

# Fundamentals of Optimization Theory With Applications to Machine Learning

Jean Gallier and Jocelyn Quaintance  
Department of Computer and Information Science  
University of Pennsylvania  
Philadelphia, PA 19104, USA  
e-mail: [jean@cis.upenn.edu](mailto:jean@cis.upenn.edu)

© Jean Gallier

December 16, 2019



# Preface

In recent years, computer vision, robotics, machine learning, and data science have been some of the key areas that have contributed to major advances in technology. Anyone who looks at papers or books in the above areas will be baffled by a strange jargon involving exotic terms such as kernel PCA, ridge regression, lasso regression, support vector machines (SVM), Lagrange multipliers, KKT conditions, *etc.* Do support vector machines chase cattle to catch them with some kind of super lasso? No! But one will quickly discover that behind the jargon which always comes with a new field (perhaps to keep the outsiders out of the club), lies a lot of “classical” linear algebra and techniques from optimization theory. And there comes the main challenge: in order to understand and use tools from machine learning, computer vision, and so on, one needs to have a firm background in linear algebra and optimization theory. To be honest, some probability theory and statistics should also be included, but we already have enough to contend with.

Many books on machine learning struggle with the above problem. How can one understand what are the dual variables of a ridge regression problem if one doesn’t know about the Lagrangian duality framework? Similarly, how is it possible to discuss the dual formulation of SVM without a firm understanding of the Lagrangian framework?

The easy way out is to sweep these difficulties under the rug. If one is just a consumer of the techniques we mentioned above, the cookbook recipe approach is probably adequate. But this approach doesn’t work for someone who really wants to do serious research and make significant contributions. To do so, we believe that one must have a solid background in linear algebra and optimization theory.

This is a problem because it means investing a great deal of time and energy studying these fields, but we believe that perseverance will be amply rewarded.

This second volume covers some elements of optimization theory and applications, especially to machine learning. This volume is divided in five parts:

- (1) Preliminaries of Optimization Theory.
- (2) Linear Optimization.
- (3) Nonlinear Optimization.
- (4) Applications to Machine Learning.

- (5) An appendix consisting of two chapters; one on Hilbert bases and the Riesz–Fischer theorem, the other one containing `Matlab` code.

Part I is devoted to some preliminaries of optimization theory. The goal of most optimization problems is to minimize (or maximize) some objective function  $J$  subject to equality or inequality constraints. Therefore it is important to understand when a function  $J$  has a minimum or a maximum (an optimum). In most optimization problems, we need to find necessary conditions for a function  $J: \Omega \rightarrow \mathbb{R}$  to have a local extremum with respect to a subset  $U$  of  $\Omega$  (where  $\Omega$  is open). This can be done in two cases:

- (1) The set  $U$  is defined by a set of equations,

$$U = \{x \in \Omega \mid \varphi_i(x) = 0, \ 1 \leq i \leq m\},$$

where the functions  $\varphi_i: \Omega \rightarrow \mathbb{R}$  are continuous (and usually differentiable).

- (2) The set  $U$  is defined by a set of inequalities,

$$U = \{x \in \Omega \mid \varphi_i(x) \leq 0, \ 1 \leq i \leq m\},$$

where the functions  $\varphi_i: \Omega \rightarrow \mathbb{R}$  are continuous (and usually differentiable).

The case of equality constraints is much easier to deal with and is treated in Chapter 4.

In the case of equality constraints, a necessary condition for a local extremum with respect to  $U$  can be given in terms of *Lagrange multipliers*.

Part II deals with the special case where the objective function is a linear form and the constraints are affine inequality and equality constraints. This subject is known as *linear programming*, and the next four chapters give an introduction to the subject.

Part III is devoted to nonlinear optimization, which is the case where the objective function  $J$  is not linear and the constraints are inequality constraints. Since it is practically impossible to say anything interesting if the constraints are not convex, we quickly consider the convex case.

Chapter 13 is devoted to some general results of optimization theory. A main theme is to find sufficient conditions that ensure that an objective function has a minimum which is achieved. We define gradient descent methods (including Newton’s method), and discuss their convergence.

Chapter 14 contains the most important results of nonlinear optimization theory. Theorem 14.6 gives necessary conditions for a function  $J$  to have a minimum on a subset  $U$  defined by convex inequality constraints in terms of the Karush–Kuhn–Tucker conditions. Furthermore, if  $J$  is also convex and if the KKT conditions hold, then  $J$  has a global minimum.

We illustrate the KKT conditions on an interesting example from machine learning the so-called *hard margin support vector machine*; see Sections 14.5 and 14.6. The problem is to separate two disjoint sets of points,  $\{u_i\}_{i=1}^p$  and  $\{v_j\}_{j=1}^q$ , using a hyperplane satisfying some optimality property (to maximize the margin).

Section 14.7 contains the most important results of the chapter. The notion of Lagrangian duality is presented and we discuss *weak duality* and *strong duality*.

In Chapter 15, we consider some deeper aspects of the theory of convex functions that are not necessarily differentiable at every point of their domain. Some substitute for the gradient is needed. Fortunately, for convex functions, there is such a notion, namely *subgradients*. A major motivation for developing this more sophisticated theory of differentiation of convex functions is to extend the Lagrangian framework to convex functions that are not necessarily differentiable.

Chapter 16 is devoted to the presentation of one of the best methods known at the present for solving optimization problems involving equality constraints, called ADMM (alternating direction method of multipliers). In fact, this method can also handle more general constraints, namely, membership in a convex set. It can also be used to solve *lasso minimization*.

