*, Vol. 32, pp. 137-146.
- [BGL06] Bonnans, J. F., Gilbert, J. C., Lemaréchal, C., and Sagastizábal, S. C., 2006. *Numerical Optimization: Theoretical and Practical Aspects*, Springer, NY.
- [BGN09] Ben-Tal, A., El Ghaoui, L., and Nemirovski, A., 2009. *Robust Optimization*, Princeton Univ. Press, Princeton, NJ.
- [BHG08] Blatt, D., Hero, A. O., Gauchman, H., 2008. “A Convergent Incremental Gradient Method with a Constant Step Size,” *SIAM J. Optimization*, Vol. 18, pp. 29-51.
- [BHT87] Bertsekas, D. P., Hossein, P., and Tseng, P., 1987. “Relaxation Methods for Network Flow Problems with Convex Arc Costs,” *SIAM J. on Control and Optimization*, Vol. 25, pp. 1219-1243.
- [BJM12] Bach, F., Jenatton, R., Mairal, J., and Obozinski, G., 2012. “Optimization with Sparsity-Inducing Penalties,” *Foundations and Trends in Machine Learning*, Vol. 4, pp. 1-106.
- [BKM14] Burachik, R. S., Kaya, C. Y., and Majeed, S. N., 2014. “A Duality Approach for Solving Control-Constrained Linear-Quadratic Optimal Control Problems,” *SIAM J. on Control and Optimization*, Vol. 52, pp. 1423-1456.
- [BLY14] Bragin, M. A., Luh, P. B., Yan, J. H., Yu, N., and Stern, G. A., 2014. “Convergence of the Surrogate Lagrangian Relaxation Method,” *J. of Optimization Theory and Applications*, on line.

- [BMN01] Ben-Tal, A., Margalit, T., and Nemirovski, A., 2001. “The Ordered Subsets Mirror Descent Optimization Method and its Use for the Positron Emission Tomography Reconstruction,” in *Inherently Parallel Algorithms in Feasibility and Optimization and their Applications* (D. Butnariu, Y. Censor, and S. Reich, eds.), Elsevier, Amsterdam, Netherlands.
- [BMR00] Birgin, E. G., Martinez, J. M., and Raydan, M., 2000. “Non-monotone Spectral Projected Gradient Methods on Convex Sets,” *SIAM J. on Optimization*, Vol. 10, pp. 1196-1211.
- [BMS99] Boltysanski, V., Martini, H., and Soltan, V., 1999. *Geometric Methods and Optimization Problems*, Kluwer, Boston.
- [BNO03] Bertsekas, D. P., Nedić, A., and Ozdaglar, A. E., 2003. *Convex Analysis and Optimization*, Athena Scientific, Belmont, MA.
- [BOT06] Bertsekas, D. P., Ozdaglar, A. E., and Tseng, P., 2006 “Enhanced Fritz John Optimality Conditions for Convex Programming,” *SIAM J. on Optimization*, Vol. 16, pp. 766-797.
- [BPC11] Boyd, S., Parikh, N., Chu, E., Peleato, B., and Eckstein, J., 2011. *Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers*, Now Publishers Inc, Boston, MA.
- [BPP13] Bhatnagar, S., Prasad, H., and Prashanth, L. A., 2013. *Stochastic Recursive Algorithms for Optimization*, Lecture Notes in Control and Information Sciences, Springer, NY.
- [BPT97a] Bertsekas, D. P., Polymenakos, L. C., and Tseng, P., 1997. “An ϵ -Relaxation Method for Separable Convex Cost Network Flow Problems,” *SIAM J. on Optimization*, Vol. 7, pp. 853-870.
- [BPT97b] Bertsekas, D. P., Polymenakos, L. C., and Tseng, P., 1997. “Epsilon-Relaxation and Auction Methods for Separable Convex Cost Network Flow Problems,” in *Network Optimization*, Pardalos, P. M., Hearn, D. W., and Hager, W. W., (Eds.), *Lecture Notes in Economics and Mathematical Systems*, Springer-Verlag, NY, pp. 103-126.
- [BSL14] Bergmann, R., Steidl, G., Laus, F., and Weinmann, A., 2014. “Second Order Differences of Cyclic Data and Applications in Variational Denoising,” *arXiv preprint arXiv:1405.5349*.
- [BSS06] Bazaraa, M. S., Sherali, H. D., and Shetty, C. M., 2006. *Nonlinear Programming: Theory and Algorithms*, 3rd Edition, Wiley, NY.
- [BST14] Bolte, J., Sabach, S., and Teboulle, M., 2014. “Proximal Alternating Linearized Minimization for Nonconvex and Nonsmooth Problems,” *Math. Programming*, Vol. 146, pp. 1-36.
- [BaB88] Barzilai, J., and Borwein, J. M., 1988. “Two Point Step Size Gradient Methods,” *IMA J. of Numerical Analysis*, Vol. 8, pp. 141-148.
- [BaB96] Bauschke, H. H., and Borwein, J. M., 1996. “On Projection Algorithms for Solving Convex Feasibility Problems,” *SIAM Review*, Vol. 38, pp. 367-426.

- [BaC11] Bauschke, H. H., and Combettes, P. L., 2011. *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, Springer, NY.
- [BaM11] Bach, F., and E. Moulines, E., 2011. “Non-Asymptotic Analysis of Stochastic Approximation Algorithms for Machine Learning,” *Advances in Neural Information Processing Systems (NIPS 2011)*.
- [BaT85] Balas, E., and Toth, P., 1985. “Branch and Bound Methods,” in *The Traveling Salesman Problem*, Lawler, E., Lenstra, J. K., Rinnoy Kan, A. H. G., and Shmoys, D. B., (Eds.), Wiley, NY, pp. 361-401.
- [BaW75] Balinski, M., and Wolfe, P., (Eds.), 1975. *Nondifferentiable Optimization*, Math. Programming Study 3, North-Holland, Amsterdam.
- [Bac14] Bacak, M., 2014. “Computing Medians and Means in Hadamard Spaces,” *arXiv preprint arXiv:1210.2145v3*.
- [Bau01] Bauschke, H. H., 2001. “Projection Algorithms: Results and Open Problems,” in *Inherently Parallel Algorithms in Feasibility and Optimization and their Applications* (D. Butnariu, Y. Censor, and S. Reich, eds.), Elsevier, Amsterdam, Netherlands.
- [BeF12] Becker, S., and Fadili, J., 2012. “A Quasi-Newton Proximal Splitting Method,” in *Advances in Neural Information Processing Systems (NIPS 2012)*, pp. 2618-2626.
- [BeG83] Bertsekas, D. P., and Gafni, E., 1983. “Projected Newton Methods and Optimization of Multicommodity Flows,” *IEEE Trans. Automat. Control*, Vol. AC-28, pp. 1090-1096.
- [BeG92] Bertsekas, D. P., and Gallager, R. G., 1992. *Data Networks*, 2nd Ed., Prentice-Hall, Englewood Cliffs, NJ.
On line at <http://web.mit.edu/dimitrib/www/datanets.html>.
- [BeL88] Becker, S., and LeCun, Y., 1988. “Improving the Convergence of Back-Propagation Learning with Second Order Methods,” in *Proceedings of the 1988 Connectionist Models Summer School*, San Matteo, CA.
- [BeL07] Bengio, Y., and LeCun, Y., 2007. “Scaling Learning Algorithms Towards AI,” *Large-Scale Kernel Machines*, Vol. 34, pp. 1-41.
- [BeM71] Bertsekas, D. P., and Mitter, S. K., 1971. “Steepest Descent for Optimization Problems with Nondifferentiable Cost Functionals,” *Proc. 5th Annual Princeton Confer. Inform. Sci. Systems*, Princeton, NJ, pp. 347-351.
- [BeM73] Bertsekas, D. P., and Mitter, S. K., 1973. “A Descent Numerical Method for Optimization Problems with Nondifferentiable Cost Functionals,” *SIAM J. on Control*, Vol. 11, pp. 637-652.
- [BeN01] Ben-Tal, A., and Nemirovskii, A., 2001. *Lectures on Modern Convex Optimization: Analysis, Algorithms, and Engineering Applications*, SIAM, Philadelphia, PA.
- [BeO02] Bertsekas, D. P., Ozdaglar, A. E., 2002. “Pseudonormality and a Lagrange Multiplier Theory for Constrained Optimization,” *J. of Opti-*

- mization Theory and Applications, Vol. 114, pp. 287-343.
- [BeS82] Bertsekas, D. P., and Sandell, N. R., 1982. "Estimates of the Duality Gap for Large-Scale Separable Nonconvex Optimization Problems," Proc. 1982 IEEE Conf. Decision and Control, pp. 782-785.
- [BeT88] Bertsekas, D. P., and Tseng, P., 1988. "Relaxation Methods for Minimum Cost Ordinary and Generalized Network Flow Problems," Operations Research, Vol. 36, pp. 93-114.
- [BeT89a] Bertsekas, D. P., and Tsitsiklis, J. N., 1989. Parallel and Distributed Computation: Numerical Methods, Prentice-Hall, Englewood Cliffs, NJ; republished by Athena Scientific, Belmont, MA, 1997. On line at <http://web.mit.edu/dimitrib/www/pdc.html>.
- [BeT89b] Ben-Tal, A., and Teboulle, M., 1989. "A Smoothing Technique for Nondifferentiable Optimization Problems," Optimization, Lecture Notes in Mathematics, Vol. 1405, pp. 1-11.
- [BeT91] Bertsekas, D. P., and Tsitsiklis, J. N., 1991. "Some Aspects of Parallel and Distributed Iterative Algorithms - A Survey," Automatica, Vol. 27, pp. 3-21.
- [BeT94a] Bertsekas, D. P., and Tseng, P., 1994. "Partial Proximal Minimization Algorithms for Convex Programming," SIAM J. on Optimization, Vol. 4, pp. 551-572.
- [BeT94b] Bertsekas, D. P., and Tseng, P., 1994. "RELAX-IV: A Faster Version of the RELAX Code for Solving Minimum Cost Flow Problems," Massachusetts Institute of Technology, Laboratory for Information and Decision Systems Report LIDS-P-2276, Cambridge, MA.
- [BeT96] Bertsekas, D. P., and Tsitsiklis, J. N., 1996. Neuro-Dynamic Programming, Athena Scientific, Belmont, MA.
- [BeT97] Bertsimas, D., and Tsitsiklis, J. N., 1997. Introduction to Linear Optimization, Athena Scientific, Belmont, MA.
- [BeT00] Bertsekas, D. P., and Tsitsiklis, J. N., 2000. "Gradient Convergence of Gradient Methods with Errors," SIAM J. on Optimization, Vol. 36, pp. 627-642.
- [BeT03] Beck, A., and Teboulle, M., 2003. "Mirror Descent and Nonlinear Projected Subgradient Methods for Convex Optimization," Operations Research Letters, Vol. 31, pp. 167-175.
- [BeT09a] Beck, A., and Teboulle, M., 2009. "Fast Gradient-Based Algorithms for Constrained Total Variation Image Denoising and Deblurring Problems," IEEE Trans. on Image Processing, Vol. 18, pp. 2419-2434.
- [BeT09b] Beck, A., and Teboulle, M., 2009. "A Fast Iterative Shrinkage-Thresholding Algorithm for Linear Inverse Problems," SIAM J. on Imaging Sciences, Vol. 2, pp. 183-202.
- [BeT10] Beck, A., and Teboulle, M., 2010. "Gradient-Based Algorithms with Applications to Signal-Recovery Problems," in Convex Optimization

- in *Signal Processing and Communications* (Y. Eldar and D. Palomar, eds.), Cambridge Univ. Press, pp. 42-88.
- [BeT13] Beck, A., and Tetrushvili, L., 2013. "On the Convergence of Block Coordinate Descent Type Methods," *SIAM J. on Optimization*, Vol. 23, pp. 2037-2060.
- [BeY09] Bertsekas, D. P., and Yu, H., 2009. "Projected Equation Methods for Approximate Solution of Large Linear Systems," *J. of Computational and Applied Mathematics*, Vol. 227, pp. 27-50.
- [BeY10] Bertsekas, D. P., and Yu, H., 2010. "Asynchronous Distributed Policy Iteration in Dynamic Programming," *Proc. of Allerton Conf. on Communication, Control and Computing*, Allerton Park, Ill, pp. 1368-1374.
- [BeY11] Bertsekas, D. P., and Yu, H., 2011. "A Unifying Polyhedral Approximation Framework for Convex Optimization," *SIAM J. on Optimization*, Vol. 21, pp. 333-360.
- [BeZ97] Ben-Tal, A., and Zibulevsky, M., 1997. "Penalty/Barrier Multiplier Methods for Convex Programming Problems," *SIAM J. on Optimization*, Vol. 7, pp. 347-366.
- [Ben09] Bengio, Y., 2009. "Learning Deep Architectures for AI," *Foundations and Trends in Machine Learning*, Vol. 2, pp. 1-127.
- [Ber71] Bertsekas, D. P., 1971. "Control of Uncertain Systems With a Set-Membership Description of the Uncertainty," Ph.D. Thesis, Dept. of EECS, MIT; may be downloaded from <http://web.mit.edu/dimitrib/www/publ.html>.
- [Ber72] Bertsekas, D. P., 1972. "Stochastic Optimization Problems with Nondifferentiable Cost Functionals with an Application in Stochastic Programming," *Proc. 1972 IEEE Conf. Decision and Control*, pp. 555-559.
- [Ber73] Bertsekas, D. P., 1973. "Stochastic Optimization Problems with Nondifferentiable Cost Functionals," *J. of Optimization Theory and Applications*, Vol. 12, pp. 218-231.
- [Ber75a] Bertsekas, D. P., 1975. "Necessary and Sufficient Conditions for a Penalty Method to be Exact," *Math. Programming*, Vol. 9, pp. 87-99.
- [Ber75b] Bertsekas, D. P., 1975. "Nondifferentiable Optimization via Approximation," *Math. Programming Study 3*, Balinski, M., and Wolfe, P., (Eds.), North-Holland, Amsterdam, pp. 1-25.
- [Ber75c] Bertsekas, D. P., 1975. "Combined Primal-Dual and Penalty Methods for Constrained Optimization," *SIAM J. on Control*, Vol. 13, pp. 521-544.
- [Ber75d] Bertsekas, D. P., 1975. "On the Method of Multipliers for Convex Programming," *IEEE Transactions on Aut. Control*, Vol. 20, pp. 385-388.
- [Ber76a] Bertsekas, D. P., 1976. "On the Goldstein-Levitin-Poljak Gradient Projection Method," *IEEE Trans. Automat. Control*, Vol. 21, pp. 174-184.

- [Ber76b] Bertsekas, D. P., 1976. "Multiplier Methods: A Survey," *Automatica*, Vol. 12, pp. 133-145.
- [Ber76c] Bertsekas, D. P., 1976. "On Penalty and Multiplier Methods for Constrained Optimization," *SIAM J. on Control and Optimization*, Vol. 14, pp. 216-235.
- [Ber77] Bertsekas, D. P., 1977. "Approximation Procedures Based on the Method of Multipliers," *J. Optimization Theory and Applications*, Vol. 23, pp. 487-510.
- [Ber79] Bertsekas, D. P., 1979. "A Distributed Algorithm for the Assignment Problem," Lab. for Information and Decision Systems Working Paper, MIT, Cambridge, MA.
- [Ber81] Bertsekas, D. P., 1981. "A New Algorithm for the Assignment Problem," *Mathematical Programming*, Vol. 21, pp. 152-171.
- [Ber82a] Bertsekas, D. P., 1982. *Constrained Optimization and Lagrange Multiplier Methods*, Academic Press, NY; republished in 1996 by Athena Scientific, Belmont, MA. On line at <http://web.mit.edu/dimitrib/www/-lagrmult.html>.
- [Ber82b] Bertsekas, D. P., 1982. "Projected Newton Methods for Optimization Problems with Simple Constraints," *SIAM J. on Control and Optimization*, Vol. 20, pp. 221-246.
- [Ber82c] Bertsekas, D. P., 1982. "Distributed Dynamic Programming," *IEEE Trans. Aut. Control*, Vol. AC-27, pp. 610-616.
- [Ber83] Bertsekas, D. P., 1983. "Distributed Asynchronous Computation of Fixed Points," *Math. Programming*, Vol. 27, pp. 107-120.
- [Ber85] Bertsekas, D. P., 1985. "A Unified Framework for Primal-Dual Methods in Minimum Cost Network Flow Problems," *Mathematical Programming*, Vol. 32, pp. 125-145.
- [Ber91] Bertsekas, D. P., 1991. *Linear Network Optimization: Algorithms and Codes*, MIT Press, Cambridge, MA.
- [Ber92] Bertsekas, D. P., 1992. "Auction Algorithms for Network Problems: A Tutorial Introduction," *Computational Optimization and Applications*, Vol. 1, pp. 7-66.
- [Ber96] Bertsekas, D. P., 1996. "Incremental Least Squares Methods and the Extended Kalman Filter," *SIAM J. on Optimization*, Vol. 6, pp. 807-822.
- [Ber97] Bertsekas, D. P., 1997. "A New Class of Incremental Gradient Methods for Least Squares Problems," *SIAM J. on Optimization*, Vol. 7, pp. 913-926.
- [Ber98] Bertsekas, D. P., 1998. *Network Optimization: Continuous and Discrete Models*, Athena Scientific, Belmont, MA.
- [Ber99] Bertsekas, D. P., 1999. *Nonlinear Programming: 2nd Edition*, Athena Scientific, Belmont, MA.

