

Algorithms for Convex Optimization

In the last few years, algorithms for convex optimization have revolutionized algorithm design, both for discrete and continuous optimization problems. For problems such as maximum flow, maximum matching, and submodular function minimization, the fastest algorithms involve essential methods such as gradient descent, mirror descent, interior point methods, and ellipsoid methods. The goal of this self-contained book is to enable researchers and professionals in computer science, operations research, data science, and machine learning to gain an in-depth understanding of these algorithms. The text emphasizes how to derive key algorithms for convex optimization from first principles and how to establish precise running time bounds. This modern text explains the success of these algorithms in problems of discrete optimization, as well as how these methods have significantly pushed the state of the art of convex optimization itself.

NISHEETH K. VISHNOI is A. Bartlett Giamatti Professor of Computer Science at Yale University. His research areas include theoretical computer science, optimization, and machine learning. He is a recipient of the Best Paper Award at IEEE FOCS in 2005, the IBM Research Pat Goldberg Memorial Award in 2006, the Indian National Science Academy Young Scientist Award in 2011, and the Best Paper Award at ACM FAccT in 2019. He was elected an ACM Fellow in 2019. He obtained a bachelor's degree in computer science and engineering from IIT Bombay and a PhD in algorithms, combinatorics, and optimization from Georgia Institute of Technology.

Algorithms for Convex Optimization

NISHEETH K. VISHNOI
Yale University



CAMBRIDGE
UNIVERSITY PRESS

University Printing House, Cambridge CB2 8BS, United Kingdom
One Liberty Plaza, 20th Floor, New York, NY 10006, USA
477 Williamstown Road, Port Melbourne, VIC 3207, Australia
314–321, 3rd Floor, Plot 3, Splendor Forum, Jasola District Centre,
New Delhi – 110025, India
103 Penang Road, #05–06/07, Visioncrest Commercial, Singapore 238467

Cambridge University Press is part of the University of Cambridge.

It furthers the University's mission by disseminating knowledge in the pursuit of education, learning, and research at the highest international levels of excellence.

www.cambridge.org

Information on this title: www.cambridge.org/9781108482028

DOI: 10.1017/9781108699211

© Nisheeth K. Vishnoi 2021

This publication is in copyright. Subject to statutory exception and to the provisions of relevant collective licensing agreements, no reproduction of any part may take place without the written permission of Cambridge University Press.

First published 2021

A catalogue record for this publication is available from the British Library.

Library of Congress Cataloging-in-Publication Data

Names: Vishnoi, Nisheeth K., 1976– author.

Title: Algorithms for convex optimization / Nisheeth K. Vishnoi.

Description: New York : Cambridge University Press, [2021] |

Includes bibliographical references and index.

Identifiers: LCCN 2020052071 (print) | LCCN 2020052072 (ebook) |

ISBN 9781108482028 (hardback) | ISBN 9781108741774 (paperback) |

ISBN 9781108699211 (epub)

Subjects: LCSH: Mathematical optimization. | Convex functions. | Convex programming.

Classification: LCC QA402.5 .V57 2021 (print) | LCC QA402.5 (ebook) |

DDC 515/.882–dc23

LC record available at <https://lcn.loc.gov/2020052071>

LC ebook record available at <https://lcn.loc.gov/2020052072>

ISBN 978-1-108-48202-8 Hardback

ISBN 978-1-108-74177-4 Paperback

Cambridge University Press has no responsibility for the persistence or accuracy of URLs for external or third-party internet websites referred to in this publication and does not guarantee that any content on such websites is, or will remain, accurate or appropriate.

Dedicated to Maya and Vayu

Contents

<i>Preface</i>	<i>page xi</i>
<i>Acknowledgments</i>	xiv
<i>Notation</i>	xv
1 Bridging Continuous and Discrete Optimization	1
1.1 An Example: The Maximum Flow Problem	2
1.2 Linear Programming	8
1.3 Fast and Exact Algorithms via Interior Point Methods	12
1.4 Ellipsoid Method beyond Succinct Linear Programs	13
2 Preliminaries	17
2.1 Derivatives, Gradients, and Hessians	17
2.2 Fundamental Theorem of Calculus	19
2.3 Taylor Approximation	19
2.4 Linear Algebra, Matrices, and Eigenvalues	20
2.5 The Cauchy-Schwarz Inequality	23
2.6 Norms	24
2.7 Euclidean Topology	25
2.8 Dynamical Systems	25
2.9 Graphs	27
2.10 Exercises	29
3 Convexity	35
3.1 Convex Sets	35
3.2 Convex Functions	36
3.3 The Usefulness of Convexity	42
3.4 Exercises	46
4 Convex Optimization and Efficiency	49
4.1 Convex Programs	49