In Section 16.4, we prove the convergence of ADMM under exactly the same assumptions as in Boyd et al. [17]. It turns out that Assumption (2) in Boyd et al. [17] implies that the matrices  $A^\top A$  and  $B^\top B$  are invertible (as we show after the proof of Theorem 16.1). This allows us to prove a convergence result stronger than the convergence result proven in Boyd et al. [17].

The next four chapters constitute Part IV, which covers some applications of optimization theory (in particular Lagrangian duality) to machine learning.

Chapter 17 is an introduction to positive definite kernels and the use of kernel functions in machine learning called a *kernel function*.

We illustrate the kernel methods on kernel PCA.

In Chapter 18 we return to the problem of separating two disjoint sets of points,  $\{u_i\}_{i=1}^p$  and  $\{v_j\}_{j=1}^q$ , but this time we do not assume that these two sets are separable. To cope with nonseparability, we allow points to invade the safety zone around the separating hyperplane, and even points on the wrong side of the hyperplane. Such a method is called *soft margin support vector machine (SVM)*. We discuss variations of this method, including  $\nu$ -SV classification. In each case we present a careful derivation of the dual. We prove rigorous results about the existence of support vectors.

In Chapter 19, we discuss *linear regression*, *ridge regression*, *lasso regression* and *elastic net regression*.

In Chapter 20 we present  $\nu$ -SV *Regression*. This method is designed in the same spirit as soft margin SVM, in the sense that it allows a margin of error. Here the errors are penalized

in the  $\ell^1$ -sense. We present a careful derivation of the dual and discuss the existence of support vectors.

The methods presented in Chapters 18, 19 and 20 have all been implemented in `Matlab`, and much of this code is given in Appendix B. Remarkably, ADMM emerges as the main engine for solving most of these optimization problems. Thus it is nice to see the continuum spanning from theoretical considerations of convergence and correctness to practical matters of implementation. It is fun to see how these abstract Lagrange multipliers yield concrete results such as the weight vector  $w$  defining the desired hyperplane in regression or SVM.

Except for a few exceptions we provide complete proofs. We did so to make this book self-contained, but also because we believe that no deep knowledge of this material can be acquired without working out some proofs. However, our advice is to skip some of the proofs upon first reading, especially if they are long and intricate.

The chapters or sections marked with the symbol  $\otimes$  contain material that is typically more specialized or more advanced, and they can be omitted upon first (or second) reading.

*Acknowledgement:* We would like to thank Christine Allen-Blanchette, Kostas Daniilidis, Carlos Esteves, Spyridon Leonardos, Stephen Phillips, João Sedoc, Stephen Shatz, Jianbo Shi, and Marcelo Siqueira, for reporting typos and for helpful comments. Thanks to Gilbert Strang. We learned much from his books which have been a major source of inspiration. Special thanks to Steven Boyd. We learned a lot from his remarkable book on convex optimization and his papers, and Part III of our book is significantly inspired by his writings. The first author also wishes to express his deepest gratitude to Philippe G. Ciarlet who was his teacher and mentor in 1970-1972 while he was a student at ENPC in Paris. Professor Ciarlet was by far his best teacher. He also knew how to instill in his students the importance of intellectual rigor, honesty, and modesty. He still has his typewritten notes on measure theory and integration, and on numerical linear algebra. The latter became his wonderful book Ciarlet [25], from which we have borrowed heavily.

# Contents

<b>1</b>	<b>Introduction</b>	<b>13</b>
<b>I</b>	<b>Preliminaries for Optimization Theory</b>	<b>23</b>
<b>2</b>	<b>Topology</b>	<b>25</b>
2.1	Metric Spaces and Normed Vector Spaces . . . . .	25
2.2	Topological Spaces . . . . .	31
2.3	Subspace and Product Topologies . . . . .	36
2.4	Continuous Functions . . . . .	41
2.5	Limits and Continuity; Uniform Continuity . . . . .	45
2.6	Continuous Linear and Multilinear Maps . . . . .	50
2.7	Complete Metric Spaces and Banach Spaces . . . . .	56
2.8	Completion of a Metric Space . . . . .	57
2.9	Completion of a Normed Vector Space . . . . .	65
2.10	The Contraction Mapping Theorem . . . . .	66
2.11	Further Readings . . . . .	67
2.12	Summary . . . . .	67
2.13	Problems . . . . .	68
<b>3</b>	<b>Differential Calculus</b>	<b>71</b>
3.1	Directional Derivatives, Total Derivatives . . . . .	71
3.2	Properties of Derivatives . . . . .	79
3.3	Jacobian Matrices . . . . .	85
3.4	The Implicit and The Inverse Function Theorems . . . . .	92
3.5	Second-Order and Higher-Order Derivatives . . . . .	99
3.6	Taylor's Formula, Faà di Bruno's Formula . . . . .	106
3.7	Further Readings . . . . .	111
3.8	Summary . . . . .	111
3.9	Problems . . . . .	112
<b>4</b>	<b>Extrema of Real-Valued Functions</b>	<b>115</b>
4.1	Local Extrema and Lagrange Multipliers . . . . .	116
4.2	Using Second Derivatives to Find Extrema . . . . .	127