- [Ber07] Bertsekas, D. P., 2007. *Dynamic Programming and Optimal Control*, Vol. I, 3rd Edition, Athena Scientific, Belmont, MA.
- [Ber09] Bertsekas, D. P., 2009. *Convex Optimization Theory*, Athena Scientific, Belmont, MA.
- [Ber10a] Bertsekas, D. P., 2010. “Extended Monotropic Programming and Duality,” Lab. for Information and Decision Systems Report LIDS-P-2692, MIT, March 2006, Revised in Feb. 2010; a version appeared in *J. of Optimization Theory and Applications*, Vol. 139, pp. 209-225.
- [Ber10b] Bertsekas, D. P., 2010. “Incremental Gradient, Subgradient, and Proximal Methods for Convex Optimization: A Survey,” Lab. for Information and Decision Systems Report LIDS-P-2848, MIT.
- [Ber11] Bertsekas, D. P., 2011. “Incremental Proximal Methods for Large Scale Convex Optimization,” *Math. Programming*, Vol. 129, pp. 163-195.
- [Ber12] Bertsekas, D. P., 2012. *Dynamic Programming and Optimal Control*, Vol. II, 4th Edition: *Approximate Dynamic Programming*, Athena Scientific, Belmont, MA.
- [Ber13] Bertsekas, D. P., 2013. *Abstract Dynamic Programming*, Athena Scientific, Belmont, MA.
- [BiF07] Bioucas-Dias, J., and Figueiredo, M. A. T., 2007. “A New TwIST: Two-Step Iterative Shrinkage/Thresholding Algorithms for Image Restoration,” *IEEE Trans. Image Processing*, Vol. 16, pp. 2992-3004.
- [BiL97] Birge, J. R., and Louveaux, 1997. *Introduction to Stochastic Programming*, Springer-Verlag, New York, NY.
- [Bis95] Bishop, C. M., 1995. *Neural Networks for Pattern Recognition*, Oxford Univ. Press, NY.
- [BoL00] Borwein, J. M., and Lewis, A. S., 2000. *Convex Analysis and Nonlinear Optimization*, Springer-Verlag, NY.
- [BoL05] Bottou, L., and LeCun, Y., 2005. “On-Line Learning for Very Large Datasets,” *Applied Stochastic Models in Business and Industry*, Vol. 21, pp. 137-151.
- [BoS00] Bonnans, J. F., and Shapiro, A., 2000. *Perturbation Analysis of Optimization Problems*, Springer-Verlag, NY.
- [BoV04] Boyd, S., and Vandenbergue, L., 2004. *Convex Optimization*, Cambridge Univ. Press, Cambridge, UK.
- [Bor08] Borkar, V. S., 2008. *Stochastic Approximation: A Dynamical Systems Viewpoint*, Cambridge Univ. Press.
- [Bot05] Bottou, L., 2005. “SGD: Stochastic Gradient Descent,” <http://leon.bottou.org/projects/sgd>.
- [Bot09] Bottou, L., 2009. “Curiously Fast Convergence of Some Stochastic Gradient Descent Algorithms,” unpublished open problem offered to the attendance of the SLDS 2009 conference.

- [Bot10] Bottou, L., 2010. "Large-Scale Machine Learning with Stochastic Gradient Descent," In Proc. of COMPSTAT 2010, pp. 177-186.
- [BrL78] Brezis, H., and Lions, P. L., 1978. "Produits Infinites de Résolvantes," Israel J. of Mathematics, Vol. 29, pp. 329-345.
- [BrS13] Brown, D. B., and Smith, J. E., 2013. "Information Relaxations, Duality, and Convex Stochastic Dynamic Programs," Working Paper, Fuqua School of Business, Durham, NC, USA.
- [Bre73] Brezis, H., 1973. "Opérateurs Maximaux Monotones et Semi-Groupes de Contractions Dans les Espaces de Hilbert," North-Holland, Amsterdam.
- [BuM13] Burachik, R. S., and Majeed, S. N., 2013. "Strong Duality for Generalized Monotropic Programming in Infinite Dimensions," J. of Mathematical Analysis and Applications, Vol. 400, pp. 541-557.
- [BuQ98] Burke, J. V., and Qian, M., 1998. "A Variable Metric Proximal Point Algorithm for Monotone Operators," SIAM J. on Control and Optimization, Vol. 37, pp. 353-375.
- [Bur91] Burke, J. V., 1991. "An Exact Penalization Viewpoint of Constrained Optimization," SIAM J. on Control and Optimization, Vol. 29, pp. 968-998.
- [CDS01] Chen, S. S., Donoho, D. L., and Saunders, M. A., 2001. "Atomic Decomposition by Basis Pursuit," SIAM Review, Vol. 43, pp. 129-159.
- [CFM75] Camerini, P. M., Fratta, L., and Maffioli, F., 1975. "On Improving Relaxation Methods by Modified Gradient Techniques," Math. Programming Studies, Vol. 3, pp. 26-34.
- [CGT00] Conn, A. R., Gould, N. I., and Toint, P. L., 2000. Trust Region Methods, SIAM, Philadelphia, PA.
- [CHY13] Chen, C., He, B., Ye, Y., and Yuan, X., 2013. "The Direct Extension of ADMM for Multi-Block Convex Minimization Problems is not Necessarily Convergent," Optimization Online.
- [CPR14] Chouzenoux, E., Pesquet, J. C., and Repetti, A., 2014. "Variable Metric Forward-Backward Algorithm for Minimizing the Sum of a Differentiable Function and a Convex Function," J. of Optimization Theory and Applications, Vol. 162, pp. 107-132.
- [CPS92] Cottle, R. W., Pang, J.-S., and Stone, R. E., 1992. The Linear Complementarity Problem, Academic Press, NY.
- [CRP12] Chandrasekaran, V., Recht, B., Parrilo, P. A., and Willsky, A. S., 2012. "The Convex Geometry of Linear Inverse Problems," Foundations of Computational Mathematics, Vol. 12, pp. 805-849.
- [CaC68] Canon, M. D., and Cullum, C. D., 1968. "A Tight Upper Bound on the Rate of Convergence of the Frank-Wolfe Algorithm," SIAM J. on Control, Vol. 6, pp. 509-516.
- [CaG74] Cantor, D. G., Gerla, M., 1974. "Optimal Routing in Packet

- Switched Computer Networks,” *IEEE Trans. on Computing*, Vol. C-23, pp. 1062-1068.
- [CaR09] Candés, E. J., and Recht, B., 2009. “Exact Matrix Completion via Convex Optimization,” *Foundations of Computational Math.*, Vol. 9, pp. 717-772.
- [CaT10] Candés, E. J., and Tao, T., 2010. “The Power of Convex Relaxation: Near-Optimal Matrix Completion,” *IEEE Trans. on Information Theory*, Vol. 56, pp. 2053-2080.
- [CeH87] Censor, Y., and Herman, G. T., 1987. “On Some Optimization Techniques in Image Reconstruction from Projections,” *Applied Numer. Math.*, Vol. 3, pp. 365-391.
- [CeS08] Cegielski, A., and Suchocka, A., 2008. “Relaxed Alternating Projection Methods,” *SIAM J. Optimization*, Vol. 19, pp. 1093-1106.
- [CeZ92] Censor, Y., and Zenios, S. A., 1992. “The Proximal Minimization Algorithm with D-Functions,” *J. Opt. Theory and Appl.*, Vol. 73, pp. 451-464.
- [CeZ97] Censor, Y., and Zenios, S. A., 1997. *Parallel Optimization: Theory, Algorithms, and Applications*, Oxford Univ. Press, NY.
- [ChG59] Cheney, E. W., and Goldstein, A. A., 1959. “Newton’s Method for Convex Programming and Tchebycheff Approximation,” *Numer. Math.*, Vol. I, pp. 253-268.
- [ChR97] Chen, G. H., and Rockafellar, R. T., 1997. “Convergence Rates in Forward-Backward Splitting,” *SIAM J. on Optimization*, Vol. 7, pp. 421-444.
- [ChT93] Chen, G., and Teboulle, M., 1993. “Convergence Analysis of a Proximal-Like Minimization Algorithm Using Bregman Functions,” *SIAM J. on Optimization*, Vol. 3, pp. 538-543.
- [ChT94] Chen, G., and Teboulle, M., 1994. “A Proximal-Based Decomposition Method for Convex Minimization Problems,” *Mathematical Programming*, Vol. 64, pp. 81-101.
- [Cha04] Chambolle, A., 2004. “An Algorithm for Total Variation Minimization and Applications,” *J. of Mathematical Imaging and Vision*, Vol. 20, pp. 89-97.
- [Che07] Chen, Y., 2007. “A Smoothing Inexact Newton Method for Minimax Problems,” *Advances in Theoretical and Applied Mathematics*, Vol. 2, pp. 137-143.
- [Cla10] Clarkson, K. L., 2010. “Coresets, Sparse Greedy Approximation, and the Frank-Wolfe Algorithm,” *ACM Transactions on Algorithms*, Vol. 6, pp. 63.
- [CoT13] Couellan, N. P., and Trafalis, T. B., 2013. “On-line SVM Learning via an Incremental Primal-Dual Technique,” *Optimization Methods and Software*, Vol. 28, pp. 256-275.

- [CoV13] Combettes, P. L., and Vu, B. C., 2013. “Variable Metric Quasi-Fejér Monotonicity,” *Nonlinear Analysis: Theory, Methods and Applications*, Vol. 78, pp. 17-31.
- [Com01] Combettes, P. L., 2001. “Quasi-Fejérian Analysis of Some Optimization Algorithms,” *Studies in Computational Mathematics*, Vol. 8, pp. 115-152.
- [Cry71] Cryer, C. W., 1971. “The Solution of a Quadratic Programming Problem Using Systematic Overrelaxation,” *SIAM J. on Control*, Vol. 9, pp. 385-392.
- [DBW12] Duchi, J. C., Bartlett, P. L., and Wainwright, M. J., 2012. “Randomized Smoothing for Stochastic Optimization,” *SIAM J. on Optimization*, Vol. 22, pp. 674-701.
- [DCD14] Defazio, A. J., Caetano, T. S., and Domke, J., 2014. “Finito: A Faster, Permutable Incremental Gradient Method for Big Data Problems,” *Proceedings of the 31st ICML, Beijing*.
- [DHS06] Dai, Y. H., Hager, W. W., Schittkowski, K., and Zhang, H., 2006. “The Cyclic Barzilai-Borwein Method for Unconstrained Optimization,” *IMA J. of Numerical Analysis*, Vol. 26, pp. 604-627.
- [DHS11] Duchi, J., Hazan, E., and Singer, Y., 2011. “Adaptive Subgradient Methods for Online Learning and Stochastic Optimization,” *J. of Machine Learning Research*, Vol. 12, pp. 2121-2159.
- [DMM06] Drineas, P., Mahoney, M. W., and Muthukrishnan, S., 2006. “Sampling Algorithms for L2 Regression and Applications,” *Proc. 17th Annual SODA*, pp. 1127-1136.
- [DMM11] Drineas, P., Mahoney, M. W., Muthukrishnan, S., and Sarlos, T., 2011. “Faster Least Squares Approximation,” *Numerische Mathematik*, Vol. 117, pp. 219-249.
- [DRT11] Dhillon, I. S., Ravikumar, P., and Tewari, A., 2011. “Nearest Neighbor Based Greedy Coordinate Descent,” in *Advances in Neural Information Processing Systems 24*, (NIPS 2011), pp. 2160-2168.
- [DaW60] Dantzig, G. B., and Wolfe, P., 1960. “Decomposition Principle for Linear Programs,” *Operations Research*, Vol. 8, pp. 101-111.
- [DaY14a] Davis, D., and Yin, W., 2014. “Convergence Rate Analysis of Several Splitting Schemes,” *arXiv preprint arXiv:1406.4834*.
- [DaY14b] Davis, D., and Yin, W., 2014. “Convergence Rates of Relaxed Peaceman-Rachford and ADMM Under Regularity Assumptions,” *arXiv preprint arXiv:1407.5210*.
- [Dan67] Danskin, J. M., 1967. *The Theory of Max-Min and its Application to Weapons Allocation Problems*, Springer, NY.
- [Dav76] Davidon, W. C., 1976. “New Least Squares Algorithms,” *J. Optimization Theory and Applications*, Vol. 18, pp. 187-197.

- [DeM71] Demjanov, V. F., and Malozemov, V. N., 1971. "On the Theory of Non-Linear Minimax Problems," Russian Math. Surveys, Vol. 26, p. 57.
- [DeR70] Demjanov, V. F., and Rubinov, A. M., 1970. Approximate Methods in Optimization Problems, American Elsevier, NY.
- [DeS96] Dennis, J. E., and Schnabel, R. B., 1996. Numerical Methods for Unconstrained Optimization and Nonlinear Equations, SIAM, Philadelphia, PA.
- [DeT91] Dennis, J. E., and Torczon, V., 1991. "Direct Search Methods on Parallel Machines," SIAM J. on Optimization, Vol. 1, pp. 448-474.
- [Dem66] Demjanov, V. F., 1966. "The Solution of Several Minimax Problems," Kibernetika, Vol. 2, pp. 58-66.
- [Dem68] Demjanov, V. F., 1968. "Algorithms for Some Minimax Problems," J. of Computer and Systems Science, Vol. 2, pp. 342-380.
- [DoE03] Donoho, D. L., Elad, M., 2003. "Optimally Sparse Representation in General (Nonorthogonal) Dictionaries via ℓ_1 Minimization," Proc. of the National Academy of Sciences, Vol. 100, pp. 2197-2202.
- [DrH04] Drezner, Z., and Hamacher, H. W., 2004. Facility Location: Applications and Theory, Springer, NY.
- [DuS83] Dunn, J. C., and Sachs, E., 1983. "The Effect of Perturbations on the Convergence Rates of Optimization Algorithms," Appl. Math. Optim., Vol. 10, pp. 143-157.
- [DuS09] Duchi, J., and Singer, Y., 2009. "Efficient Online and Batch Learning Using Forward Backward Splitting, J. of Machine Learning Research, Vol. 10, pp. 2899-2934.
- [Dun79] Dunn, J. C., 1979. "Rates of Convergence for Conditional Gradient Algorithms Near Singular and Nonsingular Extremals," SIAM J. on Control and Optimization, Vol. 17, pp. 187-211.
- [Dun80] Dunn, J. C., 1980. "Convergence Rates for Conditional Gradient Sequences Generated by Implicit Step Length Rules," SIAM J. on Control and Optimization, Vol. 18, pp. 473-487.
- [Dun81] Dunn, J. C., 1981. "Global and Asymptotic Convergence Rate Estimates for a Class of Projected Gradient Processes," SIAM J. on Control and Optimization, Vol. 19, pp. 368-400.
- [Dun87] Dunn, J. C., 1987. "On the Convergence of Projected Gradient Processes to Singular Critical Points," J. of Optimization Theory and Applications, Vol. 55, pp. 203-216.
- [Dun91] Dunn, J. C., 1991. "A Subspace Decomposition Principle for Scaled Gradient Projection Methods: Global Theory," SIAM J. on Control and Optimization, Vol. 29, pp. 219-246.
- [EcB92] Eckstein, J., and Bertsekas, D. P., 1992. "On the Douglas-Rachford Splitting Method and the Proximal Point Algorithm for Maximal Monotone Operators," Math. Programming, Vol. 55, pp. 293-318.

