

Nonlinear Programming

SECOND EDITION

Dimitri P. Bertsekas

Massachusetts Institute of Technology

WWW site for book information and orders

<http://world.std.com/~athenasc/index.html>



Athena Scientific, Belmont, Massachusetts

ATHENA SCIENTIFIC
OPTIMIZATION AND COMPUTATION SERIES

1. Dynamic Programming and Optimal Control, Vols. I and II, by Dimitri P. Bertsekas, 1995, ISBN 1-886529-11-6, 704 pages
2. Nonlinear Programming, Second Edition, by Dimitri P. Bertsekas, 1999, ISBN 1-886529-00-0, 791 pages
3. Neuro-Dynamic Programming, by Dimitri P. Bertsekas and John N. Tsitsiklis, 1996, ISBN 1-886529-10-8, 512 pages
4. Constrained Optimization and Lagrange Multiplier Methods, by Dimitri P. Bertsekas, 1996, ISBN 1-886529-04-3, 410 pages
5. Stochastic Optimal Control: The Discrete-Time Case by Dimitri P. Bertsekas and Steven E. Shreve, 1996, ISBN 1-886529-03-5, 330 pages
6. Introduction to Linear Optimization by Dimitris Bertsimas and John N. Tsitsiklis, 1997, ISBN 1-886529-19-1, 608 pages
7. Parallel and Distributed Computation: Numerical Methods by Dimitri P. Bertsekas and John N. Tsitsiklis, 1997, ISBN 1-886529-01-9, 718 pages
8. Network Flows and Monotropic Optimization by R. Tyrrell Rockafellar, 1998, ISBN 1-886529-06-X, 634 pages
9. Network Optimization: Continuous and Discrete Models by Dimitri P. Bertsekas, 1998, ISBN 1-886529-02-7, 608 pages

Contents

1. Unconstrained Optimization	p. 1
1.1. Optimality Conditions	p. 4
1.1.1. Variational Ideas	p. 4
1.1.2. Main Optimality Conditions	p. 13
1.2. Gradient Methods – Convergence	p. 22
1.2.1. Descent Directions and Stepsize Rules	p. 22
1.2.2. Convergence Results	p. 43
1.3. Gradient Methods – Rate of Convergence	p. 62
1.3.1. The Local Analysis Approach	p. 64
1.3.2. The Role of the Condition Number	p. 65
1.3.3. Convergence Rate Results	p. 75
1.4. Newton's Method and Variations	p. 88
1.5. Least Squares Problems	p. 102
1.5.1. The Gauss-Newton Method	p. 107
1.5.2. Incremental Gradient Methods*	p. 108
1.5.3. Incremental Forms of the Gauss-Newton Method*	p. 119
1.6. Conjugate Direction Methods	p. 130
1.7. Quasi-Newton Methods	p. 148
1.8. Nonderivative Methods	p. 158
1.8.1. Coordinate Descent	p. 160
1.8.2. Direct Search Methods	p. 162
1.9. Discrete-Time Optimal Control Problems*	p. 166
1.10. Some Practical Guidelines	p. 183
1.11. Notes and Sources	p. 187
 2. Optimization Over a Convex Set	 p. 191
2.1. Optimality Conditions	p. 192
2.2. Feasible Directions and the Conditional Gradient Method	p. 209
2.2.1. Descent Directions and Stepsize Rules	p. 210
2.2.2. The Conditional Gradient Method	p. 215
2.3. Gradient Projection Methods	p. 223
2.3.1. Feasible Directions and Stepsize Rules Based on Projection	p. 223

2.3.2. Convergence Analysis*	p. 234
2.4. Two-Metric Projection Methods	p. 244
2.5. Manifold Suboptimization Methods	p. 250
2.6. Affine Scaling for Linear Programming	p. 259
2.7. Block Coordinate Descent Methods*	p. 267
2.8. Notes and Sources	p. 272
3. Lagrange Multiplier Theory	p. 275
3.1. Necessary Conditions for Equality Constraints	p. 277
3.1.1. The Penalty Approach	p. 281
3.1.2. The Elimination Approach	p. 283
3.1.3. The Lagrangian Function	p. 287
3.2. Sufficient Conditions and Sensitivity Analysis	p. 295
3.2.1. The Augmented Lagrangian Approach	p. 297
3.2.2. The Feasible Direction Approach	p. 300
3.2.3. Sensitivity*	p. 301
3.3. Inequality Constraints	p. 307
3.3.1. Karush-Kuhn-Tucker Optimality Conditions	p. 309
3.3.2. Conversion to the Equality Case*	p. 312
3.3.3. Second Order Sufficiency Conditions and Sensitivity*	p. 314
3.3.4. Sufficiency Conditions and Lagrangian Minimization*	p. 315
3.3.5. Fritz John Optimality Conditions*	p. 317
3.3.6. Refinements*	p. 330
3.4. Linear Constraints and Duality*	p. 357
3.4.1. Convex Cost Functions and Linear Constraints	p. 357
3.4.2. Duality Theory: A Simple Form for Linear Constraints	p. 359
3.5. Notes and Sources	p. 367
4. Lagrange Multiplier Algorithms	p. 369
4.1. Barrier and Interior Point Methods	p. 370
4.1.1. Linear Programming and the Logarithmic Barrier*	p. 373
4.2. Penalty and Augmented Lagrangian Methods	p. 388
4.2.1. The Quadratic Penalty Function Method	p. 390
4.2.2. Multiplier Methods – Main Ideas	p. 398
4.2.3. Convergence Analysis of Multiplier Methods*	p. 407
4.2.4. Duality and Second Order Multiplier Methods*	p. 410
4.2.5. The Exponential Method of Multipliers*	p. 413
4.3. Exact Penalties – Sequential Quadratic Programming*	p. 421
4.3.1. Nondifferentiable Exact Penalty Functions	p. 422
4.3.2. Differentiable Exact Penalty Functions	p. 439
4.4. Lagrangian and Primal-Dual Interior Point Methods*	p. 446
4.4.1. First-Order Methods	p. 446
4.4.2. Newton-Like Methods for Equality Constraints	p. 450
4.4.3. Global Convergence	p. 460