4.3	Using Convexity to Find Extrema . . . . .	131
4.4	Summary . . . . .	141
4.5	Problems . . . . .	142
<b>5</b>	<b>Newton's Method and Its Generalizations</b>	<b>145</b>
5.1	Newton's Method for Real Functions of a Real Argument . . . . .	145
5.2	Generalizations of Newton's Method . . . . .	147
5.3	Summary . . . . .	156
5.4	Problems . . . . .	156
<b>6</b>	<b>Quadratic Optimization Problems</b>	<b>165</b>
6.1	Quadratic Optimization: The Positive Definite Case . . . . .	165
6.2	Quadratic Optimization: The General Case . . . . .	175
6.3	Maximizing a Quadratic Function on the Unit Sphere . . . . .	180
6.4	Summary . . . . .	185
6.5	Problems . . . . .	186
<b>7</b>	<b>Schur Complements and Applications</b>	<b>187</b>
7.1	Schur Complements . . . . .	187
7.2	SPD Matrices and Schur Complements . . . . .	190
7.3	SP Semidefinite Matrices and Schur Complements . . . . .	191
7.4	Summary . . . . .	193
7.5	Problems . . . . .	193
<b>II</b>	<b>Linear Optimization</b>	<b>195</b>
<b>8</b>	<b>Convex Sets, Cones, <math>\mathcal{H}</math>-Polyhedra</b>	<b>197</b>
8.1	What is Linear Programming? . . . . .	197
8.2	Affine Subsets, Convex Sets, Hyperplanes, Half-Spaces . . . . .	199
8.3	Cones, Polyhedral Cones, and $\mathcal{H}$ -Polyhedra . . . . .	202
8.4	Summary . . . . .	207
8.5	Problems . . . . .	208
<b>9</b>	<b>Linear Programs</b>	<b>209</b>
9.1	Linear Programs, Feasible Solutions, Optimal Solutions . . . . .	209
9.2	Basic Feasible Solutions and Vertices . . . . .	216
9.3	Summary . . . . .	223
9.4	Problems . . . . .	223
<b>10</b>	<b>The Simplex Algorithm</b>	<b>227</b>
10.1	The Idea Behind the Simplex Algorithm . . . . .	227
10.2	The Simplex Algorithm in General . . . . .	236



10.3	How to Perform a Pivoting Step Efficiently . . . . .	243
10.4	The Simplex Algorithm Using Tableaux . . . . .	247
10.5	Computational Efficiency of the Simplex Method . . . . .	255
10.6	Summary . . . . .	257
10.7	Problems . . . . .	258
<b>11</b>	<b>Linear Programming and Duality</b>	<b>261</b>
11.1	Variants of the Farkas Lemma . . . . .	261
11.2	The Duality Theorem in Linear Programming . . . . .	267
11.3	Complementary Slackness Conditions . . . . .	275
11.4	Duality for Linear Programs in Standard Form . . . . .	276
11.5	The Dual Simplex Algorithm . . . . .	279
11.6	The Primal-Dual Algorithm . . . . .	286
11.7	Summary . . . . .	296
11.8	Problems . . . . .	296
<b>III</b>	<b>NonLinear Optimization</b>	<b>301</b>
<b>12</b>	<b>Basics of Hilbert Spaces</b>	<b>303</b>
12.1	The Projection Lemma . . . . .	303
12.2	Duality and the Riesz Representation Theorem . . . . .	316
12.3	Farkas–Minkowski Lemma in Hilbert Spaces . . . . .	321
12.4	Summary . . . . .	322
12.5	Problems . . . . .	323
<b>13</b>	<b>General Results of Optimization Theory</b>	<b>325</b>
13.1	Optimization Problems; Basic Terminology . . . . .	325
13.2	Existence of Solutions of an Optimization Problem . . . . .	329
13.3	Minima of Quadratic Functionals . . . . .	334
13.4	Elliptic Functionals . . . . .	340
13.5	Iterative Methods for Unconstrained Problems . . . . .	343
13.6	Gradient Descent Methods for Unconstrained Problems . . . . .	346
13.7	Convergence of Gradient Descent with Variable Stepsize . . . . .	354
13.8	Steepest Descent for an Arbitrary Norm . . . . .	357
13.9	Newton’s Method For Finding a Minimum . . . . .	359
13.10	Conjugate Gradient Methods; Unconstrained Problems . . . . .	363
13.11	Gradient Projection for Constrained Optimization . . . . .	375
13.12	Penalty Methods for Constrained Optimization . . . . .	377
13.13	Summary . . . . .	379
13.14	Problems . . . . .	381
<b>14</b>	<b>Introduction to Nonlinear Optimization</b>	<b>385</b>

14.1	The Cone of Feasible Directions . . . . .	387
14.2	Active Constraints and Qualified Constraints . . . . .	393
14.3	The Karush–Kuhn–Tucker Conditions . . . . .	400
14.4	Equality Constrained Minimization . . . . .	411
14.5	Hard Margin Support Vector Machine; Version I . . . . .	416
14.6	Hard Margin Support Vector Machine; Version II . . . . .	421
14.7	Lagrangian Duality and Saddle Points . . . . .	429
14.8	Weak and Strong Duality . . . . .	438
14.9	Handling Equality Constraints Explicitly . . . . .	446
14.10	Dual of the Hard Margin Support Vector Machine . . . . .	450
14.11	Conjugate Function and Legendre Dual Function . . . . .	455
14.12	Some Techniques to Obtain a More Useful Dual Program . . . . .	465
14.13	Uzawa’s Method . . . . .	469
14.14	Summary . . . . .	474
14.15	Problems . . . . .	476
<b>15</b>	<b>Subgradients and Subdifferentials</b> $\otimes$	<b>479</b>
15.1	Extended Real-Valued Convex Functions . . . . .	481
15.2	Subgradients and Subdifferentials . . . . .	490
15.3	Basic Properties of Subgradients and Subdifferentials . . . . .	502
15.4	Additional Properties of Subdifferentials . . . . .	509
15.5	The Minimum of a Proper Convex Function . . . . .	513
15.6	Generalization of the Lagrangian Framework . . . . .	519
15.7	Summary . . . . .	523
15.8	Problems . . . . .	524
<b>16</b>	<b>Dual Ascent Methods; ADMM</b>	<b>527</b>
16.1	Dual Ascent . . . . .	529
16.2	Augmented Lagrangians and the Method of Multipliers . . . . .	533
16.3	ADMM: Alternating Direction Method of Multipliers . . . . .	538
16.4	Convergence of ADMM $\otimes$ . . . . .	541
16.5	Stopping Criteria . . . . .	550
16.6	Some Applications of ADMM . . . . .	552
16.7	Solving Hard Margin ( $SVM_{h_2}$ ) Using ADMM . . . . .	556
16.8	Applications of ADMM to $\ell^1$ -Norm Problems . . . . .	557
16.9	Summary . . . . .	563
16.10	Problems . . . . .	564
<b>IV</b>	<b>Applications to Machine Learning</b>	<b>565</b>
<b>17</b>	<b>Positive Definite Kernels</b>	<b>567</b>
17.1	Feature Maps and Kernel Functions . . . . .	567