- [EcS13] Eckstein, J., and Silva, P. J. S., 2013. “A Practical Relative Error Criterion for Augmented Lagrangians,” *Math. Programming*, Vol. 141, Ser. A, pp. 319-348.
- [Eck94] Eckstein, J., 1994. “Nonlinear Proximal Point Algorithms Using Bregman Functions, with Applications to Convex Programming,” *Math. of Operations Research*, Vol. 18, pp. 202-226.
- [Eck03] Eckstein, J., 2003. “A Practical General Approximation Criterion for Methods of Multipliers Based on Bregman Distances,” *Math. Programming*, Vol. 96, Ser. A, pp. 61-86.
- [Eck12] Eckstein, J., 2012. “Augmented Lagrangian and Alternating Direction Methods for Convex Optimization: A Tutorial and Some Illustrative Computational Results,” RUTCOR Research Report RRR 32-2012, Rutgers, Univ.
- [EkT76] Ekeland, I., and Temam, R., 1976. *Convex Analysis and Variational Problems*, North-Holland Publ., Amsterdam.
- [ElM75] Elzinga, J., and Moore, T. G., 1975. “A Central Cutting Plane Algorithm for the Convex Programming Problem,” *Math. Programming*, Vol. 8, pp. 134-145.
- [Erm76] Ermoliev, Yu. M., 1976. *Stochastic Programming Methods*, Nauka, Moscow.
- [Eve63] Everett, H., 1963. “Generalized Lagrange Multiplier Method for Solving Problems of Optimal Allocation of Resources,” *Operations Research*, Vol. 11, pp. 399-417.
- [FGW02] Forsgren, A., Gill, P. E., and Wright, M. H., 2002. “Interior Methods for Nonlinear Optimization,” *SIAM Review*, Vol. 44, pp. 525-597.
- [FHT10] Friedman, J., Hastie, T., and Tibshirani, R., 2010. “Regularization Paths for Generalized Linear Models via Coordinate Descent,” *J. of Statistical Software*, Vol. 33, pp. 1-22.
- [FJS98] Facchinei, F., Judice, J., and Soares, J., 1998. “An Active Set Newton Algorithm for Large-Scale Nonlinear Programs with Box Constraints,” *SIAM J. on Optimization*, Vol. 8, pp. 158-186.
- [FLP02] Facchinei, F., Lucidi, S., and Palagi, L., 2002. “A Truncated Newton Algorithm for Large Scale Box Constrained Optimization,” *SIAM J. on Optimization*, Vol. 12, pp. 1100-1125.
- [FLT02] Fukushima, M., Luo, Z.-Q., and Tseng, P., 2002. “Smoothing Functions for Second-Order-Cone Complementarity Problems,” *SIAM J. Optimization*, Vol. 12, pp. 436-460.
- [FaP03] Facchinei, F., and Pang, J.-S., 2003. *Finite-Dimensional Variational Inequalities and Complementarity Problems*, Springer Verlag, NY.
- [Fab73] Fabian, V., 1973. “Asymptotically Efficient Stochastic Approximation: The RM Case,” *Ann. Statist.*, Vol. 1, pp. 486-495.

- [FeM91] Ferris, M. C., and Mangasarian, O. L., 1991. "Finite Perturbation of Convex Programs," *Appl. Math. Optim.*, Vol. 23, pp. 263-273.
- [FeM02] Ferris, M. C., and Munson, T. S., 2002. "Interior-Point Methods for Massive Support Vector Machines," *SIAM J. on Optimization*, Vol. 13, pp. 783-804.
- [FeR12] Fercoq, O., and Richtarik, P., 2012. "Accelerated, Parallel, and Proximal Coordinate Descent," *arXiv preprint arXiv:1312.5799*.
- [Fei02] Feinberg, E. A., 2002. "Total Reward Criteria," in E. A. Feinberg and A. Shwartz, (Eds.), *Handbook of Markov Decision Processes*, Springer, N. Y.
- [Fen51] Fenchel, W., 1951. *Convex Cones, Sets, and Functions*, Mimeographed Notes, Princeton Univ.
- [FlH95] Florian, M. S., and Hearn, D., 1995. "Network Equilibrium Models and Algorithms," *Handbooks in OR and MS*, Ball, M. O., Magnanti, T. L., Monma, C. L., and Nemhauser, G. L., (Eds.), Vol. 8, North-Holland, Amsterdam, pp. 485-550.
- [FiM68] Fiacco, A. V., and McCormick, G. P., 1968. *Nonlinear Programming: Sequential Unconstrained Minimization Techniques*, Wiley, NY.
- [FiN03] Figueiredo, M. A. T., and Nowak, R. D., 2003. "An EM Algorithm for Wavelet-Based Image Restoration," *IEEE Trans. Image Processing*, Vol. 12, pp. 906-916.
- [Fle00] Fletcher, R., 2000. *Practical Methods of Optimization*, 2nd edition, Wiley, NY.
- [FoG83] Fortin, M., and Glowinski, R., 1983. "On Decomposition-Coordination Methods Using an Augmented Lagrangian," in: M. Fortin and R. Glowinski, eds., *Augmented Lagrangian Methods: Applications to the Solution of Boundary-Value Problems*, North-Holland, Amsterdam.
- [FrG13] Friedlander, M. P., and Goh, G., 2013. "Tail Bounds for Stochastic Approximation," *arXiv preprint arXiv:1304.5586*.
- [FrG14] Freund, R. M., and Grigas, P., 2014. "New Analysis and Results for the Frank-Wolfe Method," *arXiv preprint arXiv:1307.0873*, to appear in *Math. Programming*.
- [FrS00] Frommer, A., and Szyld, D. B., 2000. "On Asynchronous Iterations," *J. of Computational and Applied Mathematics*, Vol. 123, pp. 201-216.
- [FrS12] Friedlander, M. P., and Schmidt, M., 2012. "Hybrid Deterministic-Stochastic Methods for Data Fitting," *SIAM J. Sci. Comput.*, Vol. 34, pp. A1380-A1405.
- [FrT07] Friedlander, M. P., and Tseng, P., 2007. "Exact Regularization of Convex Programs," *SIAM J. on Optimization*, Vol. 18, pp. 1326-1350.
- [FrW56] Frank, M., and Wolfe, P., 1956. "An Algorithm for Quadratic Programming," *Naval Research Logistics Quarterly*, Vol. 3, pp. 95-110.

- [Fra02] Frangioni, A., 2002. "Generalized Bundle Methods," *SIAM J. on Optimization*, Vol. 13, pp. 117-156.
- [Fri56] Frisch, M. R., 1956. "La Resolution des Problemes de Programme Lineaire par la Methode du Potential Logarithmique," *Cahiers du Seminaire D'Econometrie*, Vol. 4, pp. 7-20.
- [FuM81] Fukushima, M., and Mine, H., 1981. "A Generalized Proximal Point Algorithm for Certain Non-Convex Minimization Problems," *Internat. J. Systems Sci.*, Vol. 12, pp. 989-1000.
- [Fuk92] Fukushima, M., 1992. "Application of the Alternating Direction Method of Multipliers to Separable Convex Programming Problems," *Computational Optimization and Applications*, Vol. 1, pp. 93-111.
- [GBY12] Grant, M., Boyd, S., and Ye, Y., 2012. "CVX: Matlab Software for Disciplined Convex Programming, Version 2.0 Beta," *Recent Advances in Learning and Control*, cvx.com.
- [GGM06] Gaudioso, M., Giallombardo, G., and Miglionico, G., 2006. "An Incremental Method for Solving Convex Finite Min-Max Problems," *Math. of Operations Research*, Vol. 31, pp. 173-187.
- [GHV92] Goffin, J. L., Haurie, A., and Vial, J. P., 1992. "Decomposition and Nondifferentiable Optimization with the Projective Algorithm," *Management Science*, Vol. 38, pp. 284-302.
- [GKX10] Gupta, M. D., Kumar, S., and Xiao, J. 2010. "L1 Projections with Box Constraints," *arXiv preprint arXiv:1010.0141*.
- [GLL86] Grippo, L., Lampariello, F., and Lucidi, S., 1986. "A Nonmonotone Line Search Technique for Newton's Method," *SIAM J. on Numerical Analysis*, Vol. 23, pp. 707-716.
- [GLY94] Goffin, J. L., Luo, Z.-Q., and Ye, Y., 1994. "On the Complexity of a Column Generation Algorithm for Convex or Quasiconvex Feasibility Problems," in *Large Scale Optimization: State of the Art*, Hager, W. W., Hearn, D. W., and Pardalos, P. M., (Eds.), Kluwer, Boston.
- [GLY96] Goffin, J. L., Luo, Z.-Q., and Ye, Y., 1996. "Complexity Analysis of an Interior Cutting Plane Method for Convex Feasibility Problems," *SIAM J. on Optimization*, Vol. 6, pp. 638-652.
- [GMW81] Gill, P. E., Murray, W., and Wright, M. H., 1981. *Practical Optimization*, Academic Press, NY.
- [GNS08] Griva, I., Nash, S. G., and Sofer, A., 2008. *Linear and Nonlinear Optimization*, 2nd Edition, SIAM, Philadelphia, PA.
- [GOP14] Gurbuzbalaban, M., Ozdaglar, A., and Parrilo, P., 2014. "A Globally Convergent Incremental Newton Method," *arXiv preprint arXiv:1410.5284*.
- [GPR67] Gubin, L. G., Polyak, B. T., and Raik, E. V., 1967. "The Method of Projection for Finding the Common Point in Convex Sets," *USSR Comput. Math. Phys.*, Vol. 7, pp. 1-24.

- [GSW12] Gamarnik, D., Shah, D., and Wei, Y., 2012. “Belief Propagation for Min-Cost Network Flow: Convergence and Correctness,” *Operations Research*, Vol. 60, pp. 410-428.
- [GaB84] Gafni, E. M., and Bertsekas, D. P., 1984. “Two-Metric Projection Methods for Constrained Optimization,” *SIAM J. on Control and Optimization*, Vol. 22, pp. 936-964.
- [GaM76] Gabay, D., and Mercier, B., 1976. “A Dual Algorithm for the Solution of Nonlinear Variational Problems via Finite-Element Approximations,” *Comp. Math. Appl.*, Vol. 2, pp. 17-40.
- [Gab79] Gabay, D., 1979. *Methodes Numeriques pour l’Optimization Non Lineaire*, These de Doctorat d’Etat et Sciences Mathematiques, Univ. Pierre et Marie Curie (Paris VI).
- [Gab83] Gabay, D., 1983. “Applications of the Method of Multipliers to Variational Inequalities,” in M. Fortin and R. Glowinski, eds., *Augmented Lagrangian Methods: Applications to the Solution of Boundary-Value Problems*, North-Holland, Amsterdam.
- [GeM05] Gerlach, S., and Matzenmiller, A., 2005. “Comparison of Numerical Methods for Identification of Viscoelastic Line Spectra from Static Test Data,” *International J. for Numerical Methods in Engineering*, Vol. 63, pp. 428-454.
- [Geo72] Geoffrion, A. M., 1972. “Generalized Benders Decomposition,” *J. of Optimization Theory and Applications*, Vol. 10, pp. 237-260.
- [Geo77] Geoffrion, A. M., 1977. “Objective Function Approximations in Mathematical Programming,” *Math. Programming*, Vol. 13, pp. 23-27.
- [GiB14] Giselsson, P., and Boyd, S., 2014. “Metric Selection in Douglas-Rachford Splitting and ADMM,” *arXiv preprint arXiv:1410.8479*.
- [GiM74] Gill, P. E., and Murray, W., (Eds.), 1974. *Numerical Methods for Constrained Optimization*, Academic Press, NY.
- [GLM75] Glowinski, R. and Marrocco, A., 1975. “Sur l’ Approximation par Elements Finis d’ Ordre un et la Resolution par Penalisation-Dualite d’une Classe de Problemes de Dirichlet Non Lineaires” *Revue Francaise d’Automatique Informatique Recherche Operationnelle, Analyse Numerique*, R-2, pp. 41-76.
- [GoK13] Gonzaga, C. C., and Karas, E. W., 2013. “Fine Tuning Nesterov’s Steepest Descent Algorithm for Differentiable Convex Programming,” *Math. Programming*, Vol. 138, pp. 141-166.
- [GoO09] Goldstein, T., and Osher, S., 2009. “The Split Bregman Method for L1-Regularized Problems,” *SIAM J. on Imaging Sciences*, Vol. 2, pp. 323-343.
- [GoS10] Goldstein, T., and Setzer, S., 2010. “High-Order Methods for Basis Pursuit,” *UCLA CAM Report*, 10-41.
- [GoV90] Goffin, J. L., and Vial, J. P., 1990. “Cutting Planes and Column

- Generation Techniques with the Projective Algorithm,” *J. Opt. Th. and Appl.*, Vol. 65, pp. 409-429.
- [GoV02] Goffin, J. L., and Vial, J. P., 2002. “Convex Nondifferentiable Optimization: A Survey Focussed on the Analytic Center Cutting Plane Method,” *Optimization Methods and Software*, Vol. 17, pp. 805-867.
- [GoZ12] Gong, P., and Zhang, C., 2012. “Efficient Nonnegative Matrix Factorization via Projected Newton Method,” *Pattern Recognition*, Vol. 45, pp. 3557-3565.
- [Gol64] Goldstein, A. A., 1964. “Convex Programming in Hilbert Space,” *Bull. Amer. Math. Soc.*, Vol. 70, pp. 709-710.
- [Gol85] Golshtein, E. G., 1985. “A Decomposition Method for Linear and Convex Programming Problems,” *Matecon*, Vol. 21, pp. 1077-1091.
- [Gon00] Gonzaga, C. C., 2000. “Two Facts on the Convergence of the Cauchy Algorithm,” *J. of Optimization Theory and Applications*, Vol. 107, pp. 591-600.
- [GrS99] Grippo, L., and Sciandrone, M., 1999. “Globally Convergent Block-Coordinate Techniques for Unconstrained Optimization,” *Optimization Methods and Software*, Vol. 10, pp. 587-637.
- [GrS00] Grippo, L., and Sciandrone, M., 2000. “On the Convergence of the Block Nonlinear Gauss-Seidel Method Under Convex Constraints,” *Operations Research Letters*, Vol. 26, pp. 127-136.
- [Gri94] Grippo, L., 1994. “A Class of Unconstrained Minimization Methods for Neural Network Training,” *Optim. Methods and Software*, Vol. 4, pp. 135-150.
- [Gri00] Grippo, L., 2000. “Convergent On-Line Algorithms for Supervised Learning in Neural Networks,” *IEEE Trans. Neural Networks*, Vol. 11, pp. 1284-1299.
- [Gul92] Guler, O., 1992. “New Proximal Point Algorithms for Convex Minimization,” *SIAM J. on Optimization*, Vol. 2, pp. 649-664.
- [HCW14] Hong, M., Chang, T. H., Wang, X., Razaviyayn, M., Ma, S., and Luo, Z. Q., 2014. “A Block Successive Upper Bound Minimization Method of Multipliers for Linearly Constrained Convex Optimization,” *arXiv preprint arXiv:1401.7079*.
- [HJN14] Harchaoui, Z., Juditsky, A., and Nemirovski, A., 2014. “Conditional Gradient Algorithms for Norm-Regularized Smooth Convex Optimization,” *Math. Programming*, pp. 1-38.
- [HKR95] den Hertog, D., Kaliski, J., Roos, C., and Terlaky, T., 1995. “A Path-Following Cutting Plane Method for Convex Programming,” *Annals of Operations Research*, Vol. 58, pp. 69-98.
- [HLV87] Hearn, D. W., Lawphongpanich, S., and Ventura, J. A., 1987. “Restricted Simplicial Decomposition: Computation and Extensions,” *Math. Programming Studies*, Vol. 31, pp. 119-136.