4.4.4. Primal-Dual Interior Point Methods	p. 463
4.4.5. Comparison of Various Methods	p. 471
4.5. Notes and Sources	p. 473
5. Duality and Convex Programming	p. 477
5.1. The Dual Problem	p. 479
5.1.1. Lagrange Multipliers	p. 480
5.1.2. The Weak Duality Theorem	p. 485
5.1.3. Characterization of Primal and Dual Optimal Solutions	p. 490
5.1.4. The Case of an Infeasible or Unbounded Primal Problem	p. 491
5.1.5. Treatment of Equality Constraints	p. 493
5.1.6. Separable Problems and Their Geometry	p. 494
5.1.7. Additional Issues About Duality	p. 498
5.2. Convex Cost – Linear Constraints*	p. 503
5.2.1. Proofs of Duality Theorems	p. 505
5.3. Convex Cost – Convex Constraints	p. 511
5.4. Conjugate Functions and Fenchel Duality*	p. 521
5.4.1. Monotropic Programming Duality	p. 525
5.4.2. Network Optimization	p. 529
5.4.3. Games and the Minimax Theorem	p. 531
5.4.4. The Primal Function	p. 534
5.4.5. A Dual View of Penalty Methods	p. 536
5.4.6. The Proximal and Entropy Minimization Algorithms	p. 542
5.5. Discrete Optimization and Duality	p. 558
5.5.1. Examples of Discrete Optimization Problems	p. 559
5.5.2. Branch-and-Bound	p. 567
5.5.3. Lagrangian Relaxation	p. 576
5.6. Notes and Sources	p. 587
6. Dual Methods	p. 591
6.1. Dual Derivatives and Subgradients*	p. 594
6.2. Dual Ascent Methods for Differentiable Dual Problems*	p. 600
6.2.1. Coordinate Ascent for Quadratic Programming	p. 600
6.2.2. Decomposition and Primal Strict Convexity	p. 603
6.2.3. Partitioning and Dual Strict Concavity	p. 604
6.3. Nondifferentiable Optimization Methods*	p. 609
6.3.1. Subgradient Methods	p. 610
6.3.2. Approximate and Incremental Subgradient Methods	p. 614
6.3.3. Cutting Plane Methods	p. 618
6.3.4. Ascent and Approximate Ascent Methods	p. 625
6.4. Decomposition Methods*	p. 638
6.4.1. Lagrangian Relaxation of the Coupling Constraints	p. 639
6.4.2. Decomposition by Right-Hand Side Allocation	p. 642
6.5. Notes and Sources	p. 645

Appendix A: Mathematical Background	p. 647
A.1. Vectors and Matrices	p. 648
A.2. Norms, Sequences, Limits, and Continuity	p. 649
A.3. Square Matrices and Eigenvalues	p. 656
A.4. Symmetric and Positive Definite Matrices	p. 659
A.5. Derivatives	p. 664
A.6. Contraction Mappings	p. 669
Appendix B: Convex Analysis	p. 671
B.1. Convex Sets and Functions	p. 671
B.2. Separating Hyperplanes	p. 689
B.3. Cones and Polyhedral Convexity	p. 694
B.4. Extreme Points	p. 701
B.5. Differentiability Issues	p. 707
Appendix C: Line Search Methods	p. 723
C.1. Cubic Interpolation	p. 723
C.2. Quadratic Interpolation	p. 724
C.3. The Golden Section Method	p. 726
Appendix D: Implementation of Newton's Method	p. 729
D.1. Cholesky Factorization	p. 729
D.2. Application to a Modified Newton Method	p. 731
References	p. 735
Index	p. 773

Preface

Nonlinear programming is a mature field that has experienced major developments in the last ten years. The first such development is the merging of linear and nonlinear programming algorithms through the use of interior point methods. This has resulted in a profound rethinking of how we solve linear programming problems, and in a major reassessment of how we treat constraints in nonlinear programming. A second development, less visible but still important, is the increased emphasis on large-scale problems, and the associated algorithms that take advantage of problem structure as well as parallel hardware. A third development has been the extensive use of iterative unconstrained optimization to solve the difficult least squares problems arising in the training of neural networks. As a result, simple gradient-like methods and stepsize rules have attained increased importance.