17.2	Basic Properties of Positive Definite Kernels . . . . .	573
17.3	Hilbert Space Representation of a Positive Kernel . . . . .	580
17.4	Kernel PCA . . . . .	583
17.5	Summary . . . . .	586
17.6	Problems . . . . .	587
<b>18</b>	<b>Soft Margin Support Vector Machines</b>	<b>589</b>
18.1	Soft Margin Support Vector Machines; $(SVM_{s1})$ . . . . .	592
18.2	Solving SVM $(SVM_{s1})$ Using ADMM . . . . .	607
18.3	Soft Margin Support Vector Machines; $(SVM_{s2})$ . . . . .	608
18.4	Solving SVM $(SVM_{s2})$ Using ADMM . . . . .	615
18.5	Soft Margin Support Vector Machines; $(SVM_{s2'})$ . . . . .	616
18.6	Classification of the Data Points in Terms of $\nu$ $(SVM_{s2'})$ . . . . .	626
18.7	Existence of Support Vectors for $(SVM_{s2'})$ . . . . .	629
18.8	Solving SVM $(SVM_{s2'})$ Using ADMM . . . . .	640
18.9	Soft Margin Support Vector Machines; $(SVM_{s3})$ . . . . .	644
18.10	Classification of the Data Points in Terms of $\nu$ $(SVM_{s3})$ . . . . .	651
18.11	Existence of Support Vectors for $(SVM_{s3})$ . . . . .	653
18.12	Solving SVM $(SVM_{s3})$ Using ADMM . . . . .	655
18.13	Soft Margin SVM; $(SVM_{s4})$ . . . . .	658
18.14	Solving SVM $(SVM_{s4})$ Using ADMM . . . . .	667
18.15	Soft Margin SVM; $(SVM_{s5})$ . . . . .	669
18.16	Solving SVM $(SVM_{s5})$ Using ADMM . . . . .	673
18.17	Summary and Comparison of the SVM Methods . . . . .	675
18.18	Problems . . . . .	688
<b>19</b>	<b>Ridge Regression, Lasso, Elastic Net</b>	<b>693</b>
19.1	Ridge Regression . . . . .	694
19.2	Ridge Regression; Learning an Affine Function . . . . .	697
19.3	Kernel Ridge Regression . . . . .	706
19.4	Lasso Regression ( $\ell^1$ -Regularized Regression) . . . . .	710
19.5	Lasso Regression; Learning an Affine Function . . . . .	714
19.6	Elastic Net Regression . . . . .	720
19.7	Summary . . . . .	726
19.8	Problems . . . . .	726
<b>20</b>	<b><math>\nu</math>-SV Regression</b>	<b>729</b>
20.1	$\nu$ -SV Regression; Derivation of the Dual . . . . .	729
20.2	Existence of Support Vectors . . . . .	740
20.3	Solving $\nu$ -Regression Using ADMM . . . . .	750
20.4	Kernel $\nu$ -SV Regression . . . . .	756
20.5	$\nu$ -Regression Version 2; Penalizing $b$ . . . . .	759
20.6	Summary . . . . .	766

20.7 Problems . . . . .	767
<b>V Appendix</b>	<b>769</b>
<b>A Total Orthogonal Families in Hilbert Spaces</b>	<b>771</b>
A.1 Total Orthogonal Families, Fourier Coefficients . . . . .	771
A.2 The Hilbert Space $\ell^2(K)$ and the Riesz–Fischer Theorem . . . . .	780
A.3 Summary . . . . .	789
A.4 Problems . . . . .	790
<b>B Matlab Programs</b>	<b>791</b>
B.1 Hard Margin ( $\text{SVM}_{h2}$ ) . . . . .	791
B.2 Soft Margin SVM ( $\text{SVM}_{s2'}$ ) . . . . .	795
B.3 Soft Margin SVM ( $\text{SVM}_{s3}$ ) . . . . .	803
B.4 $\nu$ -SV Regression . . . . .	808
<b>Bibliography</b>	<b>813</b>

# Chapter 1

## Introduction

This second volume covers some elements of optimization theory and applications, especially to machine learning. This volume is divided in five parts:

- (1) Preliminaries of Optimization Theory.
- (2) Linear Optimization.
- (3) Nonlinear Optimization.
- (4) Applications to Machine Learning.
- (5) An appendix consisting of two chapters; one on Hilbert bases and the Riesz–Fischer theorem, the other one containing `Matlab` code.