4.2	Computational Models	51
4.3	Membership Problem for Convex Sets	53
4.4	Solution Concepts for Optimization Problems	59
4.5	The Notion of Polynomial Time for Convex Optimization	63
4.6	Exercises	65
5	Duality and Optimality	69
5.1	Lagrangian Duality	70
5.2	The Conjugate Function	74
5.3	KKT Optimality Conditions	76
5.4	Proof of Strong Duality under Slater's Condition	78
5.5	Exercises	79
6	Gradient Descent	84
6.1	The Setup	84
6.2	Gradient Descent	85
6.3	Analysis When the Gradient Is Lipschitz Continuous	89
6.4	Application: The Maximum Flow Problem	96
6.5	Exercises	101
7	Mirror Descent and the Multiplicative Weights Update	108
7.1	Beyond the Lipschitz Gradient Condition	108
7.2	A Local Optimization Principle and Regularizers	110
7.3	Exponential Gradient Descent	112
7.4	Mirror Descent	121
7.5	Multiplicative Weights Update	125
7.6	Application: Perfect Matching in Bipartite Graphs	126
7.7	Exercises	132
8	Accelerated Gradient Descent	143
8.1	The Setup	143
8.2	Main Result on Accelerated Gradient Descent	144
8.3	Proof Strategy: Estimate Sequences	145
8.4	Construction of an Estimate Sequence	147
8.5	The Algorithm and Its Analysis	152
8.6	An Algorithm for Strongly Convex and Smooth Functions	153
8.7	Application: Linear System of Equations	155
8.8	Exercises	156
9	Newton's Method	160
9.1	Finding a Root of a Univariate Function	160
9.2	Newton's Method for Multivariate Functions	164
9.3	Newton's Method for Unconstrained Optimization	165

9.4	First Take on the Analysis	167
9.5	Newton's Method as Steepest Descent	170
9.6	Analysis Based on a Local Norm	175
9.7	Analysis Based on the Euclidean Norm	181
9.8	Exercises	182
10	An Interior Point Method for Linear Programming	185
10.1	Linear Programming	185
10.2	Constrained Optimization via Barrier Functions	187
10.3	The Logarithmic Barrier Function	188
10.4	The Central Path	189
10.5	A Path-Following Algorithm for Linear Programming	190
10.6	Analysis of the Path-Following Algorithm	194
10.7	Exercises	208
11	Variants of Interior Point Method and Self-Concordance	215
11.1	The Minimum Cost Flow Problem	215
11.2	An IPM for Linear Programming in Standard Form	219
11.3	Application: The Minimum Cost Flow Problem	226
11.4	Self-Concordant Barriers	230
11.5	Linear Programming Using Self-Concordant Barriers	232
11.6	Semidefinite Programming Using Self-Concordant Barriers	239
11.7	Convex Optimization Using Self-Concordant Barriers	242
11.8	Exercises	242
12	Ellipsoid Method for Linear Programming	248
12.1	0-1-Polytopes with Exponentially Many Constraints	248
12.2	Cutting Plane Methods	252
12.3	Ellipsoid Method	258
12.4	Analysis of Volume Drop and Efficiency for Ellipsoids	261
12.5	Application: Linear Optimization over 0-1-Polytopes	269
12.6	Exercises	273
13	Ellipsoid Method for Convex Optimization	279
13.1	Convex Optimization Using the Ellipsoid Method?	279
13.2	Application: Submodular Function Minimization	281
13.3	Application: The Maximum Entropy Problem	289
13.4	Convex Optimization Using the Ellipsoid Method	294
13.5	Variants of Cutting Plane Method	301
13.6	Exercises	304
	<i>Bibliography</i>	310
	<i>Index</i>	319

Preface

Convex optimization studies the problem of minimizing a convex function over a convex set. Convexity, along with its numerous implications, has been used to come up with efficient algorithms for many classes of convex programs. Consequently, convex optimization has broadly impacted several disciplines of science and engineering.

In the last few years, algorithms for convex optimization have revolutionized algorithm design, both for discrete and continuous optimization problems. The fastest-known algorithms for problems such as maximum flow in graphs, maximum matching in bipartite graphs, and submodular function minimization involve an essential and nontrivial use of algorithms for convex optimization such as gradient descent, mirror descent, interior point methods, and cutting plane methods. Surprisingly, algorithms for convex optimization have also been used to design counting problems over discrete objects such as matroids. Simultaneously, algorithms for convex optimization have become central to many modern machine learning applications. The demand for algorithms for convex optimization, driven by larger and increasingly complex input instances, has also significantly pushed the state of the art of convex optimization itself.