- [HMT10] Halko, N., Martinsson, P.-G., and Tropp, J. A., 2010. “Finding Structure with Randomness: Probabilistic Algorithms for Constructing Approximate Matrix Decompositions,” arXiv preprint arXiv:0909.4061.
- [HTF09] Hastie, T., Tibshirani, R., and Friedman, J., 2009. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, 2nd Edition, Springer, NY. On line at <http://statweb.stanford.edu/tibs/ElemStatLearn/>
- [HYS15] Hu, Y., Yang, X., and Sim, C. K., 2015. “Inexact Subgradient Methods for Quasi-Convex Optimization Problems,” *European Journal of Operational Research*, Vol. 240, pp. 315-327.
- [HZS13] Hou, K., Zhou, Z., So, A. M. C., and Luo, Z. Q., 2013. “On the Linear Convergence of the Proximal Gradient Method for Trace Norm Regularization,” in *Advances in Neural Information Processing Systems (NIPS 2013)*, pp. 710-718.
- [Ha90] Ha, C. D., 1990. “A Generalization of the Proximal Point Algorithm,” *SIAM J. on Control and Optimization*, Vol. 28, pp. 503-512.
- [HaB70] Haarhoff, P. C., and Buys, J. D., 1970. “A New Method for the Optimization of a Nonlinear Function Subject to Nonlinear Constraints,” *Computer J.*, Vol. 13, pp. 178-184.
- [HaH93] Hager, W. W., and Hearn, D. W., 1993. “Application of the Dual Active Set Algorithm to Quadratic Network Optimization,” *Computational Optimization and Applications*, Vol. 1, pp. 349-373.
- [HaM79] Han, S. P., and Mangasarian, O. L., 1979. “Exact Penalty Functions in Nonlinear Programming,” *Math. Programming*, Vol. 17, pp. 251-269.
- [Hay08] Haykin, S., 2008. *Neural Networks and Learning Machines*, (3rd Ed.), Prentice Hall, Englewood Cliffs, NJ.
- [HeD09] Helou, E. S., and De Pierro, A. R., 2009. “Incremental Subgradients for Constrained Convex Optimization: A Unified Framework and New Methods,” *SIAM J. on Optimization*, Vol. 20, pp. 1547-1572.
- [HeL89] Hearn, D. W., and Lawphongpanich, S., 1989. “Lagrangian Dual Ascent by Generalized Linear Programming,” *Operations Res. Letters*, Vol. 8, pp. 189-196.
- [HeM11] Henrion, D., and Malick, J., 2011. “Projection Methods for Conic Feasibility Problems: Applications to Polynomial Sum-of-Squares Decompositions,” *Optimization Methods and Software*, Vol. 26, pp. 23-46.
- [HeM12] Henrion, D., and Malick, J., 2012. “Projection Methods in Conic Optimization,” In *Handbook on Semidefinite, Conic and Polynomial Optimization*, Springer, NY, pp. 565-600.
- [Her09] Gabor, H., 2009. *Fundamentals of Computerized Tomography: Image Reconstruction from Projection*, (2nd ed.), Springer, NY.

- [Hes69] Hestenes, M. R., 1969. "Multiplier and Gradient Methods," *J. Opt. Th. and Appl.*, Vol. 4, pp. 303-320.
- [Hes75] Hestenes, M. R., 1975. *Optimization Theory: The Finite Dimensional Case*, Wiley, NY.
- [Hil93] Hiriart-Urruty, J.-B., and Lemarechal, C., 1993. *Convex Analysis and Minimization Algorithms*, Vols. I and II, Springer-Verlag, Berlin and NY.
- [Hil57] Hildreth, C., 1957. "A Quadratic Programming Procedure," *Naval Res. Logist. Quart.*, Vol. 4, pp. 79-85. See also "Erratum," *Naval Res. Logist. Quart.*, Vol. 4, p. 361.
- [HoK71] Hoffman, K., and Kunze, R., 1971. *Linear Algebra*, Pearson, Englewood Cliffs, NJ.
- [HoL13] Hong, M., and Luo, Z. Q., 2013. "On the Linear Convergence of the Alternating Direction Method of Multipliers," *arXiv preprint arXiv:1208.3922*.
- [Hoh77] Hohenbalken, B. von, 1977. "Simplicial Decomposition in Nonlinear Programming," *Math. Programming*, Vol. 13, pp. 49-68.
- [Hol74] Holloway, C. A., 1974. "An Extension of the Frank and Wolfe Method of Feasible Directions," *Math. Programming*, Vol. 6, pp. 14-27.
- [IPS03] Iusem, A. N., Pennanen, T., and Svaiter, B. F., 2003. "Inexact Variants of the Proximal Point Algorithm Without Monotonicity," *SIAM J. on Optimization*, Vol. 13, pp. 1080-1097.
- [IST94] Iusem, A. N., Svaiter, B. F., and Teboulle, M., 1994. "Entropy-Like Proximal Methods in Convex Programming," *Math. of Operations Research*, Vol. 19, pp. 790-814.
- [IbF96] Ibaraki, S., and Fukushima, M., 1996. "Partial Proximal Method of Multipliers for Convex Programming Problems," *J. of Operations Research Society of Japan*, Vol. 39, pp. 213-229.
- [IuT95] Iusem, A. N., and Teboulle, M., 1995. "Convergence Rate Analysis of Nonquadratic Proximal Methods for Convex and Linear Programming," *Math. of Operations Research*, Vol. 20, pp. 657-677.
- [Ius99] Iusem, A. N., 1999. "Augmented Lagrangian Methods and Proximal Point Methods for Convex Minimization," *Investigacion Operativa*, Vol. 8, pp. 11-49.
- [Ius03] Iusem, A. N., 2003. "On the Convergence Properties of the Projected Gradient Method for Convex Optimization," *Computational and Applied Mathematics*, Vol. 22, pp. 37-52.
- [JFY09] Joachims, T., Finley, T., and Yu, C.-N. J., 2009. "Cutting-Plane Training of Structural SVMs," *Machine Learning*, Vol. 77, pp. 27-59.
- [JRJ09] Johansson, B., Rabi, M., and Johansson, M., 2009. "A Randomized Incremental Subgradient Method for Distributed Optimization in Networked Systems," *SIAM J. on Optimization*, Vol. 20, pp. 1157-1170.

- [Jag13] Jaggi, M., 2013. "Revisiting Frank-Wolfe: Projection-Free Sparse Convex Optimization," Proc. of ICML 2013.
- [JiZ14] Jiang, B., and Zhang, S., 2014. "Iteration Bounds for Finding the ϵ -Stationary Points for Structured Nonconvex Optimization," arXiv preprint arXiv:1410.4066.
- [JoY09] Joachims, T., and Yu, C.-N. J., 2009. "Sparse Kernel SVMs via Cutting-Plane Training," Machine Learning, Vol. 76, pp. 179-193.
- [JoZ13] Johnson, R., and Zhang, T., 2013. "Accelerating Stochastic Gradient Descent Using Predictive Variance Reduction," Advances in Neural Information Processing Systems 26 (NIPS 2013).
- [Joa06] Joachims, T., 2006. "Training Linear SVMs in Linear Time," International Conference on Knowledge Discovery and Data Mining, pp. 217-226.
- [JuN11a] Juditsky, A., and Nemirovski, A., 2011. "First Order Methods for Nonsmooth Convex Large-Scale Optimization, I: General Purpose Methods," in Optimization for Machine Learning, by Sra, S., Nowozin, S., and Wright, S. J. (eds.), MIT Press, Cambridge, MA, pp. 121-148.
- [JuN11b] Juditsky, A., and Nemirovski, A., 2011. "First Order Methods for Nonsmooth Convex Large-Scale Optimization, II: Utilizing Problem's Structure," in Optimization for Machine Learning, by Sra, S., Nowozin, S., and Wright, S. J. (eds.), MIT Press, Cambridge, MA, pp. 149-183.
- [KaW94] Kall, P., and Wallace, S. W., 1994. Stochastic Programming, Wiley, Chichester, UK.
- [Kac37] Kaczmarz, S., 1937. "Approximate Solution of Systems of Linear Equations," Bull. Acad. Pol. Sci., Lett. A 35, pp. 335-357 (in German); English transl.: Int. J. Control, Vol. 57, pp. 1269-1271, 1993.
- [Kar84] Karmarkar, N., 1984. "A New Polynomial-Time Algorithm for Linear Programming," In Proc. of the 16th Annual ACM Symp. on Theory of Computing, pp. 302-311.
- [Kel60] Kelley, J. E., 1960. "The Cutting-Plane Method for Solving Convex Programs," J. Soc. Indust. Appl. Math., Vol. 8, pp. 703-712.
- [Kel99] Kelley, C. T., 1999. Iterative Methods for Optimization, Siam, Philadelphia, PA.
- [Kib80] Kibardin, V. M., 1980. "Decomposition into Functions in the Minimization Problem," Automation and Remote Control, Vol. 40, pp. 1311-1323.
- [Kiw04] Kiwiel, K. C., 2004. "Convergence of Approximate and Incremental Subgradient Methods for Convex Optimization," SIAM J. on Optimization, Vol. 14, pp. 807-840.
- [KoB72] Kort, B. W., and Bertsekas, D. P., 1972. "A New Penalty Function Method for Constrained Minimization," Proc. 1972 IEEE Confer. Decision Control, New Orleans, LA, pp. 162-166.

- [KoB76] Kort, B. W., and Bertsekas, D. P., 1976. "Combined Primal-Dual and Penalty Methods for Convex Programming," *SIAM J. on Control and Optimization*, Vol. 14, pp. 268-294.
- [KoN93] Kortanek, K. O., and No, H., 1993. "A Central Cutting Plane Algorithm for Convex Semi-Infinite Programming Problems," *SIAM J. on Optimization*, Vol. 3, pp. 901-918.
- [Kor75] Kort, B. W., 1975. "Combined Primal-Dual and Penalty Function Algorithms for Nonlinear Programming," Ph.D. Thesis, Dept. of Engineering-Economic Systems, Stanford Univ., Stanford, Ca.
- [Kra55] Krasnosel'skii, M. A., 1955. "Two Remarks on the Method of Successive Approximations," *Uspehi Mat. Nauk*, Vol. 10, pp. 123-127.
- [KuC78] Kushner, H. J., and Clark, D. S., 1978. *Stochastic Approximation Methods for Constrained and Unconstrained Systems*, Springer-Verlag, NY.
- [KuY03] Kushner, H. J., and Yin, G., 2003. *Stochastic Approximation and Recursive Algorithms and Applications*, Springer-Verlag, NY.
- [LBB98] LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P., 1998. "Gradient-Based Learning Applied to Document Recognition," *Proceedings of the IEEE*, Vol. 86, pp. 2278-2324.
- [LJS12] Lacoste-Julien, S., Jaggi, M., Schmidt, M., and Pletscher, P., 2012. "Block-Coordinate Frank-Wolfe Optimization for Structural SVMs," *arXiv preprint arXiv:1207.4747*.
- [LLS12] Lee, J., Sun, Y., and Saunders, M., 2012. "Proximal Newton-Type Methods for Convex Optimization," *NIPS 2012*.
- [LLS14] Lee, J., Sun, Y., and Saunders, M., 2014. "Proximal Newton-Type Methods for Minimizing Composite Functions," *arXiv preprint arXiv:1206.1623*.
- [LLX14] Lin, Q., Lu, Z., and Xiao, L., 2014. "An Accelerated Proximal Coordinate Gradient Method and its Application to Regularized Empirical Risk Minimization," *arXiv preprint arXiv:1407.1296*.
- [LLZ09] Langford, J., Li, L., and Zhang, T., 2009. "Sparse Online Learning via Truncated Gradient," In *Advances in Neural Information Processing Systems (NIPS 2009)*, pp. 905-912.
- [LMS92] Lustig, I. J., Marsten, R. E., and Shanno, D. F., 1992. "On Implementing Mehrotra's Predictor-Corrector Interior-Point Method for Linear Programming," *SIAM J. on Optimization*, Vol. 2, pp. 435-449.
- [LMY12] Lu, Z., Monteiro, R. D. C., and Yuan, M., 2012. "Convex Optimization Methods for Dimension Reduction and Coefficient Estimation in Multivariate Linear Regression," *Mathematical Programming*, Vol. 131, pp. 163-194.
- [LPS98] Larsson, T., Patriksson, M., and Stromberg, A.-B., 1998. "Ergodic Convergence in Subgradient Optimization," *Optimization Methods*

and Software, Vol. 9, pp. 93-120.

[LRW98] Lagarias, J. C., Reeds, J. A., Wright, M. H., and Wright, P. E., 1998. "Convergence Properties of the Nelder-Mead Simplex Method in Low Dimensions," *SIAM J. on Optimization*, Vol. 9, pp. 112-147.

[LVB98] Lobo, M. S., Vandenberghe, L., Boyd, S., and Lebrete, H., 1998. "Applications of Second-Order Cone Programming," *Linear Algebra and Applications*, Vol. 284, pp. 193-228.

[LWS14] Liu, J., Wright, S. J., and Sridhar, S., 2014. "An Asynchronous Parallel Randomized Kaczmarz Algorithm," Univ. of Wisconsin Report, arXiv preprint arXiv:1401.4780.

[LaD60] Land, A. H., and Doig, A. G., 1960. "An Automatic Method for Solving Discrete Programming Problems," *Econometrica*, Vol. 28, pp. 497-520.

[LaS87] Lawrence, J., and Spingarn, J. E., 1987. "On Fixed Points of Non-expansive Piecewise Isometric Mappings," *Proc. London Math. Soc.*, Vol. 55, pp. 605-624.

[LaT85] Lancaster, P., and Tismenetsky, M., 1985. *The Theory of Matrices*, Academic Press, NY.

[Lan15] Landi, G., 2016. "A Modified Newton Projection Method for ℓ_1 -Regularized Least Squares Image Deblurring," *J. of Mathematical Imaging and Vision*, Vol. 51, pp. 195-208.

[LeL10] Leventhal, D., and Lewis, A. S., 2010. "Randomized Methods for Linear Constraints: Convergence Rates and Conditioning," *Math. of Operations Research*, Vol. 35, pp. 641-654.

[LeP65] Levitin, E. S., and Poljak, B. T., 1965. "Constrained Minimization Methods," *Ž. Vyčisl. Mat. i Mat. Fiz.*, Vol. 6, pp. 787-823.

[LeS93] Lemaréchal, C., and Sagastizábal, C., 1993. "An Approach to Variable Metric Bundle Methods," in *Systems Modelling and Optimization*, Proc. of the 16th IFIP-TC7 Conference, Compiègne, Henry, J., and Yvon, J.-P., (Eds.), *Lecture Notes in Control and Information Sciences* 197, pp. 144-162.

[LeS99] Lee, D., and Seung, H., 1999. "Learning the Parts of Objects by Non-Negative Matrix Factorization," *Nature*, Vol. 401, pp. 788-791.

[LeS13] Lee, Y. T., and Sidford, A., 2013. "Efficient Accelerated Coordinate Descent Methods and Faster Algorithms for Solving Linear Systems," *Proc. 2013 IEEE 54th Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 147-156.

[LeW11] Lee, S., and Wright, S. J., 2011. "Approximate Stochastic Sub-gradient Estimation Training for Support Vector Machines," Univ. of Wisconsin Report, arXiv preprint arXiv:1111.0432.

[Lem74] Lemarechal, C., 1974. "An Algorithm for Minimizing Convex Functions," in *Information Processing '74*, Rosenfeld, J. L., (Ed.), North-Holland,

Amsterdam, pp. 552-556.

[Lem75] Lemarechal, C., 1975. "An Extension of Davidon Methods to Non-differentiable Problems," Math. Programming Study 3, Balinski, M., and Wolfe, P., (Eds.), North-Holland, Amsterdam, pp. 95-109.

[Lem89] Lemaire, B., 1989. "The Proximal Algorithm," in New Methods in Optimization and Their Industrial Uses, J.-P. Penot, (ed.), Birkhauser, Basel, pp. 73-87.

[LiM79] Lions, P. L., and Mercier, B., 1979. "Splitting Algorithms for the Sum of Two Nonlinear Operators," SIAM J. on Numerical Analysis, Vol. 16, pp. 964-979.