Part I is devoted to some preliminaries of optimization theory. The goal of most optimization problems is to minimize (or maximize) some objective function  $J$  subject to equality or inequality constraints. Therefore it is important to understand when a function  $J$  has a minimum or a maximum (an optimum). If the function  $J$  is sufficiently differentiable, then a necessary condition for a function to have an optimum typically involves the derivative of the function  $J$ , and if  $J$  is real-valued, its gradient  $\nabla J$ .

Thus it is desirable to review some basic notions of topology and calculus, in particular, to have a firm grasp of the notion of derivative of a function between normed vector spaces. Partial derivatives  $\partial f/\partial A$  of functions whose range and domain are spaces of matrices tend to be used casually, even though in most cases a correct definition is never provided. It is possible, and simple, to define rigorously derivatives, gradients, and directional derivatives of functions defined on matrices and to avoid these nonsensical partial derivatives.

Chapter 2 contains a review of basic topological notions used in analysis. We pay particular attention to complete metric spaces and complete normed vector spaces. In fact, we provide a detailed construction of the completion of a metric space (and of a normed vector space) using equivalence classes of Cauchy sequences. Chapter 3 is devoted to some notions

of differential calculus, in particular, directional derivatives, total derivatives, gradients, Hessians, and the inverse function theorem.

Chapter 4 deals with extrema of real-valued functions. In most optimization problems, we need to find necessary conditions for a function  $J: \Omega \rightarrow \mathbb{R}$  to have a local extremum with respect to a subset  $U$  of  $\Omega$  (where  $\Omega$  is open). This can be done in two cases:

- (1) The set  $U$  is defined by a set of equations,

$$U = \{x \in \Omega \mid \varphi_i(x) = 0, \ 1 \leq i \leq m\},$$

where the functions  $\varphi_i: \Omega \rightarrow \mathbb{R}$  are continuous (and usually differentiable).

- (2) The set  $U$  is defined by a set of inequalities,

$$U = \{x \in \Omega \mid \varphi_i(x) \leq 0, \ 1 \leq i \leq m\},$$

where the functions  $\varphi_i: \Omega \rightarrow \mathbb{R}$  are continuous (and usually differentiable).

In (1), the equations  $\varphi_i(x) = 0$  are called *equality constraints*, and in (2), the inequalities  $\varphi_i(x) \leq 0$  are called *inequality constraints*. The case of equality constraints is much easier to deal with and is treated in Chapter 4.

If the functions  $\varphi_i$  are convex and  $\Omega$  is convex, then  $U$  is convex. This is a very important case that we will discuss later. In particular, if the functions  $\varphi_i$  are affine, then the equality constraints can be written as  $Ax = b$ , and the inequality constraints as  $Ax \leq b$ , for some  $m \times n$  matrix  $A$  and some vector  $b \in \mathbb{R}^m$ . We will also discuss the case of affine constraints later.

In the case of equality constraints, a necessary condition for a local extremum with respect to  $U$  can be given in terms of *Lagrange multipliers*. In the case of inequality constraints, there is also a necessary condition for a local extremum with respect to  $U$  in terms of generalized Lagrange multipliers and the *Karush–Kuhn–Tucker* conditions. This will be discussed in Chapter 14.

In Chapter 5 we discuss Newton’s method and some of its generalizations (the Newton–Kantorovich theorem). These are methods to find the zeros of a function.

Chapter 6 covers the special case of determining when a quadratic function has a minimum, subject to affine equality constraints. A complete answer is provided in terms of the notion of symmetric positive semidefinite matrices.

The Schur complement is introduced in Chapter 7. We give a complete proof of a criterion for a matrix to be positive definite (or positive semidefinite) stated in Boyd and Vandenberghe [18] (Appendix B).

Part II deals with the special case where the objective function is a linear form and the constraints are affine inequality and equality constraints. This subject is known as linear

programming, and the next four chapters give an introduction to the subject. Although linear programming has been supplanted by convex programming and its variants, it is still a great workhorse. It is also a great warm up for the general treatment of Lagrangian duality. We pay particular attention to versions of Farkas' lemma, which is at the heart of duality in linear programming.

Part III is devoted to nonlinear optimization, which is the case where the objective function  $J$  is not linear and the constraints are inequality constraints. Since it is practically impossible to say anything interesting if the constraints are not convex, we quickly consider the convex case.

In optimization theory one often deals with function spaces of infinite dimension. Typically, these spaces either are Hilbert spaces or can be completed as Hilbert spaces. Thus it is important to have some minimum knowledge about Hilbert spaces, and we feel that this minimum knowledge includes the projection lemma, the fact that a closed subset has an orthogonal complement, the Riesz representation theorem, and a version of the Farkas–Minkowski lemma. Chapter 12 covers these topics. A more detailed introduction to Hilbert spaces is given in Appendix A.

Chapter 13 is devoted to some general results of optimization theory. A main theme is to find sufficient conditions that ensure that an objective function has a minimum which is achieved. We define the notion of a coercive function. The most general result is Theorem 13.2, which applies to a coercive convex function on a convex subset of a separable Hilbert space. In the special case of a coercive quadratic functional, we obtain the Lions–Stampacchia theorem (Theorem 13.6), and the Lax–Milgram theorem (Theorem 13.7). We define elliptic functionals, which generalize quadratic functions defined by symmetric positive definite matrices. We define gradient descent methods, and discuss their convergence. A gradient descent method looks for a descent direction and a stepsize parameter, which is obtained either using an exact line search or a backtracking line search. A popular technique to find the search direction is steepest descent. In addition to steepest descent for the Euclidean norm, we discuss steepest descent for an arbitrary norm. We also consider a special case of steepest descent, Newton's method. This method converges faster than the other gradient descent methods, but it is quite expensive since it requires computing and storing Hessians. We also present the method of conjugate gradients and prove its correctness. We briefly discuss the method of gradient projection and the penalty method in the case of constrained optima.