The goal of this book is to enable a reader to gain an in-depth understanding of algorithms for convex optimization. The emphasis is to derive key algorithms for convex optimization from first principles and to establish precise running time bounds in terms of the input length. Given the broad applicability of these methods, it is not possible for a single book to show the applications of these methods to all of them. This book shows applications to fast algorithms for various discrete optimization and counting problems. The applications selected in this book serve the purpose of illustrating a rather surprising bridge between continuous and discrete optimization.

The structure of the book. The book has roughly four parts. Chapters 3, 4, and 5 provide an introduction to convexity, models of computation and notions of efficiency in convex optimization, and duality. Chapters 6, 7, and 8 introduce first-order methods such as gradient descent, mirror descent and the multiplicative weights update method, and accelerated gradient descent, respectively. Chapters 9, 10, and 11 present Newton’s method and various interior point methods for linear programming. Chapters 12 and 13 present cutting plane methods such as the ellipsoid method for linear and general convex programs. Chapter 1 summarizes the book via a brief history of the interplay between continuous and discrete optimization: how the search for fast algorithms for discrete problems is leading to improvements in algorithms for convex optimization.

Many chapters contain applications ranging from finding maximum flows, minimum cuts, and perfect matchings in graphs, to linear optimization over 0-1-polytopes, to submodular function minimization, to computing maximum entropy distributions over combinatorial polytopes.

The book is self-contained and starts with a review of calculus, linear algebra, geometry, dynamical systems, and graph theory in Chapter 2. Exercises posed in this book not only play an important role in checking one’s understanding; sometimes important methods and concepts are introduced and developed entirely through them. Examples include the Frank-Wolfe method, coordinate descent, stochastic gradient descent, online convex optimization, the min-max theorem for zero-sum games, the Winnow algorithm for classification, bandit optimization, the conjugate gradient method, primal-dual interior point method, and matrix scaling.

How to use this book. This book can be used either as a textbook for a stand-alone advanced undergraduate or beginning graduate-level course, or as a supplement to an introductory course on convex optimization or algorithm design. The intended audience includes advanced undergraduate students, graduate students, and researchers from theoretical computer science, discrete optimization, operations research, statistics, and machine learning. To make this book accessible to a broad audience with different backgrounds, the writing style deliberately emphasizes the intuition, sometimes at the expense of rigor.

A course for a theoretical computer science or discrete optimization audience could cover the entire book. A course on convex optimization can omit the applications to discrete optimization and can, instead, include applications as per the choice of the instructor. Finally, an introductory course on convex optimization for machine learning could include material from Chapters 2 to 7.

Beyond convex optimization? This book should also prepare the reader for working in areas beyond convex optimization, e.g., nonconvex optimization and geodesic convex optimization, which are currently in their formative years.

Nonconvex optimization. One property of convex functions is that a “local” minimum is also a “global” minimum. Thus, algorithms for convex optimization, essentially, find a local minimum. Interestingly, this viewpoint has led to convex optimization methods being very successful for nonconvex optimization problems, especially those that arise in machine learning. Unlike convex programs, some of which can be **NP**-hard to optimize, most interesting classes of nonconvex optimization problems are **NP**-hard. Hence, in many of these applications, we define a suitable notion of local minimum and look for methods that can take us to one. Thus, algorithms for convex optimization are important for nonconvex optimization as well; see the survey by Jain and Kar (2017).

Geodesic convex optimization. Sometimes, a function that is nonconvex in a Euclidean space turns out to be convex if we introduce a suitable Riemannian metric on the underlying space and redefine convexity with respect to the “straight lines” – geodesics – induced by the metric. Such functions are called geodesically convex and arise in optimization problems over Riemannian manifolds such as matrix Lie groups; see the survey by Vishnoi (2018). The theory of efficient algorithms for geodesic convex optimization is under construction, and the paper by Bürgisser et al. (2019) presents some recent progress.

Acknowledgments

The contents of this book have been developed over several courses – for both undergraduate and graduate students – that I have taught, starting in Fall 2014 and is closest to that of a course taught in Fall 2019 at Yale. I am grateful to all the students and other attendees of these courses for their questions and comments that have made me reflect on the topic and improve the presentation. I am thankful to Slobodan Mitrovic, Damian Straszak, Jakub Tarnawski, and George Zakhour for being some of the first to take this course and scribing my initial lectures on this topic. Special thanks to Damian for scribing a significant fraction of my lectures, sometimes adding his own insights. I am indebted to Somenath Biswas, Elisa Celis, Yan Zhong Ding, and Anay Mehrotra for carefully reading a draft of this book and giving numerous valuable comments and suggestions.