[LiP87] Lin, Y. Y., and Pang, J.-S., 1987. "Iterative Methods for Large Convex Quadratic Programs: A Survey," SIAM J. on Control and Optimization, Vol. 18, pp. 383-411.

[LiW14] Liu, J., and Wright, S. J., 2014. "Asynchronous Stochastic Coordinate Descent: Parallelism and Convergence Properties," Univ. of Wisconsin Report, arXiv preprint arXiv:1403.3862.

[Lin07] Lin, C. J., 2007. "Projected Gradient Methods for Nonnegative Matrix Factorization," Neural Computation, Vol. 19, pp. 2756-2779.

[Lit66] Litvakov, B. M., 1966. "On an Iteration Method in the Problem of Approximating a Function from a Finite Number of Observations," Avtom. Telemekh., No. 4, pp. 104-113.

[Lju77] Ljung, L., 1977. "Analysis of Recursive Stochastic Algorithms," IEEE Trans. on Automatic Control, Vol. 22, pp. 551-575.

[LuT91] Luo, Z. Q., and Tseng, P., 1991. "On the Convergence of a Matrix-Splitting Algorithm for the Symmetric Monotone Linear Complementarity Problem," SIAM J. on Control and Optimization, Vol. 29, pp. 1037-1060.

[LuT92] Luo, Z. Q., and Tseng, P., 1992. "On the Convergence of the Coordinate Descent Method for Convex Differentiable Minimization," J. Optim. Theory Appl., Vol. 72, pp. 7-35.

[LuT93a] Luo, Z. Q., and Tseng, P., 1993. "On the Convergence Rate of Dual Ascent Methods for Linearly Constrained Convex Minimization," Math. of Operations Research, Vol. 18, pp. 846-867.

[LuT93b] Luo, Z. Q., and Tseng, P., 1993. "Error Bound and Reduced-Gradient Projection Algorithms for Convex Minimization over a Polyhedral Set," SIAM J. on Optimization, Vol. 3, pp. 43-59.

[LuT93c] Luo, Z. Q., and Tseng, P., 1993. "Error Bounds and Convergence Analysis of Feasible Descent Methods: A General Approach," Annals of Operations Research, Vol. 46, pp. 157-178.

[LuT94a] Luo, Z. Q., and Tseng, P., 1994. "Analysis of an Approximate Gradient Projection Method with Applications to the Backpropagation Algorithm," Optimization Methods and Software, Vol. 4, pp. 85-101.

[LuT94b] Luo, Z. Q., and Tseng, P., 1994. "On the Rate of Convergence of a

- Distributed Asynchronous Routing Algorithm,” *IEEE Trans. on Automatic Control*, Vol. 39, pp. 1123-1129.
- [LuT13] Luss, R., and Teboulle, M., 2013. “Conditional Gradient Algorithms for Rank-One Matrix Approximations with a Sparsity Constraint,” *SIAM Review*, Vol. 55, pp. 65-98.
- [LuY08] Luenberger, D. G., and Ye, Y., 2008. *Linear and Nonlinear Programming*, 3rd Edition, Springer, NY.
- [Lue84] Luenberger, D. G., 1984. *Introduction to Linear and Nonlinear Programming*, 2nd Edition, Addison-Wesley, Reading, MA.
- [Luo91] Luo, Z. Q., 1991. “On the Convergence of the LMS Algorithm with Adaptive Learning Rate for Linear Feedforward Networks,” *Neural Computation*, Vol. 3, pp. 226-245.
- [Luq84] Luque, F. J., 1984. “Asymptotic Convergence Analysis of the Proximal Point Algorithm,” *SIAM J. on Control and Optimization*, Vol. 22, pp. 277-293.
- [MRS10] Mosik-Aoyama, D., Roughgarden, T., and Shah, D., 2010. “Fully Distributed Algorithms for Convex Optimization Problems,” *SIAM J. on Optimization*, Vol. 20, pp. 3260-3279.
- [MSQ98] Mifflin, R., Sun, D., and Qi, L., 1998. “Quasi-Newton Bundle-Type Methods for Nondifferentiable Convex Optimization,” *SIAM J. on Optimization*, Vol. 8, pp. 583-603.
- [MYF03] Moriyama, H., Yamashita, N., and Fukushima, M., 2003. “The Incremental Gauss-Newton Algorithm with Adaptive Stepsize Rule,” *Computational Optimization and Applications*, Vol. 26, pp. 107-141.
- [MaM01] Mangasarian, O. L., Musicant, D. R., 2001. “Lagrangian Support Vector Machines,” *J. of Machine Learning Research*, Vol. 1, pp. 161-177.
- [MaS94] Mangasarian, O. L., and Solodov, M. V., 1994. “Serial and Parallel Backpropagation Convergence Via Nonmonotone Perturbed Minimization,” *Opt. Methods and Software*, Vol. 4, pp. 103-116.
- [Mai13] Mairal, J., 2013. “Optimization with First-Order Surrogate Functions,” *arXiv preprint arXiv:1305.3120*.
- [Mai14] Mairal, J., 2014. “Incremental Majorization-Minimization Optimization with Application to Large-Scale Machine Learning,” *arXiv preprint arXiv:1402.4419*.
- [Man53] Mann, W. R., 1953. “Mean Value Methods in Iteration,” *Proc. Amer. Math. Soc.*, Vol. 4, pp. 506-510.
- [Mar70] Martinet, B., 1970. “Regularisation d’Inéquations Variationnelles par Approximations Successives,” *Revue Fran. d’Automatique et Infomatique Rech. Opérationnelle*, Vol. 4, pp. 154-159.
- [Mar72] Martinet, B., 1972. “Détermination Approchée d’un Point Fixe d’une Application Pseudo-Contractante. Cas de l’Application Prox,” *Comptes Rendus de l’Académie des Sciences, Paris, Serie A* 274, pp. 163-165.

- [Meh92] Mehrotra, S., 1992. "On the Implementation of a Primal-Dual Interior Point Method," *SIAM J. on Optimization*, Vol. 2, pp. 575-601.
- [Mey07] Meyn, S., 2007. *Control Techniques for Complex Networks*, Cambridge Univ. Press, NY.
- [MiF81] Mine, H., Fukushima, M. 1981. "A Minimization Method for the Sum of a Convex Function and a Continuously Differentiable Function," *J. of Optimization Theory and Applications*, Vol. 33, pp. 9-23.
- [Mif96] Mifflin, R., 1996. "A Quasi-Second-Order Proximal Bundle Algorithm," *Math. Programming*, Vol. 73, pp. 51-72.
- [Min62] Minty, G. J., 1962. "Monotone (Nonlinear) Operators in Hilbert Space," *Duke J. of Math.*, Vol. 29, pp. 341-346.
- [Min64] Minty, G. J., 1964. "On the Monotonicity of the Gradient of a Convex Function," *Pacific J. of Math.*, Vol. 14, pp. 243-247.
- [Min86] Minoux, M., 1986. *Math. Programming: Theory and Algorithms*, Wiley, NY.
- [MoT89] Moré, J. J., and Toraldo, G., 1989. "Algorithms for Bound Constrained Quadratic Programming Problems," *Numer. Math.*, Vol. 55, pp. 377-400.
- [NBB01] Nedić, A., Bertsekas, D. P., and Borkar, V., 2001. "Distributed Asynchronous Incremental Subgradient Methods," *Proc. of 2000 Haifa Workshop "Inherently Parallel Algorithms in Feasibility and Optimization and Their Applications,"* by D. Butnariu, Y. Censor, and S. Reich, Eds., Elsevier, Amsterdam.
- [NJL09] Nemirovski, A., Juditsky, A., Lan, G., and Shapiro, A., 2009. "Robust Stochastic Approximation Approach to Stochastic Programming," *SIAM J. on Optimization*, Vol. 19, pp. 1574-1609.
- [NSW14] Needell, D., Srebro, N., and Ward, R., 2014. "Stochastic Gradient Descent and the Randomized Kaczmarz Algorithm," *arXiv preprint arXiv:1310.5715v3*.
- [NaT02] Nazareth, L., and Tseng, P., 2002. "Gilding the Lily: A Variant of the Nelder-Mead Algorithm Based on Golden-Section Search," *Computational Optimization and Applications*, Vol. 22, pp. 133-144.
- [NaZ05] Narkiss, G., and Zibulevsky, M., 2005. "Sequential Subspace Optimization Method for Large-Scale Unconstrained Problems," *Technion-IIT, Department of Electrical Engineering*.
- [NeB00] Nedić, A., and Bertsekas, D. P., 2000. "Convergence Rate of Incremental Subgradient Algorithms," *Stochastic Optimization: Algorithms and Applications*, S. Uryasev and P. M. Pardalos, Eds., Kluwer, pp. 263-304.
- [NeB01] Nedić, A., and Bertsekas, D. P., 2001. "Incremental Subgradient Methods for Nondifferentiable Optimization," *SIAM J. on Optimization*, Vol. 12, pp. 109-138.

- [NeB10] Nedić, A., and Bertsekas, D. P., 2010. "The Effect of Deterministic Noise in Subgradient Methods," *Math. Programming, Ser. A*, Vol. 125, pp. 75-99.
- [NeC13] Necoara, I., and Clipici, D., 2013. "Distributed Coordinate Descent Methods for Composite Minimization," *arXiv preprint arXiv:1312.5302*.
- [NeM65] Nelder, J. A., and Mead, R., 1965. "A Simplex Method for Function Minimization," *Computer J.*, Vol. 7, pp. 308-313.
- [NeN94] Nesterov, Y., and Nemirovskii, A., 1994. *Interior Point Polynomial Algorithms in Convex Programming*, SIAM, Studies in Applied Mathematics 13, Philadelphia, PA.
- [NeO09a] Nedić, A., and Ozdaglar, A., 2009. "Distributed Subgradient Methods for Multi-Agent Optimization," *IEEE Trans. on Aut. Control*, Vol. 54, pp. 48-61.
- [NeO09b] Nedić, A., and Ozdaglar, A., 2009. "Subgradient Methods for Saddle-Point Problems," *J. of Optimization Theory and Applications*, Vol. 142, pp. 205-228.
- [NeW88] Nemhauser, G. L., and Wolsey, L. A., 1988. *Integer and Combinatorial Optimization*, Wiley, NY.
- [NeY83] Nemirovsky, A., and Yudin, D. B., 1983. *Problem Complexity and Method Efficiency*, Wiley, NY.
- [Ned10] Nedić, A., 2010. "Random Projection Algorithms for Convex Set Intersection Problems," *Proc. 2010 IEEE Conference on Decision and Control*, Atlanta, Georgia, pp. 7655-7660.
- [Ned11] Nedić, A., 2011. "Random Algorithms for Convex Minimization Problems," *Math. Programming, Ser. B*, Vol. 129, pp. 225-253.
- [Nee10] Needell, D., 2010. "Randomized Kaczmarz Solver for Noisy Linear Systems," *BIT Numerical Mathematics*, Vol. 50, pp. 395-403.
- [Nes83] Nesterov, Y., 1983. "A Method for Unconstrained Convex Minimization Problem with the Rate of Convergence $O(1/k^2)$," *Doklady AN SSSR*, Vol. 269, pp. 543-547; translated as *Soviet Math. Dokl.*
- [Nes95] Nesterov, Y., 1995. "Complexity Estimates of Some Cutting Plane Methods Based on Analytic Barrier," *Math. Programming*, Vol. 69, pp. 149-176.
- [Nes04] Nesterov, Y., 2004. *Introductory Lectures on Convex Optimization*, Kluwer Academic Publisher, Dordrecht, The Netherlands.
- [Nes05] Nesterov, Y., 2005. "Smooth Minimization of Nonsmooth Functions," *Math. Programming*, Vol. 103, pp. 127-152.
- [Nes12] Nesterov, Y., 2012. "Efficiency of Coordinate Descent Methods on Huge-Scale Optimization Problems," *SIAM J. on Optimization*, Vol. 22, pp. 341-362.
- [Nev75] Neveu, J., 1975. *Discrete Parameter Martingales*, North-Holland, Amsterdam, The Netherlands.

- [NoW06] Nocedal, J., and Wright, S. J., 2006. Numerical Optimization, 2nd Edition, Springer, NY.
- [Noc80] Nocedal, J., 1980. "Updating Quasi-Newton Matrices with Limited Storage," Math. of Computation, Vol. 35, pp. 773-782.
- [OBG05] Osher, S., Burger, M., Goldfarb, D., Xu, J., and Yin, W., 2005. "An Iterative Regularization Method for Total Variation-Based Image Restoration," Multiscale Modeling and Simulation, Vol. 4, pp. 460-489.
- [OJW05] Olafsson, A., Jeraj, R., and Wright, S. J., 2005. Optimization of Intensity-Modulated Radiation Therapy with Biological Objectives," Physics in Medicine and Biology, Vol. 50, pp. 53-57.
- [OMV00] Ounorou, A., Mahey, P., and Vial, J. P., 2000. "A Survey of Algorithms for Convex Multicommodity Flow Problems," Management Science, Vol. 46, pp. 126-147.
- [OrR70] Ortega, J. M., and Rheinboldt, W. C., 1970. Iterative Solution of Nonlinear Equations in Several Variables, Academic Press, NY.
- [OvG14] Ovcharova, N., and Gwinner, J., 2014. "A Study of Regularization Techniques of Nondifferentiable Optimization in View of Application to Hemivariational Inequalities," J. of Optimization Theory and Applications, Vol. 162, pp. 754-778.
- [OzB03] Ozdaglar, A. E., and Bertsekas, D. P., 2003. "Routing and Wavelength Assignment in Optical Networks," IEEE Trans. on Networking, Vol. 11, pp. 259-272.
- [PKP09] Predd, J. B., Kulkarni, S. R., and Poor, H. V., 2009. "A Collaborative Training Algorithm for Distributed Learning," IEEE Transactions on Information Theory, Vol. 55, pp. 1856-1871.
- [PaE10] Palomar, D. P., and Eldar, Y. C., (Eds.), 2010. Convex Optimization in Signal Processing and Communications, Cambridge Univ. Press, NY.
- [PaT94] Paatero, P., and Tapper, U., 1994. "Positive Matrix Factorization: A Non-Negative Factor Model with Optimal Utilization of Error Estimates of Data Values," Environmetrics, Vol. 5, pp. 111-126.
- [PaY84] Pang, J. S., Yu, C. S., 1984. "Linearized Simplicial Decomposition Methods for Computing Traffic Equilibria on Networks," Networks, Vol. 14, pp. 427-438.
- [Paa97] Paatero, P., 1997. "Least Squares Formulation of Robust Non-Negative Factor Analysis," Chemometrics and Intell. Laboratory Syst., Vol. 37, pp. 23-35.
- [Pan84] Pang, J. S., 1984. "On the Convergence of Dual Ascent Methods for Large-Scale Linearly Constrained Optimization Problems," Unpublished manuscript, The Univ. of Texas at Dallas.
- [Pap81] Papavassilopoulos, G., 1981. "Algorithms for a Class of Nondifferentiable Problems," J. of Optimization Theory and Applications, Vol. 34,

pp. 41-82.

[Pas79] Passty, G. B., 1979. "Ergodic Convergence to a Zero of the Sum of Monotone Operators in Hilbert Space," *J. Math. Anal. Appl.*, Vol. 72, pp. 383-390.

[Pat93] Patriksson, M., 1993. "Partial Linearization Methods in Nonlinear Programming," *J. of Optimization Theory and Applications*, Vol. 78, pp. 227-246.

[Pat98] Patriksson, M., 1998. "Cost Approximation: A Unified Framework of Descent Algorithms for Nonlinear Programs," *SIAM J. Optimization*, Vol. 8, pp. 561-582.

[Pat99] Patriksson, M., 1999. *Nonlinear Programming and Variational Inequality Problems: A Unified Approach*, Springer, NY.

[Pat01] Patriksson, M., 2001. "Simplicial Decomposition Algorithms," *Encyclopedia of Optimization*, Springer, pp. 2378-2386.

[Pat04] Patriksson, M., 2004. "Algorithms for Computing Traffic Equilibria," *Networks and Spatial Economics*, Vol. 4, pp. 23-38.