The purpose of this book is to provide an up-to-date, comprehensive, and rigorous account of nonlinear programming at the beginning graduate student level. In addition to the classical topics, such as descent algorithms, Lagrange multiplier theory, and duality, some of the important recent developments are covered: interior point methods for linear and nonlinear programs, major aspects of large-scale optimization, and least squares problems and neural network training.

A further noteworthy feature of the book is that it treats Lagrange multipliers and duality using two different and complementary approaches: a variational approach based on the implicit function theorem, and a convex analysis approach based on geometrical arguments. The former approach applies to a broader class of problems, while the latter is more elegant and more powerful for the convex programs to which it applies.

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

Chapter 1: This chapter covers unconstrained optimization: main concepts, optimality conditions, and algorithms. The material is classic, but there are discussions of topics frequently left untreated, such as the behavior of algorithms for singular problems, neural network training, and discrete-time optimal control.

based on projection on the subdifferential and ϵ -subdifferential, respectively, were first proposed by Bertsekas and Mitter [BeM71], [BeM73]. Bundle methods, proposed by Lemarechal [Lem74], [Lem75], and Wolfe [Wol75], provided effective implementations of ϵ -ascent ideas, and stimulated a great deal of subsequent research on nondifferentiable optimization; see e.g. the book by Hiriart-Urruty and Lemarechal [HiL93].

The texts by Auslender [Aus76], Shapiro [Sha79], Evtushenko [Evt85], Shor [Sho85], Minoux [Min86], Poljak [Pol87], Hiriart-Urruty and Lemarechal [HiL93], and Shor [Sho98] give extensive accounts of subgradient methods that complement our treatment and give many references.

Cutting plane methods were introduced by Cheney and Goldstein [ChG59], and by Kelley [Kel60]. For analysis of proximal cutting plane and related methods, see Ruszczyński [Rus89], Lemaréchal and Sagastizábal [LeS93], Mifflin [Mif96], Bonnans et. al. [BGL95], Kiwiel [Kiw97b], Burke and Qian [BuQ98], and Mifflin, Sun, and Qi [MSQ98].

Central cutting plane methods were introduced by Elzinga and Moore [EIM75]. More recent proposals, some of which relate to interior point methods, are discussed in Goffin and Vial [GoV90], Goffin, Haurie, and Vial [GHV92], Ye [Ye92], Kortanek and No [KoN93], Goffin, Luo, and Ye [GLY94], Atkinson and Vaidya [AtV95], den Hertog et. al. [HKR95], Nesterov [Nes95], Goffin, Luo, and Ye [GLY96]. For a textbook treatment, see Ye [Ye97], and for a recent survey, see Goffin and Vial [GoV99].

Section 6.4: Three historically important references on decomposition methods are Dantzig and Wolfe [DaW60], Benders [Ben62], and Everett [Eve63]. The early text by Lasdon [Las70] on large-scale optimization was particularly influential; see also Geoffrion [Geo70], [Geo74].

The theoretical and applications literature on large-scale optimization and decomposition is quite voluminous. We provide a few references that complement the material we have covered in this chapter: Stephanopoulos and Westerberg [StW75], Kennington and Shalaby [KeS77], Bertsekas [Ber79a], Meyer [Mey79], Cohen [Coh80], Fortin and Glowinski [FoG83], Birge [Bir85], Golshtein [Gol85], Tanikawa and Mukai [TaM85], Spingarn [Spi85], Minoux [Min86], Ruszczyński [Rus86], Sen and Serali [SeS86], Bertsekas and Tsitsiklis [BeT89], Hearn and Lawphongpanich [HeL89], Rockafellar [Roc90], Toint and Tuytens [ToT90], Ferris and Mangasarian [FeM91], Kim and Nazareth [KiN91], Rockafellar and Wets [RoW91], Tseng [Tse91b], [Tse91c], Auslender [Aus92], Eckstein and Bertsekas [EcB92], Fukushima [Fuk92], Gaudioso and Monaco [GaM92], Mulvey and Ruszczyński [MuR92], Pinar and Zenios [PiZ92], Nagurney [Nag93], Patriksson [Pat93a], [Pat93b], Tseng [Tse93], Eckstein [Eck94b], Migdalas [Mig94], Pinar and Zenios [PiZ94], Mahey, Oualibouch, and Tao [MOT95], Mulvey and Ruszczyński [MuR95], Zhu [Zhu95], Censor and Zenios [1997], Kontogiorgis and Meyer [KoM98], Patriksson [Pat98], Zhao and Luh [ZhL98].