Chapter 14 contains the most important results of nonlinear optimization theory. We begin by defining the cone of feasible directions and then state a necessary condition for a function to have local minimum on a set  $U$  that is not necessarily convex in terms of the cone of feasible directions. The cone of feasible directions is not always convex, but it is if the constraints are inequality constraints. An inequality constraint  $\varphi(u) \leq 0$  is said to be *active* if  $\varphi(u) = 0$ . One can also define the notion of *qualified constraint*. Theorem 14.5 gives necessary conditions for a function  $J$  to have a minimum on a subset  $U$  defined by qualified inequality constraints in terms of the Karush–Kuhn–Tucker conditions (for short

KKT conditions), which involve nonnegative Lagrange multipliers. The proof relies on a version of the Farkas–Minkowski lemma. Some of the KKT conditions assert that  $\lambda_i \varphi_i(u) = 0$ , where  $\lambda_i \geq 0$  is the Lagrange multiplier associated with the constraint  $\varphi_i \leq 0$ . To some extent, this implies that active constraints are more important than inactive constraints, since if  $\varphi_i(u) < 0$  is an inactive constraint, then  $\lambda_i = 0$ . In general, the KKT conditions are useless unless the constraints are convex. In this case, there is a manageable notion of qualified constraint given by Slater’s conditions. Theorem 14.6 gives necessary conditions for a function  $J$  to have a minimum on a subset  $U$  defined by convex inequality constraints in terms of the Karush–Kuhn–Tucker conditions. Furthermore, if  $J$  is also convex and if the KKT conditions hold, then  $J$  has a global minimum.

In Section 14.4, we apply Theorem 14.6 to the special case where the constraints are equality constraints, which can be expressed as  $Ax = b$ . In the special case where the convex objective function  $J$  is a convex quadratic functional of the form

$$J(x) = \frac{1}{2}x^\top Px + q^\top x + r,$$

where  $P$  is a  $n \times n$  symmetric positive semidefinite matrix, the necessary and sufficient conditions for having a minimum are expressed by a linear system involving a matrix called the KKT matrix. We discuss conditions that guarantee that the KKT matrix is invertible, and how to solve the KKT system. We also briefly discuss variants of Newton’s method dealing with equality constraints.

We illustrate the KKT conditions on an interesting example, the so-called hard margin support vector machine; see Sections 14.5 and 14.6. The problem is a classification problem, or more accurately a separation problem. Suppose we have two nonempty disjoint finite sets of  $p$  blue points  $\{u_i\}_{i=1}^p$  and  $q$  red points  $\{v_j\}_{j=1}^q$  in  $\mathbb{R}^n$ . Our goal is to find a hyperplane  $H$  of equation  $w^\top x - b = 0$  (where  $w \in \mathbb{R}^n$  is a nonzero vector and  $b \in \mathbb{R}$ ), such that all the blue points  $u_i$  are in one of the two open half-spaces determined by  $H$ , and all the red points  $v_j$  are in the other open half-space determined by  $H$ .

If the two sets are indeed separable, then in general there are infinitely many hyperplanes separating them. Vapnik had the idea to find a hyperplane that maximizes the smallest distance between the points and the hyperplane. Such a hyperplane is indeed unique and is called a maximal hard margin hyperplane, or hard margin support vector machine. The support vectors are those for which the constraints are active.

Section 14.7 contains the most important results of the chapter. The notion of Lagrangian duality is presented. Given a primal optimization problem ( $P$ ) consisting in minimizing an objective function  $J(v)$  with respect to some inequality constraints  $\varphi_i(v) \leq 0$ ,  $i = 1, \dots, m$ , we define the *dual function*  $G(\mu)$  as the result of minimizing the Lagrangian

$$L(v, \mu) = J(v) + \sum_{i=1}^m \mu_i \varphi_i(v)$$



with respect to  $v$ , with  $\mu \in \mathbb{R}_+^m$ . The dual program (D) is then to maximize  $G(\mu)$  with respect to  $\mu \in \mathbb{R}_+^m$ . It turns out that  $G$  is a concave function, and the dual program is an unconstrained maximization. This is actually a misleading statement because  $G$  is generally a partial function, so maximizing  $G(\mu)$  is equivalent to a constrained maximization problem in which the constraints specify the domain of  $G$ , but in many cases, we obtain a dual program simpler than the primal program. If  $d^*$  is the optimal value of the dual program and if  $p^*$  is the optimal value of the primal program, we always have

$$d^* \leq p^*,$$

which is known as *weak duality*. Under certain conditions,  $d^* = p^*$ , that is, the duality gap is zero, in which case we say that *strong duality* holds. Also, under certain conditions, a solution of the dual yields a solution of the primal, and if the primal has an optimal solution, then the dual has an optimal solution, but beware that the converse is generally false (see Theorem 14.17). We also show how to deal with equality constraints, and discuss the use of conjugate functions to find the dual function. Our coverage of Lagrangian duality is quite thorough, but we do not discuss more general orderings such as the semidefinite ordering. For these topics which belong to convex optimization, the reader is referred to Boyd and Vandenberghe [18].