[Pen02] Pennanen, T., 2002. "Local Convergence of the Proximal Point Algorithm and Multiplier Methods Without Monotonicity," *Math. of Operations Research*, Vol. 27, pp. 170-191.

[Pfl96] Pflug, G., 1996. *Optimization of Stochastic Models. The Interface Between Simulation and Optimization*, Kluwer, Boston.

[PiZ94] Pinar, M., and Zenios, S., 1994. "On Smoothing Exact Penalty Functions for Convex Constrained Optimization," *SIAM J. on Optimization*, Vol. 4, pp. 486-511.

[PoJ92] Poljak, B. T., and Juditsky, A. B., 1992. "Acceleration of Stochastic Approximation by Averaging," *SIAM J. on Control and Optimization*, Vol. 30, pp. 838-855.

[PoT73] Poljak, B. T., and Tsypkin, Y. Z., 1973. "Pseudogradient Adaptation and Training Algorithms," *Automation and Remote Control*, Vol. 12, pp. 83-94.

[PoT74] Poljak, B. T., and Tretjakov, N. V., 1974. "An Iterative Method for Linear Programming and its Economic Interpretation," *Matecon*, Vol. 10, pp. 81-100.

[PoT80] Poljak, B. T., and Tsypkin, Y. Z., 1980. "Adaptive Estimation Algorithms (Convergence, Optimality, Stability)," *Automation and Remote Control*, Vol. 40, pp. 378-389.

[PoT81] Poljak, B. T., and Tsypkin, Y. Z., 1981. "Optimal Pseudogradient Adaptation Algorithms," *Automation and Remote Control*, Vol. 41, pp. 1101-1110.

[PoT97] Polyak, R., and Teboulle, M., 1997. "Nonlinear Rescaling and Proximal-Like Methods in Convex Optimization," *Math. Programming*, Vol. 76, pp. 265-284.

- [Pol64] Poljak, B. T., 1964. "Some Methods of Speeding up the Convergence of Iteration Methods," *Z. Vyčisl. Mat. i Mat. Fiz.*, Vol. 4, pp. 1-17.
- [Pol71] Polak, E., 1971. *Computational Methods in Optimization: A Unified Approach*, Academic Press, NY.
- [Pol78] Poljak, B. T., 1978. "Nonlinear Programming Methods in the Presence of Noise," *Math. Programming*, Vol. 14, pp. 87-97.
- [Pol79] Poljak, B. T., 1979. "On Bertsekas' Method for Minimization of Composite Functions," *Internat. Symp. Systems Opt. Analysis*, Benoussan, A., and Lions, J. L., (Eds.), Springer-Verlag, Berlin and NY, pp. 179-186.
- [Pol87] Poljak, B. T., 1987. *Introduction to Optimization*, Optimization Software Inc., NY.
- [Pol88] Polyak, R. A., 1988. "Smooth Optimization Methods for Minimax Problems," *SIAM J. on Control and Optimization*, Vol. 26, pp. 1274-1286.
- [Pol92] Polyak, R. A., 1992. "Modified Barrier Functions (Theory and Methods)," *Math. Programming*, Vol. 54, pp. 177-222.
- [Pow69] Powell, M. J. D., 1969. "A Method for Nonlinear Constraints in Minimizing Problems," in *Optimization*, Fletcher, R., (Ed.), Academic Press, NY, pp. 283-298.
- [Pow73] Powell, M. J. D., 1973. "On Search Directions for Minimization Algorithms," *Math. Programming*, Vol. 4, pp. 193-201.
- [Pow11] Powell, W. B., 2011. *Approximate Dynamic Programming: Solving the Curses of Dimensionality*, 2nd Ed., Wiley, NY.
- [Pre95] Prekopa, A., 1995. *Stochastic Programming*, Kluwer, Boston.
- [Psh65] Pshenichnyi, B. N., 1965. "Dual Methods in Extremum Problems," *Kibernetika*, Vol. 1, pp. 89-95.
- [Pyt98] Pytlak, R., 1998. "An Efficient Algorithm for Large-Scale Nonlinear Programming Problems with Simple Bounds on the Variables," *SIAM J. on Optimization*, Vol. 8, pp. 532-560.
- [QSG13] Qin, Z., Scheinberg, K., and Goldfarb, D., 2013. "Efficient Block-Coordinate Descent Algorithms for the Group Lasso," *Math. Programming Computation*, Vol. 5, pp. 143-169.
- [RFP10] Recht, B., Fazel, M., and Parrilo, P. A., 2010. "Guaranteed Minimum-Rank Solutions of Linear Matrix Equations via Nuclear Norm Minimization," *SIAM Review*, Vol. 52, pp. 471-501.
- [RGV14] Richard, E., Gaïffas, S., and Vayatis, N., 2014. "Link Prediction in Graphs with Autoregressive Features," *J. of Machine Learning Research*, Vol. 15, pp. 565-593.
- [RHL13] Razaviyayn, M., Hong, M., and Luo, Z. Q., 2013. "A Unified Convergence Analysis of Block Successive Minimization Methods for Nonsmooth Optimization," *SIAM J. on Optimization*, Vol. 23, pp. 1126-1153.
- [RHW86] Rumelhart, D. E., Hinton, G. E., and Williams, R. J., 1986. "Learning Internal Representation by Error Backpropagation," in *Parallel*

- Distributed Processing-Explorations in the Microstructure of Cognition, by Rumelhart and McClelland, (eds.), MIT Press, Cambridge, MA, pp. 318-362.
- [RHW88] Rumelhart, D. E., Hinton, G. E., and Williams, R. J., 1988. "Learning Representations by Back-Propagating Errors," in *Cognitive Modeling*, by T. A. Polk, and C. M. Seifert, (eds.), MIT Press, Cambridge, MA, pp. 213-220.
- [RHZ14] Razaviyayn, M., Hong, M., and Luo, Z. Q., 2013. "A Unified Convergence Analysis of Block Successive Minimization Methods for Nonsmooth Optimization," *SIAM J. on Optimization*, Vol. 23, pp. 1126-1153.
- [RNV09] Ram, S. S., Nedić, A., and Veeravalli, V. V., 2009. "Incremental Stochastic Subgradient Algorithms for Convex Optimization," *SIAM J. on Optimization*, Vol. 20, pp. 691-717.
- [RNV10] Ram, S. S., Nedić, A., and Veeravalli, V. V., 2010. "Distributed Stochastic Subgradient Projection Algorithms for Convex Optimization," *J. of Optimization Theory and Applications*, Vol. 147, pp. 516-545.
- [ROF92] Rudin, L. I., Osher, S., and Fatemi, E., 1992. "Nonlinear Total Variation Based Noise Removal Algorithms," *Physica D: Nonlinear Phenomena*, Vol. 60, pp. 259-268.
- [RRW11] Recht, B., Re, C., Wright, S. J., and Niu, F., 2011. "Hogwild: A Lock-Free Approach to Parallelizing Stochastic Gradient Descent," in *Advances in Neural Information Processing Systems (NIPS 2011)*, pp. 693-701.
- [RSW13] Rao, N., Shah, P., Wright, S., and Nowak, R., 2013. "A Greedy Forward-Backward Algorithm for Atomic Norm Constrained Minimization," in *Proc. 2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 5885-5889.
- [RTV06] Roos, C., Terlaky, T., and Vial, J. P., 2006. *Interior Point Methods for Linear Optimization*, Springer, NY.
- [RXB11] Recht, B., Xu, W., and Hassibi, B., 2011. "Null Space Conditions and Thresholds for Rank Minimization," *Math. Programming*, Vol. 127, pp. 175-202.
- [RaN04] Rabbat, M. G., and Nowak, R. D., 2004. "Distributed Optimization in Sensor Networks," in *Proc. Inf. Processing Sensor Networks*, Berkeley, CA, pp. 20-27.
- [RaN05] Rabbat M. G., and Nowak R. D., 2005. "Quantized Incremental Algorithms for Distributed Optimization," *IEEE J. on Select Areas in Communications*, Vol. 23, pp. 798-808.
- [Ray93] Raydan, M., 1993. "On the Barzilai and Borwein Choice of Steplength for the Gradient Method," *IMA J. of Numerical Analysis*, Vol. 13, pp. 321-326.
- [Ray97] Raydan, M., 1997. "The Barzilai and Borwein Gradient Method for the Large Scale Unconstrained Minimization Problem," *SIAM J. on*

Optimization, Vol. 7, pp. 26-33.

[ReR13] Recht, B., and Ré, C., 2013. "Parallel Stochastic Gradient Algorithms for Large-Scale Matrix Completion," *Math. Programming Computation*, Vol. 5, pp. 201-226.

[Rec11] Recht, B., 2011. "A Simpler Approach to Matrix Completion," *The J. of Machine Learning Research*, Vol. 12, pp. 3413-3430.

[RiT14] Richtarik, P., and Takac, M., 2014. "Iteration Complexity of Randomized Block-Coordinate Descent Methods for Minimizing a Composite Function," *Math. Programming*, Vol. 144, pp. 1-38.

[RoS71] Robbins, H., and Siegmund, D. O., 1971. "A Convergence Theorem for Nonnegative Almost Supermartingales and Some Applications," *Optimizing Methods in Statistics*, pp. 233-257; see "Herbert Robbins Selected Papers," Springer, NY, 1985, pp. 111-135.

[RoW91] Rockafellar, R. T., and Wets, R. J.-B., 1991. "Scenarios and Policy Aggregation in Optimization Under Uncertainty," *Math. of Operations Research*, Vol. 16, pp. 119-147.

[RoW98] Rockafellar, R. T., and Wets, R. J.-B., 1998. *Variational Analysis*, Springer-Verlag, Berlin.

[Rob99] Robinson, S. M., 1999. "Linear Convergence of Epsilon-Subgradient Descent Methods for a Class of Convex Functions," *Math. Programming, Ser. A*, Vol. 86, pp. 41-50.

[Roc66] Rockafellar, R. T., 1966. "Characterization of the Subdifferentials of Convex Functions," *Pacific J. of Mathematics*, Vol. 17, pp. 497-510.

[Roc70] Rockafellar, R. T., 1970. *Convex Analysis*, Princeton Univ. Press, Princeton, NJ.

[Roc73] Rockafellar, R. T., 1973. "A Dual Approach to Solving Nonlinear Programming Problems by Unconstrained Optimization," *Math. Programming*, pp. 354-373.

[Roc76a] Rockafellar, R. T., 1976. "Monotone Operators and the Proximal Point Algorithm," *SIAM J. on Control and Optimization*, Vol. 14, pp. 877-898.

[Roc76b] Rockafellar, R. T., 1976. "Augmented Lagrangians and Applications of the Proximal Point Algorithm in Convex Programming," *Math. of Operations Research*, Vol. 1, pp. 97-116.

[Roc76c] Rockafellar, R. T., 1976. "Solving a Nonlinear Programming Problem by Way of a Dual Problem," *Symp. Matematica*, Vol. 27, pp. 135-160.

[Roc84] Rockafellar, R. T., 1984. *Network Flows and Monotropic Optimization*, Wiley, NY; republished by Athena Scientific, Belmont, MA, 1998.

[Rud76] Rudin, W., 1976. *Real Analysis*, McGraw-Hill, NY.

[Rup85] Ruppert, D., 1985. "A Newton-Raphson Version of the Multivariate Robbins-Monro Procedure," *The Annals of Statistics*, Vol. 13, pp. 236-245.

- [Rus86] Ruszczyński, A., 1986. “A Regularized Decomposition Method for Minimizing a Sum of Polyhedral Functions,” *Math. Programming*, Vol. 35, pp. 309-333.
- [Rus06] Ruszczyński, A., 2006. *Nonlinear Optimization*, Princeton Univ. Press, Princeton, NJ.
- [SBC91] Saarinen, S., Bramley, R. B., and Cybenko, G., 1991. “Neural Networks, Backpropagation and Automatic Differentiation,” in *Automatic Differentiation of Algorithms*, by A. Griewank and G. F. Corliss, (eds.), SIAM, Philadelphia, PA, pp. 31-42.
- [SBK64] Shah, B., Buehler, R., and Kempthorne, O., 1964. “Some Algorithms for Minimizing a Function of Several Variables,” *J. Soc. Indust. Appl. Math.*, Vol. 12, pp. 74-92.
- [SBT12] Shah, P., Bhaskar, B. N., Tang, G., and Recht, B., 2012. “Linear System Identification via Atomic Norm Regularization,” *arXiv preprint arXiv:1204.0590*.
- [SDR09] Shapiro, A., Dentcheva, D., and Ruszczyński, A., 2009. *Lectures on Stochastic Programming: Modeling and Theory*, SIAM, Phila., PA.
- [SFR09] Schmidt, M., Fung, G., and Rosales, R., 2009. “Optimization Methods for ℓ_1 -Regularization,” Univ. of British Columbia, Technical Report TR-2009-19.
- [SKS12] Schmidt, M., Kim, D., and Sra, S., 2012. “Projected Newton-Type Methods in Machine Learning,” in *Optimization for Machine Learning*, by Sra, S., Nowozin, S., and Wright, S. J., (eds.), MIT Press, Cambridge, MA, pp. 305-329.
- [SLB13] Schmidt, M., Le Roux, N., and Bach, F., 2013. “Minimizing Finite Sums with the Stochastic Average Gradient,” *arXiv preprint arXiv:1309.2388*.
- [SNW12] Sra, S., Nowozin, S., and Wright, S. J., 2012. *Optimization for Machine Learning*, MIT Press, Cambridge, MA.
- [SRB11] Schmidt, M., Roux, N. L., and Bach, F. R., 2011. “Convergence Rates of Inexact Proximal-Gradient Methods for Convex Optimization,” In *Advances in Neural Information Processing Systems*, pp. 1458-1466.
- [SSS07] Shalev-Shwartz, S., Singer, Y., Srebro, N., and Cotter, A., 2007. “Pegasos: Primal Estimated Subgradient Solver for SVM,” in *ICML 07*, New York, NY, pp. 807-814.
- [SaB96] Savari, S. A., and Bertsekas, D. P., 1996. “Finite Termination of Asynchronous Iterative Algorithms,” *Parallel Computing*, Vol. 22, pp. 39-56.
- [SaT13] Saha, A., and Tewari, A., 2013. “On the Nonasymptotic Convergence of Cyclic Coordinate Descent Methods,” *SIAM J. on Optimization*, Vol. 23, pp. 576-601.
- [Sak66] Sakrisson, D. T., 1966. “Stochastic Approximation: A Recursive

- Method for Solving Regression Problems,” in *Advances in Communication Theory and Applications*, 2, A. V. Balakrishnan, ed., Academic Press, NY, pp. 51-106.
- [Say14] Sayed, A. H., 2014. “Adaptation, Learning, and Optimization over Networks,” *Foundations and Trends in Machine Learning*, Vol. 7, no. 4-5, pp. 311-801.
- [ScF14] Schmidt, M., and Friedlander, M. P., 2014. “Coordinate Descent Converges Faster with the Gauss-Southwell Rule than Random Selection,” *Advances in Neural Information Processing Systems 27 (NIPS 2014)*.
- [Sch82] Schnabel, R. B., 1982. “Determining Feasibility of a Set of Non-linear Inequality Constraints,” *Math. Programming Studies*, Vol. 16, pp. 137-148.
- [Sch86] Schrijver, A., 1986. *Theory of Linear and Integer Programming*, Wiley, NY.
- [Sch10] Schmidt, M., 2010. “Graphical Model Structure Learning with L1-Regularization,” PhD Thesis, Univ. of British Columbia.
- [Sch14a] Schmidt, M., 2014. “Convergence Rate of Stochastic Gradient with Constant Step Size,” Computer Science Report, Univ. of British Columbia.
- [Sch14b] Schmidt, M., 2014. “Convergence Rate of Proximal Gradient with General Step-Size,” Dept. of Computer Science, Unpublished Note, Univ. of British Columbia.
- [ShZ12] Shamir, O., and Zhang, T., 2012. “Stochastic Gradient Descent for Non-Smooth Optimization: Convergence Results and Optimal Averaging Schemes,” arXiv preprint arXiv:1212.1824.
- [Sha79] Shapiro, J. E., 1979. *Mathematical Programming Structures and Algorithms*, Wiley, NY.
- [Sho85] Shor, N. Z., 1985. *Minimization Methods for Nondifferentiable Functions*, Springer-Verlag, Berlin.
- [Sho98] Shor, N. Z., 1998. *Nondifferentiable Optimization and Polynomial Problems*, Kluwer Academic Publishers, Dordrecht, Netherlands.
- [SmS04] Smola, A. J., and Scholkopf, B., 2004. “A Tutorial on Support Vector Regression,” *Statistics and Computing*, Vol. 14, pp. 199-222.
- [SoZ98] Solodov, M. V., and Zavriev, S. K., 1998. “Error Stability Properties of Generalized Gradient-Type Algorithms,” *J. Opt. Theory and Appl.*, Vol. 98, pp. 663-680.
- [Sol98] Solodov, M. V., 1998. “Incremental Gradient Algorithms with Stepsizes Bounded Away from Zero,” *Computational Optimization and Applications*, Vol. 11, pp. 23-35.
- [Spa03] Spall, J. C., 2003. *Introduction to Stochastic Search and Optimization: Estimation, Simulation, and Control*, J. Wiley, Hoboken, NJ.
- [Spa12] Spall, J. C., 2012. “Cyclic Seesaw Process for Optimization and Identification,” *J. of Optimization Theory and Applications*, Vol. 154, pp.