APPENDIX A: *Mathematical Background*

In this appendix, we collect definitions, notational conventions, and several results from linear algebra and analysis that are used extensively in nonlinear programming. Only a few proofs are given. Additional proofs can be found in Appendix A of the book by Bertsekas and Tsitsiklis [BeT89], which provides a similar but more extended summary of linear algebra and analysis. Related and additional material can be found in the books by Hager [Hag88], Hoffman and Kunze [HoK71], Lancaster and Tismenetsky [LaT85], and Strang [Str76] (linear algebra), and the books by Ash [Ash72], Ortega and Rheinboldt [OrR70], and Rudin [Rud76] (analysis).

Notation

If S is a set and x is an element of S , we write $x \in S$. A set can be specified in the form $S = \{x \mid x \text{ satisfies } P\}$, as the set of all elements satisfying property P . The union of two sets S and T is denoted by $S \cup T$ and their intersection by $S \cap T$. The symbols \exists and \forall have the meanings “there exists” and “for all,” respectively. The set of real numbers (also referred to as scalars) is denoted by \mathbb{R} .

If a and b are real numbers or $+\infty$, $-\infty$, we denote by $[a, b]$ the set of numbers x satisfying $a \leq x \leq b$ (including the possibility $x = +\infty$ or $x = -\infty$). A rounded, instead of square, bracket denotes strict inequality in the definition. Thus (a, b) , $[a, b)$, and (a, b) denote the set of all x satisfying $a < x \leq b$, $a \leq x < b$, and $a < x < b$, respectively.

If f is a function, we use the notation $f : A \mapsto B$ to indicate the fact that f is defined on a set A (its *domain*) and takes values in a set B (its *range*).

A.1 VECTORS AND MATRICES

We denote by \mathbb{R} the real line and by \mathbb{R}^n the set of n -dimensional real vectors. For any $x \in \mathbb{R}^n$, we use x_i to indicate its i th *coordinate*, also called its i th *component*.

Vectors in \mathbb{R}^n will be viewed as column vectors, unless the contrary is explicitly stated. For any $x \in \mathbb{R}^n$, x' denotes the transpose of x , which is an n -dimensional row vector. The *inner product* of two vectors $x, y \in \mathbb{R}^n$ is defined by $x'y = \sum_{i=1}^n x_i y_i$. Any two vectors $x, y \in \mathbb{R}^n$ satisfying $x'y = 0$ are called *orthogonal*.

If w is a vector in \mathbb{R}^n , the notations $w > 0$ and $w \geq 0$ indicate that all coordinates of w are positive or nonnegative, respectively. For any two vectors w, v , the notation $w > v$ means that $w - v > 0$. The notations $w \geq v$, $w < v$, etc., are to be interpreted accordingly.

Subspaces and Linear Independence

A subset S of \mathbb{R}^n is called a *subspace* of \mathbb{R}^n if $ax + by \in S$ for every $x, y \in S$ and every $a, b \in \mathbb{R}$. A *linear manifold* in \mathbb{R}^n is a translated subspace, that is, a set of the form

$$y + S = \{y + x \mid x \in S\},$$

where y is a vector in \mathbb{R}^n and S is a subspace of \mathbb{R}^n . The *span* of a finite collection $\{x_1, \dots, x_m\}$ of elements of \mathbb{R}^n is the subspace consisting of all vectors y of the form $y = \sum_{k=1}^m a_k x_k$, where each a_k is a scalar.

The vectors $x_1, \dots, x_m \in \mathbb{R}^n$ are called *linearly independent* if there exists no set of scalars a_1, \dots, a_m such that $\sum_{k=1}^m a_k x_k = 0$, unless $a_k = 0$ for each k . An equivalent definition is that $x_1 \neq 0$ and for every $k > 1$, the vector x_k does not belong to the span of x_1, \dots, x_{k-1} .

Given a subspace S of \mathbb{R}^n containing at least one nonzero vector, a *basis* for S is a collection of vectors that are linearly independent and whose span is equal to S . Every basis of a given subspace has the same number of vectors. This number is called the *dimension* of S . By convention, the subspace $\{0\}$ is said to have dimension zero. The *dimension of a linear manifold* $y + S$ is the dimension of the corresponding subspace S . An important fact is that every subspace of nonzero dimension has an *orthogonal basis*, that is, a basis consisting of mutually orthogonal vectors.

Matrices

For any matrix A , we use A_{ij} , $[A]_{ij}$, or a_{ij} to denote its ij th entry. The *transpose* of A , denoted by A' , is defined by $[A']_{ij} = a_{ji}$. For any two matrices A and B of compatible dimensions, we have $(AB)' = B'A'$.