In Chapter 15, we consider some deeper aspects of the theory of convex functions that are not necessarily differentiable at every point of their domain. Some substitute for the gradient is needed. Fortunately, for convex functions, there is such a notion, namely *subgradients*. Geometrically, given a (proper) convex function  $f$ , the subgradients at  $x$  are vectors normal to supporting hyperplanes to the epigraph of the function at  $(x, f(x))$ . The *subdifferential*  $\partial f(x)$  to  $f$  at  $x$  is the set of all subgradients at  $x$ . A crucial property is that  $f$  is differentiable at  $x$  iff  $\partial f(x) = \{\nabla f_x\}$ , where  $\nabla f_x$  is the gradient of  $f$  at  $x$ . Another important property is that a (proper) convex function  $f$  attains its minimum at  $x$  iff  $0 \in \partial f(x)$ . A major motivation for developing this more sophisticated theory of “differentiation” of convex functions is to extend the Lagrangian framework to convex functions that are not necessarily differentiable.

Experience shows that the applicability of convex optimization is significantly increased by considering extended real-valued functions, namely functions  $f: S \rightarrow \mathbb{R} \cup \{-\infty, +\infty\}$ , where  $S$  is some subset of  $\mathbb{R}^n$  (usually convex). This is reminiscent of what happens in measure theory, where it is natural to consider functions that take the value  $+\infty$ .

In Section 15.1, we introduce extended real-valued functions, which are functions that may also take the values  $\pm\infty$ . In particular, we define proper convex functions, and the closure of a convex function. Subgradients and subdifferentials are defined in Section 15.2. We discuss some properties of subgradients in Section 15.3 and Section 15.4. In particular, we relate subgradients to one-sided directional derivatives. In Section 15.5, we discuss the problem of finding the minimum of a proper convex function and give some criteria in terms of subdifferentials. In Section 15.6, we sketch the generalization of the results presented in

Chapter 14 about the Lagrangian framework to programs allowing an objective function and inequality constraints which are convex but not necessarily differentiable.

This chapter relies heavily on Rockafellar [61]. We tried to distill the body of results needed to generalize the Lagrangian framework to convex but not necessarily differentiable functions. Some of the results in this chapter are also discussed in Bertsekas [9, 12, 10].

Chapter 16 is devoted to the presentation of one of the best methods known at the present for solving optimization problems involving equality constraints, called ADMM (alternating direction method of multipliers). In fact, this method can also handle more general constraints, namely, membership in a convex set. It can also be used to solve *lasso minimization*.

In this chapter, we consider the problem of minimizing a convex function  $J$  (not necessarily differentiable) under the equality constraints  $Ax = b$ . In Section 16.1, we discuss the dual ascent method. It is essentially gradient descent applied to the dual function  $G$ , but since  $G$  is maximized, gradient descent becomes gradient ascent.

In order to make the minimization step of the dual ascent method more robust, one can use the trick of adding the penalty term  $(\rho/2) \|Au - b\|_2^2$  to the Lagrangian. We obtain the *augmented Lagrangian*

$$L_\rho(u, \lambda) = J(u) + \lambda^\top (Au - b) + (\rho/2) \|Au - b\|_2^2,$$

with  $\lambda \in \mathbb{R}^m$ , and where  $\rho > 0$  is called the *penalty parameter*. We obtain the minimization Problem  $(P_\rho)$ ,

$$\begin{aligned} & \text{minimize} && J(u) + (\rho/2) \|Au - b\|_2^2 \\ & \text{subject to} && Au = b, \end{aligned}$$

which is equivalent to the original problem.

The benefit of adding the penalty term  $(\rho/2) \|Au - b\|_2^2$  is that by Proposition 15.37, Problem  $(P_\rho)$  has a unique optimal solution under mild conditions on  $A$ . Dual ascent applied to the dual of  $(P_\rho)$  is called the *method of multipliers* and is discussed in Section 16.2.

The new twist in ADMM is to split the function  $J$  into two independent parts, as  $J(x, z) = f(x) + g(z)$ , and to consider the Minimization Problem  $(P_{\text{admm}})$ ,

$$\begin{aligned} & \text{minimize} && f(x) + g(z) \\ & \text{subject to} && Ax + Bz = c, \end{aligned}$$

for some  $p \times n$  matrix  $A$ , some  $p \times m$  matrix  $B$ , and with  $x \in \mathbb{R}^n$ ,  $z \in \mathbb{R}^m$ , and  $c \in \mathbb{R}^p$ . We also assume that  $f$  and  $g$  are convex.

As in the method of multipliers, we form the augmented Lagrangian

$$L_\rho(x, z, \lambda) = f(x) + g(z) + \lambda^\top (Ax + Bz - c) + (\rho/2) \|Ax + Bz - c\|_2^2,$$

with  $\lambda \in \mathbb{R}^p$  and for some  $\rho > 0$ . The major difference with the method of multipliers is that instead of performing a minimization step jointly over  $x$  and  $z$ , ADMM first performs an  $x$ -minimization step and then a  $z$ -minimization step. Thus  $x$  and  $z$  are updated in an alternating or sequential fashion, which accounts for the term *alternating direction*. Because the Lagrangian is augmented, some mild conditions on  $A$  and  $B$  imply that these minimization steps are guaranteed to terminate. ADMM is presented in Section 16.3.

In Section 16.4, we prove the convergence of ADMM under exactly the same assumptions as in Boyd et al. [17]. It turns out that Assumption (2) in Boyd et al. [17] implies that the matrices  $A^\top A$  and  $B^\top B$  are invertible (as we show after the proof of Theorem 16.1). This allows us to prove a convergence result stronger than the convergence result proven in Boyd et al. [17]. In particular, we prove that *all* of the sequences  $(x^k)$ ,  $(z^k)$ , and  $(\lambda^k)$  converge to optimal solutions  $(\tilde{x}, \tilde{z})$ , and  $\tilde{\lambda}$ .