187-208.

- [Spi83] Spingarn, J. E., 1983. "Partial Inverse of a Monotone Operator," *Applied Mathematics and Optimization*, Vol. 10, pp. 247-265.
- [Spi85] Spingarn, J. E., 1985. "Applications of the Method of Partial Inverses to Convex Programming: Decomposition," *Math. Programming*, Vol. 32, pp. 199-223.
- [StV09] Strohmer, T., and Vershynin, R., 2009. "A Randomized Kaczmarz Algorithm with Exponential Convergence," *J. Fourier Anal. Appl.*, Vol. 15, pp. 262-278.
- [StW70] Stoer, J., and Witzgall, C., 1970. *Convexity and Optimization in Finite Dimensions*, Springer-Verlag, Berlin.
- [StW75] Stephanopoulos, G., and Westerberg, A. W., 1975. "The Use of Hestenes' Method of Multipliers to Resolve Dual Gaps in Engineering System Optimization," *J. Optimization Theory and Applications*, Vol. 15, pp. 285-309.
- [Str76] Strang, G., 1976. *Linear Algebra and Its Applications*, Academic Press, NY.
- [Str97] Stromberg, A-B., 1997. *Conditional Subgradient Methods and Ergodic Convergence in Nonsmooth Optimization*, Ph.D. Thesis, Univ. of Linköping, Sweden.
- [SuB98] Sutton, R. S., and Barto, A. G., 1998. *Reinforcement Learning*, MIT Press, Cambridge, MA.
- [TBA86] Tsitsiklis, J. N., Bertsekas, D. P., and Athans, M., 1986. "Distributed Asynchronous Deterministic and Stochastic Gradient Optimization Algorithms," *IEEE Trans. on Automatic Control*, Vol. AC-31, pp. 803-812.
- [TBT90] Tseng, P., Bertsekas, D. P., and Tsitsiklis, J. N., 1990. "Partially Asynchronous, Parallel Algorithms for Network Flow and Other Problems," *SIAM J. on Control and Optimization*, Vol. 28, pp. 678-710.
- [TVS10] Teo, C. H., Vishwanathan, S. V. N., Smola, A. J., and Le, Q. V., 2010. "Bundle Methods for Regularized Risk Minimization," *The J. of Machine Learning Research*, Vol. 11, pp. 311-365.
- [TaP13] Talischi, C., and Paulino, G. H., 2013. "A Consistent Operator Splitting Algorithm and a Two-Metric Variant: Application to Topology Optimization," *arXiv preprint arXiv:1307.5100*.
- [Teb92] Teboulle, M., 1992. "Entropic Proximal Mappings with Applications to Nonlinear Programming," *Math. of Operations Research*, Vol. 17, pp. 1-21.
- [Teb97] Teboulle, M., 1997. "Convergence of Proximal-Like Algorithms," *SIAM J. Optim.*, Vol. 7, pp. 1069-1083.
- [Ter96] Terlaky, T. (Ed.), 1996. *Interior Point Methods of Mathematical Programming*, Springer, NY.

- [Tib96] Tibshirani, R., 1996. "Regression Shrinkage and Selection via the Lasso," *J. of the Royal Statistical Society, Series B (Methodological)*, Vol. 58, pp. 267-288.
- [Tod01] Todd, M. J., 2001. "Semidefinite Optimization," *Acta Numerica*, Vol. 10, pp. 515-560.
- [TsB87] Tseng, P., and Bertsekas, D. P., 1987. "Relaxation Methods for Problems with Strictly Convex Separable Costs and Linear Constraints," *Math. Programming*, Vol. 38, pp. 303-321.
- [TsB90] Tseng, P., and Bertsekas, D. P., 1990. "Relaxation Methods for Monotropic Programs," *Math. Programming*, Vol. 46, pp. 127-151.
- [TsB91] Tseng, P., and Bertsekas, D. P., 1991. "Relaxation Methods for Problems with Strictly Convex Costs and Linear Constraints," *Math. of Operations Research*, Vol. 16, pp. 462-481.
- [TsB93] Tseng, P., and Bertsekas, D. P., 1993. "On the Convergence of the Exponential Multiplier Method for Convex Programming," *Math. Programming*, Vol. 60, pp. 1-19.
- [TsB00] Tseng, P., and Bertsekas, D. P., 2000. "An Epsilon-Relaxation Method for Separable Convex Cost Generalized Network Flow Problems," *Math. Programming*, Vol. 88, pp. 85-104.
- [TsY09] Tseng, P. and Yun S., 2009. "A Coordinate Gradient Descent Method for Nonsmooth Separable Minimization," *Math. Programming*, Vol. 117, pp. 387-423.
- [Tse91a] Tseng, P., 1991. "Decomposition Algorithm for Convex Differentiable Minimization," *J. of Optimization Theory and Applications*, Vol. 70, pp. 109-135.
- [Tse91b] Tseng, P., 1991. "Applications of a Splitting Algorithm to Decomposition in Convex Programming and Variational Inequalities," *SIAM J. on Control and Optimization*, Vol. 29, pp. 119-138.
- [Tse93] Tseng, P., 1993. "Dual Coordinate Ascent Methods for Non-Strictly Convex Minimization," *Math. Programming*, Vol. 59, pp. 231-247.
- [Tse95] Tseng, P., 1995. "Fortified-Descent Simplicial Search Method," Report, Dept. of Math., Univ. of Washington, Seattle, Wash.; also in *SIAM J. on Optimization*, Vol. 10, 1999, pp. 269-288.
- [Tse98] Tseng, P., 1998. "Incremental Gradient(-Projection) Method with Momentum Term and Adaptive Stepsize Rule," *SIAM J. on Control and Optimization*, Vol. 8, pp. 506-531.
- [Tse00] Tseng, P., 2000. "A Modified Forward-Backward Splitting Method for Maximal Monotone Mappings," *SIAM J. on Control and Optimization*, Vol. 38, pp. 431-446.
- [Tse01a] Tseng, P., 2001. "Convergence of Block Coordinate Descent Methods for Nondifferentiable Minimization," *J. Optim. Theory Appl.*, Vol. 109, pp. 475-494.

- [Tse01b] Tseng, P., 2001. “An Epsilon Out-of-Kilter Method for Monotropic Programming,” *Math. of Operations Research*, Vol. 26, pp. 221-233.
- [Tse04] Tseng, P., 2004. “An Analysis of the EM Algorithm and Entropy-Like Proximal Point Methods,” *Math. Operations Research*, Vol. 29, pp. 27-44.
- [Tse08] Tseng, P., 2008. “On Accelerated Proximal Gradient Methods for Convex-Concave Optimization,” Report, Math. Dept., U. of Washington.
- [Tse09] Tseng, P., 2009. “Some Convex Programs Without a Duality Gap,” *Math. Programming*, Vol. 116, pp. 553-578.
- [Tse10] Tseng, P., 2010. “Approximation Accuracy, Gradient Methods, and Error Bound for Structured Convex Optimization,” *Math. Programming*, Vol. 125, pp. 263-295.
- [VKG14] Vetterli, M., Kovacevic, J., and Goyal, V. K., 2014. *Foundations of Signal Processing*, Cambridge Univ. Press, NY.
- [VMR88] Vogl, T. P., Mangis, J. K., Rigler, A. K., Zink, W. T., and Alkon, D. L., 1988. “Accelerating the Convergence of the Back-Propagation Method,” *Biological Cybernetics*, Vol. 59, pp. 257-263.
- [VaF08] Van Den Berg, E., and Friedlander, M. P., 2008. “Probing the Pareto Frontier for Basis Pursuit Solutions,” *SIAM J. on Scientific Computing*, Vol. 31, pp. 890-912.
- [Van01] Vanderbei, R. J., 2001. *Linear Programming: Foundations and Extensions*, Springer, NY.
- [VeH93] Ventura, J. A., and Hearn, D. W., 1993. “Restricted Simplicial Decomposition for Convex Constrained Problems,” *Math. Programming*, Vol. 59, pp. 71-85.
- [Ven67] Venter, J. H., 1967. “An Extension of the Robbins-Monro Procedure,” *Ann. Math. Statist.*, Vol. 38, pp. 181-190.
- [WDS13] Weinmann, A., Demaret, L., and Storath, M., 2013. “Total Variation Regularization for Manifold-Valued Data,” *arXiv preprint arXiv:1312.7710*.
- [WFL14] Wang, M., Fang, E., and Liu, H., 2014. “Stochastic Compositional Gradient Descent: Algorithms for Minimizing Compositions of Expected-Value Functions,” *Optimization Online*.
- [WHM13] Wang, X., Hong, M., Ma, S., Luo, Z. Q., 2013. “Solving Multiple-Block Separable Convex Minimization Problems Using Two-Block Alternating Direction Method of Multipliers,” *arXiv preprint arXiv:1308.5294*.
- [WSK14] Wytock, M., Sra, S., and Kolter, J. K., 2014. “Fast Newton Methods for the Group Fused Lasso,” *Proc. of 2014 Conf. on Uncertainty in Artificial Intelligence*.
- [WSV00] Wolkowicz, H., Saigal, R., and Vanderbergue, L., (eds), 2000. *Handbook of Semidefinite Programming*, Kluwer, Boston.
- [WaB13a] Wang, M., and Bertsekas, D. P., 2013. “Incremental Constraint

- Projection-Proximal Methods for Nonsmooth Convex Optimization,” Lab. for Information and Decision Systems Report LIDS-P-2907, MIT, to appear in SIAM J. on Optimization.
- [WaB13b] Wang, M., and Bertsekas, D. P., 2013. “Convergence of Iterative Simulation-Based Methods for Singular Linear Systems,” *Stochastic Systems*, Vol. 3, pp. 38-95.
- [WaB13c] Wang, M., and Bertsekas, D. P., 2013. “Stabilization of Stochastic Iterative Methods for Singular and Nearly Singular Linear Systems,” *Math. of Operations Research*, Vol. 39, pp. 1-30.
- [WaB14] Wang, M., and Bertsekas, D. P., 2014. “Incremental Constraint Projection Methods for Variational Inequalities,” *Mathematical Programming*, pp. 1-43.
- [Was04] Wasserman, L., 2004. *All of Statistics: A Concise Course in Statistical Inference*, Springer, NY.
- [Wat92] Watson, G. A., 1992. “Characterization of the Subdifferential of Some Matrix Norms,” *Linear Algebra and its Applications*, Vol. 170, pp. 33-45.
- [WeO13] Wei, E., and Ozdaglar, A., 2013. “On the $O(1/k)$ Convergence of Asynchronous Distributed Alternating Direction Method of Multipliers,” arXiv preprint arXiv:1307.8254.
- [WiH60] Widrow, B., and Hoff, M. E., 1960. “Adaptive Switching Circuits,” *Institute of Radio Engineers, Western Electronic Show and Convention, Convention Record, Part 4*, pp. 96-104.
- [Wol75] Wolfe, P., 1975. “A Method of Conjugate Subgradients for Minimizing Nondifferentiable Functions,” *Math. Programming Study 3*, Balinski, M., and Wolfe, P., (Eds.), North-Holland, Amsterdam, pp. 145-173.
- [Wri97] Wright, S. J., 1997. *Primal-Dual Interior Point Methods*, SIAM, Philadelphia, PA.
- [Wri14] Wright, S. J., 2014. “Coordinate Descent Algorithms,” *Optimization Online*.
- [XiZ14] Xiao, L., and Zhang, T., 2014. “A Proximal Stochastic Gradient Method with Progressive Variance Reduction,” arXiv preprint arXiv:1403.4699.
- [Xia10] Xiao L., 2010. “Dual Averaging Methods for Regularized Stochastic Learning and Online Optimization,” *J. of Machine Learning Research*, Vol. 11, pp. 2534-2596.
- [YBR08] Yu, H., Bertsekas, D. P., and Rousu, J., 2008. “An Efficient Discriminative Training Method for Generative Models,” *Extended Abstract, the 6th International Workshop on Mining and Learning with Graphs (MLG)*.
- [YGT93] Ye, Y., Guler, O., Tapia, R. A., and Zhang, Y., 1993. “A Quadratically Convergent $O(\sqrt{nL})$ -Iteration Algorithm for Linear Programming,”

- Math. Programming, Vol. 59, pp. 151-162.
- [YNS10] Yousefian, F., Nedić, A., and Shanbhag, U. V., 2010. "Convex Nondifferentiable Stochastic Optimization: A Local Randomized Smoothing Technique," Proc. American Control Conference (ACC), pp. 4875-4880.
- [YNS12] Yousefian, F., Nedić, A., and Shanbhag, U. V., 2012. "On Stochastic Gradient and Subgradient Methods with Adaptive Steplength Sequences," Automatica, Vol. 48, pp. 56-67.
- [YOG08] Yin, W., Osher, S., Goldfarb, D., and Darbon, J., 2008. "Bregman Iterative Algorithms for ℓ_1 -Minimization with Applications to Compressed Sensing," SIAM J. on Imaging Sciences, Vol. 1, pp. 143-168.
- [YSQ14] You, K., Song, S., and Qiu, L., 2014. "Randomized Incremental Least Squares for Distributed Estimation Over Sensor Networks," Preprints of the 19th World Congress The International Federation of Automatic Control Cape Town, South Africa.
- [Ye92] Ye, Y., 1992. "A Potential Reduction Algorithm Allowing Column Generation," SIAM J. on Optimization, Vol. 2, pp. 7-20.
- [Ye97] Ye, Y., 1997. Interior Point Algorithms: Theory and Analysis, Wiley Interscience, NY.
- [YuR07] Yu, H., and Rousu, J., 2007. "An Efficient Method for Large Margin Parameter Optimization in Structured Prediction Problems," Technical Report C-2007-87, Univ. of Helsinki.
- [ZJL13] Zhang, H., Jiang, J., and Luo, Z. Q., 2013. "On the Linear Convergence of a Proximal Gradient Method for a Class of Nonsmooth Convex Minimization Problems," J. of the Operations Research Society of China, Vol. 1, pp. 163-186.
- [ZLW99] Zhao, X., Luh, P. B., and Wang, J., 1999. "Surrogate Gradient Algorithm for Lagrangian Relaxation," J. Optimization Theory and Applications, Vol. 100, pp. 699-712.
- [ZMJ13] Zhang, L., Mahdavi, M., and Jin, R., 2013. "Linear Convergence with Condition Number Independent Access of Full Gradients," Advances in Neural Information Processing Systems 26 (NIPS 2013), pp. 980-988.
- [ZTD92] Zhang, Y., Tapia, R. A., and Dennis, J. E., 1992. "On the Superlinear and Quadratic Convergence of Primal-Dual Interior Point Linear Programming Algorithms," SIAM J. on Optimization, Vol. 2, pp. 304-324.
- [Zal02] Zalinescu, C., 2002. Convex Analysis in General Vector Spaces, World Scientific, Singapore.
- [Zan69] Zangwill, W. I., 1969. Nonlinear Programming, Prentice-Hall, Englewood Cliffs, NJ.
- [Zou60] Zoutendijk, G., 1960. Methods of Feasible Directions, Elsevier Publ. Co., Amsterdam.
- [Zou76] Zoutendijk, G., 1976. Mathematical Programming Methods, North Holland, Amsterdam.