Let A be a square matrix. We say that A is *symmetric* if $A' = A$. We say that A is *diagonal* if $[A]_{ij} = 0$ whenever $i \neq j$. It is *lower triangular*

if $[A]_{ij} = 0$ whenever $i < j$. It is *upper triangular* if its transpose is lower triangular. We use I to denote the identity matrix. The *determinant* of A is denoted by $\det(A)$.

Let A be an $m \times n$ matrix. The *range space* of A is the set of all vectors $y \in \mathbb{R}^m$ such that $y = Ax$ for some $x \in \mathbb{R}^n$. The *null space* or *kernel* of A is the set of all vectors $x \in \mathbb{R}^n$ such that $Ax = 0$. It is seen that the range space and the null space of A are subspaces. The *rank* of A is the minimum of the dimensions of the range space of A and the range space of the transpose A' . Clearly A and A' have the same rank. We say that A has *full rank*, if its rank is equal to $\min\{m, n\}$. It can be seen that A has full rank if and only if either the rows of A are linearly independent, or the columns of A are linearly independent.

A.2 NORMS, SEQUENCES, LIMITS, AND CONTINUITY

Definition A.1: A *norm* $\|\cdot\|$ on \mathbb{R}^n is a mapping that assigns a scalar $\|x\|$ to every $x \in \mathbb{R}^n$ and that has the following properties:

- (a) $\|x\| \geq 0$ for all $x \in \mathbb{R}^n$.
- (b) $\|cx\| = |c| \cdot \|x\|$ for every $c \in \mathbb{R}$ and every $x \in \mathbb{R}^n$.
- (c) $\|x\| = 0$ if and only if $x = 0$.
- (d) $\|x + y\| \leq \|x\| + \|y\|$ for all $x, y \in \mathbb{R}^n$.

The *Euclidean norm* is defined by

$$\|x\| = (x'x)^{1/2} = \left(\sum_{i=1}^n |x_i|^2 \right)^{1/2}.$$

The space \mathbb{R}^n , equipped with this norm, is called a *Euclidean space*. We will use the Euclidean norm almost exclusively in this book. In particular, *in the absence of a clear indication to the contrary*, $\|\cdot\|$ will denote the *Euclidean norm*. Two important results for the Euclidean norm are:

Proposition A.1: (Pythagorean Theorem) If x and y are orthogonal then

$$\|x + y\|^2 = \|x\|^2 + \|y\|^2.$$

Proposition A.2: (Schwartz inequality) For any two vectors x and y , we have

$$|x'y| \leq \|x\| \cdot \|y\|,$$

with equality holding if and only if $x = \alpha y$ for some scalar α .

Two other important norms are the *maximum norm* $\|\cdot\|_\infty$ (also called *sup-norm* or ℓ_∞ -norm), defined by

$$\|x\|_\infty = \max_i |x_i|,$$

and the ℓ_1 -norm $\|\cdot\|_1$, defined by

$$\|x\|_1 = \sum_{i=1}^n |x_i|.$$

Sequences

We use both subscripts and superscripts in sequence notation. Generally, we use superscript notation for sequences of vectors generated by iterative algorithms whenever we need to reserve the subscript notation for indexing coordinates or components of vectors and functions.

A sequence $\{x_k \mid k = 1, 2, \dots\}$ (or $\{x_k\}$ for short) of scalars is said to *converge* to a scalar x if for every $\epsilon > 0$ there exists some K (depending on ϵ) such that $|x_k - x| < \epsilon$ for every $k \geq K$. A sequence $\{x_k\}$ is said to converge to ∞ (respectively, $-\infty$) if for every b there exists some K (depending on b) such that $x_k \geq b$ (respectively, $x_k \leq b$) for all $k \geq K$. If a sequence $\{x_k\}$ converges to some x (possibly infinite), we say that x is the *limit* of $\{x_k\}$; symbolically, $x_k \rightarrow x$ or $\lim_{k \rightarrow \infty} x_k = x$. A sequence $\{x_k\}$ is called a *Cauchy sequence* if for every $\epsilon > 0$, there exists some K (depending on ϵ) such that $|x_k - x_m| < \epsilon$ for all $k \geq K$ and $m \geq K$.

A sequence $\{x_k\}$ is said to be *bounded above* (respectively, *below*) if there exists some scalar b such that $x_k \leq b$ (respectively, $x_k \geq b$) for all k . It is said to be *bounded* if it is bounded above and bounded below. The sequence $\{x_k\}$ is said to be *nonincreasing* (respectively, *nondecreasing*) if $x_{k+1} \leq x_k$ (respectively, $x_{k+1} \geq x_k$) for all k . If $\{x_k\}$ converges to x and is nonincreasing (nondecreasing) we also use the notation $x_k \downarrow x$ ($x_k \uparrow x$, respectively).

Proposition A.3: Every nonincreasing or nondecreasing scalar sequence converges to a possibly infinite number. If it is also bounded, then it converges to a finite real number.