In Section 16.5, we discuss stopping criteria. In Section 16.6, we present some applications of ADMM, in particular, minimization of a proper closed convex function  $f$  over a closed convex set  $C$  in  $\mathbb{R}^n$  and quadratic programming. The second example provides one of the best methods for solving quadratic problems, in particular, the SVM problems discussed in Chapter 18. Section 16.8 gives applications of ADMM to  $\ell^1$ -norm problems, in particular, lasso regularization which plays an important role in machine learning.

The next four chapters constitute Part IV, which covers some applications of optimization theory (in particular Lagrangian duality) to machine learning.

Chapter 17 is an introduction to positive definite kernels and the use of kernel functions in machine learning.

Let  $X$  be a nonempty set. If the set  $X$  represents a set of highly nonlinear data, it may be advantageous to map  $X$  into a space  $F$  of much higher dimension called the *feature space*, using a function  $\varphi: X \rightarrow F$  called a *feature map*. This idea is that  $\varphi$  “unwinds” the description of the objects in  $F$  in an attempt to make it linear. The space  $F$  is usually a vector space equipped with an inner product  $\langle -, - \rangle$ . If  $F$  is infinite dimensional, then we assume that it is a Hilbert space.

Many algorithms that analyze or classify data make use of the inner products  $\langle \varphi(x), \varphi(y) \rangle$ , where  $x, y \in X$ . These algorithms make use of the function  $\kappa: X \times X \rightarrow \mathbb{C}$  given by

$$\kappa(x, y) = \langle \varphi(x), \varphi(y) \rangle, \quad x, y \in X,$$

called a *kernel function*.

The kernel trick is to pretend that we have a feature embedding  $\varphi: X \rightarrow F$  (actually unknown), but to only use inner products  $\langle \varphi(x), \varphi(y) \rangle$  that can be evaluated using the original data through the known kernel function  $\kappa$ . It turns out that the functions of the form  $\kappa$  as above can be defined in terms of a condition which is reminiscent of positive semidefinite matrices (see Definition 17.2). Furthermore, every function satisfying Definition 17.2 arises from a suitable feature map into a Hilbert space; see Theorem 17.8.

We illustrate the kernel methods on kernel PCA (see Section 17.4).

In Chapter 18 we return to the problem of separating two disjoint sets of points,  $\{u_i\}_{i=1}^p$  and  $\{v_j\}_{j=1}^q$ , but this time we do not assume that these two sets are separable. To cope with nonseparability, we allow points to invade the safety zone around the separating hyperplane, and even points on the wrong side of the hyperplane. Such a method is called soft margin support vector machine. We discuss variations of this method, including  $\nu$ -SV classification. In each case we present a careful derivation of the dual and we explain how to solve it using ADMM. We prove rigorous results about the existence of support vectors.

In Chapter 19 we discuss linear regression. This problem can be cast as a learning problem. We observe a sequence of (distinct) pairs  $((x_1, y_1), \dots, (x_m, y_m))$  called a *set of training data*, where  $x_i \in \mathbb{R}^n$  and  $y_i \in \mathbb{R}$ , viewed as input-output pairs of some unknown function  $f$  that we are trying to infer. The simplest kind of function is a linear function  $f(x) = x^\top w$ , where  $w \in \mathbb{R}^n$  is a vector of coefficients usually called a *weight vector*. Since the problem is overdetermined and since our observations may be subject to errors, we can't solve for  $w$  exactly as the solution of the system  $Xw = y$ , so instead we solve the least-squares problem of minimizing  $\|Xw - y\|_2^2$ , where  $X$  is the  $m \times n$  matrix whose *rows* are the row vectors  $x_i^\top$ . In general there are still infinitely many solutions so we add a regularizing term. If we add the term  $K \|w\|_2^2$  to the objective function  $J(w) = \|Xw - y\|_2^2$ , then we have *ridge regression*. This problem is discussed in Section 19.1.

We derive the dual program. The dual has a unique solution which yields a solution of the primal. However, the solution of the dual is given in terms of the matrix  $XX^\top$  (whereas the solution of the primal is given in terms of  $X^\top X$ ), and since our data points  $x_i$  are represented by the rows of the matrix  $X$ , we see that this solution only involves inner products of the  $x_i$ . This observation is the core of the idea of kernel functions, which we introduce. We also explain how to solve the problem of learning an affine function  $f(x) = x^\top w + b$ .

In general the vectors  $w$  produced by ridge regression have few zero entries. In practice it is highly desirable to obtain sparse solutions, that is, vectors  $w$  with many components equal to zero. This can be achieved by replacing the regularizing term  $K \|w\|_2^2$  by the regularizing term  $K \|w\|_1$ ; that is, to use the  $\ell^1$ -norm instead of the  $\ell^2$ -norm; see Section 19.4. This method has the exotic name of *lasso regression*. This time there is no closed-form solution, but this is a convex optimization problem and there are efficient iterative methods to solve it. We show that ADMM provides an efficient solution.

Lasso has some undesirable properties, in particular when the dimension of the data is much larger than the number of data. In order to alleviate these problems, *elastic net regression* penalizes  $w$  with *both* an  $\ell^2$  regularizing term  $K \|w\|_2^2$  and an  $\ell^1$  regularizing term  $\tau \|w\|_1$ . The method of elastic net blends ridge regression and lasso and attempts to retain their best properties; see Section 19.6. It can also be solved using ADMM but it appears to be much slower than lasso when  $K$  is small and the dimension of the data is much larger than the number of data.

In Chapter 20 we present  $\nu$ -SV Regression. This method is designed in the same spirit

as soft margin SVM, in the sense that it allows a margin of error. Here the errors are penalized in the  $\ell^1$ -sense. We discuss several variations of the method and show how to solve them using ADMM. We present a careful derivation of the dual and discuss the existence of support vectors.