INDEX

A

ADMM 111, 280, 292, 295, 298, 337, 427
Affine function 445
Affine hull 472
Affine set 445
Aggregated gradient method 91, 94, 428
Alternating direction method 111, 280
Analytic center 425
Approximation algorithms 36, 54
Armijo rule 69, 123, 125, 317
Asymptotic sequence 481
Asynchronous computation 33, 104, 376
Asynchronous gradient method 104, 106
Atomic norm 35
Auction algorithm 180, 375
Augmented Lagrangian function 260, 283, 290
Augmented Lagrangian method 109, 115, 120, 261, 294, 326, 337, 362, 384, 389, 430
Averaging of iterates 120, 157, 177

B

Backpropagation 119
Backtracking rule 69, 123
Ball center 426
Barrier function 412
Barrier method 413
Basis 446
Basis pursuit 34, 286
Batching 93
Block coordinate descent 75, 268, 281, 429, 438-442

Bolzano-Weierstrass Theorem 453
Boundary of a set 454
Boundary point 454
Bounded sequence 451, 452
Bounded set 453
Branch-and-bound 7
Bregman distance 388
Bundle methods 110, 187, 272, 295, 385

C

Caratheodory's Theorem 472
Cartesian product 445
Cauchy sequence 452
Central cutting plane methods 425, 432
Chain rule 142, 513
Classification 29
Closed ball 453
Closed function 469
Closed halfspace 484
Closed set 453
Closed set intersections 481
Closed sphere 453
Closedness under linear transformations 483
Closedness under vector sums 483
Closure of a function 476
Closure of a set 453
Closure point 453
Co-finite function 411
Coercive function 496
Compact set 453
Complementary slackness 508
Component of a vector 444
Composition of functions 142, 455, 469, 514
Concave closure 502

Concave function 468
 Condition number 60, 122, 315
 Conditional gradient method 71, 107, 191, 192, 374
 Cone 468
 Cone decomposition 219
 Cone generated by a set 472, 489, 492
 Confusion region 89, 93
 Conic duality 14, 19, 23, 511
 Conic programming 13, 217, 224, 231, 232, 423, 432, 511
 Conjugate Subgradient Theorem 201, 513
 Conjugacy Theorem 487
 Conjugate direction method 64, 68, 320
 Conjugate function 487
 Conjugate gradient method 64
 Constancy space 480
 Constant stepsize rule 56, 59, 69, 153, 304, 349
 Constraint aggregation 230
 Constraint qualification 507
 Continuity 455, 475
 Continuous differentiability 141, 172, 457
 Contraction mapping 56, 296, 312, 458
 Convergent sequence 451, 452
 Convex closure 476
 Convex combination 472
 Convex function 468, 469
 Convex hull 472
 Convex programming 4, 507
 Convex set 468
 Convexification of a function 476
 Coordinate 444
 Coordinate descent method 75, 104, 369, 429, 439-442
 Crossing function 499
 Cutting plane method 107, 182, 211, 228, 270
 Cyclic coordinate descent 76, 370, 376

Cyclic incremental method 84, 96, 166, 343

D

Danskin's Theorem 146, 172
 Dantzig-Wolfe decomposition 107, 111, 229, 295
 Decomposition algorithm 7, 77, 289, 363
 Decomposition of a convex set 479
 Derivative 456
 Descent algorithms 54
 Descent direction 58, 71
 Descent inequality 122, 305
 Diagonal dominance 106
 Diagonal scaling 63, 101, 333, 338
 Differentiability 457
 Differentiable convex function 470
 Differentiation theorems 457, 458, 513, 514
 Dimension of a convex set 472
 Dimension of a subspace 446
 Dimension of an affine set 472
 Diminishing stepsize rule 69, 127, 157, 174, 316, 351
 Direct search methods 83
 Direction of recession 478
 Directional derivative 137, 170-173, 515
 Distributed computation 8, 33, 104, 365, 376
 Domain 444
 Domain one-dimensional 409
 Double conjugate 487
 Dual cone 14, 511
 Dual function 3, 499
 Dual pair representation 219
 Dual problem 2, 147, 164, 499, 506, 507
 Dual proximal algorithm 257, 336, 384
 Duality gap estimate 9
 Duality theory 2, 498, 505
 Dynamic programming 36, 380

Dynamic stepsize rule 159, 175, 177

E

ϵ -complementary slackness 180
 ϵ -descent algorithm 83, 396, 400, 431
 ϵ -descent direction 400
 ϵ -relaxation method 375
 ϵ -subdifferential 162, 397
 ϵ -subgradient 162, 164, 169, 180, 397
 ϵ -subgradient method 162, 179
 EMP 197, 220, 396, 406
 Effective domain 468
 Entropic descent 396
 Entropy function 383
 Entropy minimization algorithm 383
 Epigraph 468
 Essentially one-dimensional 408
 Euclidean norm 450
 Eventually constant stepsize rule 308, 334
 Exact penalty 39, 45, 365, 369
 Existence of dual optimal solutions 503
 Existence of optimal solutions 483, 495
 Exponential augmented Lagrangian method 116, 134, 389, 431
 Exponential loss 30
 Exponential smoothing 116, 134, 391
 Extended Kalman filter 103, 120
 Extended monotropic programming 83, 197, 220, 229, 396, 406, 431
 Extended real number 443
 Extended real-valued function 468
 Exterior penalty method 110
 Extrapolation 63-66, 322, 338, 427, 428
 Extreme point 490

F

Farkas' Lemma 492, 505, 506

Farout region 89
 Feasibility problem 34, 283, 429
 Feasible direction 71
 Feasible direction methods 71
 Feasible solution 2, 494
 Fejér Convergence Theorem 126, 158, 465
 Fejér monotonicity 464
 Fenchel duality 10, 510
 Fenchel inequality 512
 Finitely generated cone 492
 Forward-backward algorithm 427
 Forward image 445
 Frank-Wolfe algorithm 71, 107, 191, 374
 Fritz John optimality conditions 6
 Full rank 447
 Fundamental Theorem of Linear Programming 494

G

GPA algorithm 200-202, 229
 Gauss-Southwell order 376, 441
 Generalized polyhedral approximation 107, 196, 201, 229
 Generalized simplicial decomposition 209, 229
 Generated cone 472, 489, 492
 Geometric convergence 57
 Global maximum 495
 Global minimum 495
 Gradient 456
 Gradient method 56, 59
 Gradient method distributed 104, 106
 Gradient method with momentum 63, 92
 Gradient projection 73, 82, 136, 302, 374, 385, 396, 427, 434

H

Halfspace 484
 Heavy ball method 63, 92
 Hessian matrix 457
 Hierarchical decomposition 77

Hinge loss 30
 Hyperplane 484
 Hyperplane separation 484-487

I

Ill-conditioning 60, 109, 413
 Image 445, 446
 Improper function 468
 Incremental Gauss-Newton method 103, 120
 Incremental Newton method 97, 101, 118, 119
 Incremental aggregated method 91, 94
 Incremental constraint projection method 102, 365, 429
 Incremental gradient method 84, 105, 118, 119, 130-132
 Incremental gradient with momentum 92
 Incremental method 25, 83, 166, 320
 Incremental proximal method 341, 385, 429
 Incremental subgradient method 84, 166, 341, 385, 428
 Indicator function 487
 Inner linearization 107, 182, 188, 194, 296
 Infeasible problem 494
 Infimum 444
 Inner approximation 402
 Inner product 444
 Instability 186, 191, 269
 Integer programming 6
 Interior of a set 412, 453
 Interior point 453
 Interior point method 108, 412, 415, 423, 432
 Interpolated iteration 249, 253, 298, 459
 Inverse barrier 412
 Inverse image 445, 446

J**K**

Kaczmarz method 85, 98, 131
 Krasnosel'skii-Mann Theorem 252, 285, 300, 459

L

ℓ_1 -norm 451
 ℓ_∞ -norm 450
 LMS method 119
 Lagrange multiplier 507
 Lagrangian function 3, 507
 Lasso problem 27
 Least absolute value deviations 27, 288
 Least mean squares method 119
 Left-continuous function 455
 Level set 469
 Limit 451
 Limit point 451, 453
 Limited memory quasi-Newton method 63, 338
 Line minimization 60, 65, 69, 320
 Line segment principle 473
 Lineality space 479
 Linear-conic problems 15, 16
 Linear convergence 57
 Linear equation 445
 Linear function 445
 Linear inequality 445
 Linear programming 16, 415, 434
 Linear programming duality 506
 Linear regularity 369
 Linear transformation preservation of closedness 483
 Linearly independent vectors 446
 Lipschitz continuity 141, 455, 512
 Local convergence 68
 Local maximum 495
 Local minimum 495
 Location theory 32
 Logarithmic barrier 412, 416
 Logistic loss 30
 Lower limit 452

Lower semicontinuous function 455, 469

M

Majorization-maximization algorithm 392

Matrix completion 28, 35

Matrix factorization 30, 373

Max crossing problem 499

Max function 470

Maximal monotone mapping 255

Maximum likelihood 31

Maximum norm 450

Maximum point 444, 495

Mean Value Theorem 457, 458

Merit function 54, 417

Min common problem 499

Min common/max crossing framework 499

Minimax duality 502, 516

Minimax duality gap 9

Minimax duality theorems 516

Minimax equality 12, 516-518

Minimax problems 9, 12, 35, 113, 147, 164, 215, 217

Minimax theory 498, 502, 516

Minimizer 495

Minimum point 444, 495

Minkowski-Weyl Theorem 492

Minkowski-Weyl representation 492

Mirror descent 82, 385, 395

Momentum term 63, 92

Monotone mapping 255

Monotonically nondecreasing sequence 451

Monotonically nonincreasing sequence 451

Monotropic programming 83, 197, 208, 431

Multicommodity flow 38, 193, 217

Multiplier method 109, 261, 267

N

Negative halfspace 484

Neighborhood 453

Nelder-Mead algorithm 83

Nested sequence 481

Network optimization 37, 189, 193, 208, 217, 375

Newton's method 67, 74, 75, 97, 338, 416, 424

Nonexpansive mapping 57, 249, 459

Nonlinear Farkas' Lemma 505

Nonmonotonic stepsize rule 70

Nonnegative combination 472

Nonquadratic regularization 242, 294, 382, 393

Nonsingular matrix 447

Nonstationary iteration 57, 461

Nonvertical Hyperplane Theorem 486

Nonvertical hyperplane 486

Norm 450

Norm equivalence 454

Normal cone 145

Normal of a hyperplane 484

Nuclear norm 28, 35

Null step 212

Nullspace 447

O

Open ball 453

Open halfspace 484

Open set 453

Open sphere 453

Optimality conditions 3, 144, 470, 508-511, 514

Orthogonal complement 446

Orthogonal vectors 444

Outer approximation 402

Outer linearization 107, 182, 194

Overrelaxation 253

P

PARTAN 64

Parallel subspace 445

Parallel projections method 76, 438

Parallel tangents method 64

Partitioning 9

Partial cutting plane method 187

Partial derivative 456
 Partial minimization 496
 Partial proximal algorithm 297
 Partially asynchronous algorithm 105, 377
 Penalty method 38, 108, 120, 326
 Penalty parameter 38, 40, 109
 Perturbation function 501
 Polar Cone Theorem 488, 492
 Polar cone 488, 491
 Polyhedral Proper Separation Theorem 486
 Polyhedral approximation 107, 182, 196, 217, 385
 Polyhedral cone 492
 Polyhedral function 493
 Polyhedral set 468, 489
 Positive combination 472
 Positive definite matrix 449
 Positive halfspace 484
 Positive semidefinite matrix 449
 Positively homogeneous function 488, 489
 Power regularization 393
 Predictor-corrector method 421, 432
 Primal-dual method 416, 420
 Primal function 501
 Primal problem 2, 499
 Projection Theorem 471
 Prolongation Lemma 473
 Proper Separation Theorem 486
 Proper function 468
 Properly separating hyperplane 485, 486
 Proximal Newton algorithm 338, 428
 Proximal algorithm 80, 110, 120, 234, 293, 307, 374, 384, 393, 439
 Proximal cutting plane method 270
 Proximal gradient algorithm 82, 112, 133, 330, 336, 385, 428, 436-438
 Proximal inner linearization 278
 Proximal operator 248, 296
 Proximal simplicial decomposition 280
 Pythagorean Theorem 450

Q

Quadratic penalty function 39
 Quadratic programming 21, 483
 Quasi-Newton method 68, 74, 338

R

Randomized coordinate descent 376, 440
 Randomized incremental method 84, 86, 131, 344, 353, 429
 Range 444
 Range space 447
 Rank 447
 Recession Cone Theorem 478
 Recession cone of a function 480
 Recession cone of a set 478
 Recession direction 478
 Recession function 480
 Reflection operator 249
 Regression 26
 Regularization 26-31, 117, 133, 232, 235, 287, 361, 382, 393
 Relative boundary 473
 Relative boundary point 473
 Relative interior 473
 Relative interior point 473
 Relatively open 473
 Relaxation method 375
 Reshuffling 96
 Restricted simplicial decomposition 192, 204
 Retractive sequence 481
 Retractive set 482
 Right-continuous function 455
 Robust optimization 20, 47

S

Saddle Point Theorem 518
 Saddle point 498, 517-518
 Scaling 61
 Schwarz inequality 450
 Second order cone programming 17, 49, 50, 230, 423
 Second order expansions 458

- Second order method of multipliers 267
 - Self-dual cone 17
 - Semidefinite programming 17, 22, 424
 - Sensitivity 501
 - Separable problems 7, 289, 362, 405
 - Separating Hyperplane Theorem 485
 - Separating hyperplane 485
 - Sequence definitions 451
 - Serious step 272
 - Set intersection theorems 482
 - Shapley-Folkman Theorem 9
 - Sharp minimum 80, 176, 179, 242-244, 436
 - Shrinkage operation 133, 287, 330, 362
 - Side constraints 213, 217
 - Simplicial decomposition method 72, 107, 182, 188, 193, 209, 221, 228-230, 278, 280, 320
 - Single commodity network flow 38, 208, 375
 - Singular matrix 447
 - Slater condition 8, 505
 - Smoothing 113, 168, 326, 427
 - Span 446
 - Sphere 453
 - Splitting algorithm 113, 120, 427
 - Square root of a matrix 450
 - Stationary iteration 56
 - Steepest descent direction 59, 128
 - Steepest descent method 59, 78, 401
 - Stepsize rules 69, 308, 316
 - Stochastic Newton method 120
 - Stochastic approximation 94
 - Stochastic gradient method 93, 167
 - Stochastic programming 32
 - Stochastic subgradient method 93, 167
 - Strict Separation Theorem 485
 - Strict local minimum 495
 - Strictly convex function 468
 - Strictly separating hyperplane 485
 - Strong convexity 312, 435, 440, 471
 - Strong duality 3, 499
 - Strong monotonicity 296
 - Subdifferential 136, 141, 143, 146, 167, 171, 511-514
 - Subgradient 136, 147, 148, 167, 511-514
 - Subgradient methods 78, 136, 179
 - Subsequence 452
 - Subspace 445
 - Superlinear convergence 68, 100, 338, 393, 437
 - Supermartingale Convergence 462-464
 - Support function 488
 - Support vector machine 30, 48
 - Supporting Hyperplane Theorem 484
 - Supporting hyperplane 484
 - Supremum 444
 - Symmetric matrix 449
- T**
- Theorems of the alternative 52
 - Tikhonov regularization 111
 - Total variation denoising 28
 - Totally asynchronous algorithm 105, 376
 - Triangle inequality 450
 - Trust region 70, 296
 - Two-metric projection method 74, 102, 129, 189, 374
- U**
- Uniform sampling 96, 344
 - Upper limit 452
 - Upper semicontinuous function 455, 469
- V**
- Vector sum preservation of closedness 483
 - Vertical hyperplane 486
- W**
- Weak Duality Theorem 3, 7, 499, 507

Weber point 32, 52

Weierstrass' Theorem 456, 495

Weighted sup-norm 380

Y

Z

Zero-sum games 9