The *supremum* of a nonempty set A of scalars, denoted by $\sup A$, is defined as the smallest scalar x such that $x \geq y$ for all $y \in A$. If no such scalar exists, we say that the supremum of A is ∞ . Similarly, the *infimum* of A , denoted by $\inf A$, is defined as the largest scalar x such that $x \leq y$ for all $y \in A$, and is equal to $-\infty$ if no such scalar exists. Given a scalar sequence $\{x_k\}$, the supremum of the sequence, denoted by $\sup_k x_k$, is defined as $\sup\{x_k \mid k = 1, 2, \dots\}$. The infimum of a sequence is similarly defined. Given a sequence $\{x_k\}$, let $y_m = \sup\{x_k \mid k \geq m\}$, $z_m = \inf\{x_k \mid k \geq m\}$. The sequences $\{y_m\}$ and $\{z_m\}$ are nonincreasing and nondecreasing, respectively, and therefore have a (possibly infinite) limit (Prop. A.3). The limit of y_m is denoted by $\limsup_{m \rightarrow \infty} x_m$ and the limit of z_m is denoted by $\liminf_{m \rightarrow \infty} x_m$.

Proposition A.4: Let $\{x_k\}$ be a scalar sequence.

(a) There holds

$$\inf_k x_k \leq \liminf_{k \rightarrow \infty} x_k \leq \limsup_{k \rightarrow \infty} x_k \leq \sup_k x_k.$$

(b) $\{x_k\}$ converges if and only if $\liminf_{k \rightarrow \infty} x_k = \limsup_{k \rightarrow \infty} x_k$ and, in that case, both of these quantities are equal to the limit of x_k .

(c) If $x_k \leq y_k$, then

$$\begin{aligned} \liminf_{k \rightarrow \infty} x_k &\leq \liminf_{k \rightarrow \infty} y_k, \\ \limsup_{k \rightarrow \infty} x_k &\leq \limsup_{k \rightarrow \infty} y_k. \end{aligned}$$

A sequence $\{x_k\}$ of vectors in \mathbb{R}^n is said to converge to some $x \in \mathbb{R}^n$ if the i th coordinate of x_k converges to the i th coordinate of x for every i . We use the notations $x_k \rightarrow x$ and $\lim_{k \rightarrow \infty} x_k = x$ to indicate convergence for vector sequences as well. The sequence $\{x_k\}$ is called *bounded* (or a *Cauchy sequence*) if each of its corresponding coordinate sequences is bounded (or a *Cauchy sequence*, respectively).

Definition A.2: We say that a vector $x \in \mathbb{R}^n$ is a *limit point* of a sequence $\{x_k\}$ in \mathbb{R}^n if there exists a subsequence of $\{x_k\}$ that converges to x . Let A be a subset of \mathbb{R}^n . We say that $x \in \mathbb{R}^n$ is a *limit point* of A if there exists a sequence $\{x_k\}$, consisting of elements of A , that converges to x .

Proposition A.5:

- (a) A bounded sequence of vectors in \mathbb{R}^n converges if and only if it has a unique limit point.
- (b) A sequence in \mathbb{R}^n converges if and only if it is a Cauchy sequence.
- (c) Every bounded sequence in \mathbb{R}^n has at least one limit point.
- (d) Let $\{x_k\}$ be a scalar sequence. If $\limsup_{k \rightarrow \infty} x_k$ ($\liminf_{k \rightarrow \infty} x_k$) is finite, then it is the largest (respectively, smallest) limit point of $\{x_k\}$.

 $o(\cdot)$ Notation

If p is a positive integer and $h : \mathbb{R}^n \mapsto \mathbb{R}^m$, then we write

$$h(x) = o(\|x\|^p)$$

if and only if

$$\lim_{x_k \rightarrow 0} \frac{h(x_k)}{\|x_k\|^p} = 0,$$

for all sequences $\{x_k\}$, with $x_k \neq 0$ for all k , that converge to 0.

Closed and Open Sets

Definition A.3: A set $A \subset \mathbb{R}^n$ is called *closed* if it contains all of its limit points. It is called *open* if its complement (the set $\{x \mid x \notin A\}$) is closed. It is called *bounded* if there exists some $c \in \mathbb{R}$ such that the magnitude of any coordinate of any element of A is less than c . The subset A is called *compact* if every sequence of elements of A has a subsequence that converges to an element of A . A *neighborhood* of a vector x is an open set containing x . If $A \subset \mathbb{R}^n$ and $x \in A$, we say that x is an *interior* point of A if there exists a neighborhood of x that is contained in A . A vector $x \in A$ which is not an interior point of A is said to be a *boundary* point of A .

For any norm $\|\cdot\|$ in \mathbb{R}^n , and any $\epsilon > 0$ and $x^* \in \mathbb{R}^n$, consider the sets

$$\{x \mid \|x - x^*\| < \epsilon\}, \quad \{x \mid \|x - x^*\| \leq \epsilon\}.$$

The first set is open and is called an *open sphere* centered at x^* , while the second set is closed and is called a *closed sphere* centered at x^* . Sometimes the terms *open ball* and *closed ball* are used, respectively.