Finally, this book has been influenced by several classic works: *Geometric Algorithms and Combinatorial Optimization* by Grötschel et al. (1988), *Convex Optimization* by Boyd and Vandenberghe (2004), *Introductory Lectures on Convex Optimization* by Nesterov (2014), and *The Multiplicative Weights Update Method: A Meta-algorithm and Applications* by Arora et al. (2012).

Notation

Numbers and sets:

- The set of natural numbers, integers, rationals, and real numbers are denoted by \mathbb{N} , \mathbb{Z} , \mathbb{Q} , and \mathbb{R} , respectively. $\mathbb{Z}_{\geq 0}$, $\mathbb{Q}_{\geq 0}$, and $\mathbb{R}_{\geq 0}$ denote the set of nonnegative integers, rationals, and reals, respectively.
- For a positive integer n , we denote by $[n]$ the set $\{1, 2, \dots, n\}$.
- For a set $S \subseteq [n]$, we use $1_S \in \mathbb{R}^n$ to denote the indicator vector of S defined as $1_S(i) = 1$ for all $i \in S$ and $1_S(i) = 0$ otherwise.
- For a set $S \subseteq [n]$ of cardinality k , we sometimes write \mathbb{R}^S to denote \mathbb{R}^k .

Vectors, matrices, inner products, and norms:

- Vectors are denoted by x and y . A vector $x \in \mathbb{R}^n$ is a column vector but is usually written as $x = (x_1, \dots, x_n)$. The transpose of a vector x is denoted by x^\top .
- The standard basis vectors in \mathbb{R}^n are denoted by e_1, \dots, e_n , where e_i is the vector whose i th entry is one and the remaining entries are zero.
- For vectors $x, y \in \mathbb{R}^n$, by $x \geq y$, we mean that $x_i \geq y_i$ for all $i \in [n]$.
- For a vector $x \in \mathbb{R}^n$, we use $\text{Diag}(x)$ to denote the $n \times n$ matrix whose (i, i) th entry is x_i for $1 \leq i \leq n$ and is zero on all other entries.
- When it is clear from context, 0 and 1 are also used to denote vectors with all 0 entries and all 1 entries, respectively.
- For vectors x and y , their inner product is denoted by $\langle x, y \rangle$ or $x^\top y$.
- For a vector x , its ℓ_2 or Euclidean norm is denoted by $\|x\|_2 := \sqrt{\langle x, x \rangle}$. We sometimes also refer to the ℓ_1 or Manhattan distance norm $\|x\|_1 := \sum_{i=1}^n |x_i|$. The ℓ_∞ -norm is defined as $\|x\|_\infty := \max_{i=1}^n |x_i|$.
- The outer product of a vector x with itself is denoted by xx^\top .
- Matrices are denoted by capitals, e.g., A and L . The transpose of A is denoted by A^\top .

- The trace of an $n \times n$ matrix A is $\text{Tr}(A) := \sum_{i=1}^n A_{ii}$. The determinant of an $n \times n$ matrix A is $\det(A) = \sum_{\sigma \in S_n} \text{sgn}(\sigma) \prod_{i=1}^n A_{i\sigma(i)}$. Here S_n is the set of all permutations of n elements and $\text{sgn}(\sigma)$ is the number of transpositions in a permutation σ , i.e., the number of pairs $i < j$ such that $\sigma(i) > \sigma(j)$.

Graphs:

- A graph G has a vertex set V and an edge set E . All graphs are assumed to be undirected unless stated otherwise. If the graph is weighted, there is a weight function $w : E \rightarrow \mathbb{R}_{\geq 0}$.
- A graph is said to be simple if there is at most one edge between two vertices and there are no edges whose endpoints are the same vertex.
- Typically, n is reserved for the number of vertices $|V|$ and m for the number of edges $|E|$.

Probability:

- $\mathbb{E}_{\mathcal{D}}[\cdot]$ denotes the expectation and $\Pr_{\mathcal{D}}[\cdot]$ denotes the probability over a distribution \mathcal{D} . The subscript is dropped when clear from context.

Running times:

- Standard big-O notation is used to describe the limiting behavior of a function. \tilde{O} denotes that potential poly-logarithmic factors have been omitted, i.e., $f = \tilde{O}(g)$ is equivalent to $f = O(g \log^k(g))$ for some constant k .