Proposition A.6:

- (a) The union of finitely many closed sets is closed.
- (b) The intersection of closed sets is closed.
- (c) The union of open sets is open.
- (d) The intersection of finitely many open sets is open.
- (e) A set is open if and only if all of its elements are interior points.
- (f) Every subspace of \mathbb{R}^n is closed.
- (g) A subset of \mathbb{R}^n is compact if and only if it is closed and bounded.

Continuity

Let A be a subset of \mathbb{R}^m and let $f : A \mapsto \mathbb{R}^n$ be some function. Let x be a limit point of A . If the sequence $\{f(x_k)\}$ has a common limit z for every sequence $\{x_k\}$ of elements of A such that $\lim_{k \rightarrow \infty} x_k = x$, we write $\lim_{y \rightarrow x} f(y) = z$.

If A is a subset of \mathbb{R} and x is a limit point of A , the notation $\lim_{y \uparrow x} f(y)$ [respectively, $\lim_{y \downarrow x} f(y)$] will stand for the limit of $f(x_k)$, where $\{x_k\}$ is any sequence of elements of A converging to x and satisfying $x_k \leq x$ (respectively, $x_k \geq x$), assuming that the limit exists and is independent of the choice of the sequence $\{x_k\}$.

Definition A.4: Let A be a subset of \mathbb{R}^m .

- (a) A function $f : A \mapsto \mathbb{R}^n$ is said to be *continuous* at a point $x \in A$ if $\lim_{y \rightarrow x} f(y) = f(x)$. It is said to be continuous on A (or over A) if it is continuous at every point $x \in A$.
- (b) A real valued function $f : A \mapsto \mathbb{R}$ is called *upper semicontinuous* (respectively, *lower semicontinuous*) at a vector $x \in A$ if $f(x) \geq \limsup_{k \rightarrow \infty} f(x_k)$ [respectively, $f(x) \leq \liminf_{k \rightarrow \infty} f(x_k)$] for every sequence $\{x_k\}$ of elements of A converging to x .
- (c) A real valued function $f : A \mapsto \mathbb{R}$ is called *coercive* if

$$\lim_{k \rightarrow \infty} f(x_k) = \infty$$

for every sequence $\{x_k\}$ of elements of A such that $\|x_k\| \rightarrow \infty$ for some norm $\|\cdot\|$.

- (d) Let A be a subset of \mathbb{R} . A function $f : A \mapsto \mathbb{R}^n$ is called *right-continuous* (respectively, *left-continuous*) at a point $x \in A$ if $\lim_{y \downarrow x} f(y) = f(x)$ [respectively, $\lim_{y \uparrow x} f(y) = f(x)$].

It is easily seen that when A is a subset of \mathbb{R} , a nondecreasing and right-continuous (respectively, left-continuous) function $f : A \mapsto \mathbb{R}$ is upper (respectively, lower) semicontinuous.

Proposition A.7:

- (a) The composition of two continuous functions is continuous.
- (b) Any vector norm on \mathbb{R}^n is a continuous function.
- (c) Let $f : \mathbb{R}^m \mapsto \mathbb{R}^n$ be continuous, and let $A \subset \mathbb{R}^m$ be open (respectively, closed). Then the set $\{x \in \mathbb{R}^m \mid f(x) \in A\}$ is open (respectively, closed).

An important property of compactness in connection with optimization problems is the following theorem.

Proposition A.8: (Weierstrass' Theorem) Let A be a nonempty subset of \mathbb{R}^n and let $f : A \mapsto \mathbb{R}$ be lower semicontinuous at all points of A . Assume that one of the following three conditions holds:

- (1) A is compact.
- (2) A is closed and f is coercive.
- (3) There exists a scalar γ such that the level set

$$\{x \in A \mid f(x) \leq \gamma\}$$

is nonempty and compact.

Then, there exists a vector $x \in A$ such that $f(x) = \inf_{z \in A} f(z)$.

Proof: Assume condition (1). Let $\{z_k\} \subset A$ be a sequence such that

$$\lim_{k \rightarrow \infty} f(z_k) = \inf_{z \in A} f(z).$$

Since A is bounded, this sequence has at least one limit point x [Prop.

A.5(c)]. Since A is closed, x belongs to A , while the lower semicontinuity of f implies that $f(x) \leq \lim_{k \rightarrow \infty} f(z_k) = \inf_{z \in A} f(z)$. Therefore, we must have $f(x) = \inf_{z \in A} f(z)$.

Assume condition (2). Consider a sequence $\{z_k\}$ as in the proof of part (a). Since f is coercive, $\{z_k\}$ must be bounded and the proof proceeds like the proof of part (a).

Assume condition (3). If the given γ is equal to $\inf_{z \in A} f(z)$, the set of minima of f over A is $\{x \in A \mid f(x) \leq \gamma\}$, and since by assumption this set is nonempty, we are done. If $\gamma > \inf_{z \in A} f(z)$, consider a sequence $\{z_k\}$ as in the proof of part (a). Then, for all k sufficiently large, z_k must belong to the set $\{x \in A \mid f(x) \leq \gamma\}$. Since this set is compact, $\{z_k\}$ must be bounded and the proof proceeds like the proof of part (a). **Q.E.D.**

Note that with appropriate adjustments, the above proposition applies to the existence of maxima of f over A . In particular, if f is upper semicontinuous at all points of A and A is compact, then there exists a vector $y \in A$ such that $f(y) = \sup_{z \in A} f(z)$. Vectors $x \in A$ or $y \in A$ that attain the minimum or the maximum of a function f over a set A , respectively, *even if they are not unique*, are denoted by

$$x = \arg \min_{z \in A} f(z), \quad y = \arg \max_{z \in A} f(z).$$

Proposition A.9: For any two norms $\|\cdot\|$ and $\|\cdot\|'$ on \mathbb{R}^n , there exists some positive constant $c \in \mathbb{R}$ such that $\|x\| \leq c\|x\|'$ for all $x \in \mathbb{R}^n$.

Proof: Let a be the minimum of $\|x\|'$ over the set of all $x \in \mathbb{R}^n$ such that $\|x\| = 1$. The latter set is closed and bounded and, therefore, the minimum is attained at some \tilde{x} (Prop. A.8) that must be nonzero since $\|\tilde{x}\| = 1$. For any $x \in \mathbb{R}^n$, $x \neq 0$, the $\|\cdot\|$ norm of $x/\|x\|$ is equal to 1. Therefore,

$$0 < a = \|\tilde{x}\|' \leq \left\| \frac{x}{\|x\|} \right\|' = \frac{\|x\|'}{\|x\|}, \quad \forall x \neq 0,$$

which proves the desired result with $c = 1/a$. **Q.E.D.**

The preceding proposition is referred to as the *norm equivalence property* in \mathbb{R}^n . It shows that if a sequence converges with respect to one norm, it converges with respect to all other norms. From this we obtain the following.

Proposition A.10: If a subset of \mathbb{R}^n is open (respectively, closed, bounded, or compact) for some norm, it is open (respectively, closed, bounded, or compact), for all other norms.

A norm $\|\cdot\|$ on the set of $n \times n$ matrices is a real-valued mapping that has the same properties as vector norms do when the matrix is viewed as an element of \mathbb{R}^{n^2} . The norm of an $n \times n$ matrix A is denoted by $\|A\|$.

We are mainly interested in *induced norms*, which are constructed as follows. Given any vector norm $\|\cdot\|$, the corresponding induced matrix norm, also denoted by $\|\cdot\|$, is defined by

$$\|A\| = \max_{\{x \in \mathbb{R}^n \mid \|x\|=1\}} \|Ax\|. \quad (\text{A.1})$$

The set over which the maximization takes place above is closed [Prop. A.7(c)] and bounded; the function being maximized is continuous [Prop. A.7(b)] and therefore the maximum is attained (Prop. A.8). It is easily verified that for any vector norm, Eq. (A.1) defines a bona fide matrix norm having all the required properties.

Note that by the Schwartz inequality (Prop. A.2), we have

$$\|A\| = \max_{\|x\|=1} \|Ax\| = \max_{\|y\|=\|x\|=1} |y'Ax|.$$

By reversing the roles of x and y in the above relation and by using the equality $y'Ax = x'A'y$, it follows that

$$\|A\| = \|A'\|. \quad (\text{A.2})$$

A.3 SQUARE MATRICES AND EIGENVALUES

Definition A.5: A square matrix A is called *singular* if its determinant is zero. Otherwise it is called *nonsingular* or *invertible*.

Proposition A.11:

(a) Let A be an $n \times n$ matrix. The following are equivalent:

- (i) The matrix A is nonsingular.
- (ii) The matrix A' is nonsingular.
- (iii) For every nonzero $x \in \mathbb{R}^n$, we have $Ax \neq 0$.
- (iv) For every $y \in \mathbb{R}^n$, there exists a unique $x \in \mathbb{R}^n$ such that $Ax = y$.

- (v) There exists an $n \times n$ matrix B such that $AB = I = BA$.
- (vi) The columns of A are linearly independent.
- (vii) The rows of A are linearly independent.
- (b) Assuming that A is nonsingular, the matrix B of statement (v) (called the *inverse* of A and denoted by A^{-1}) is unique.
- (c) For any two square invertible matrices A and B of the same dimensions, we have $(AB)^{-1} = B^{-1}A^{-1}$.

Let A and B be square matrices and let C be a matrix of appropriate dimension. Then we have

$$(A + CBC')^{-1} = A^{-1} - A^{-1}C(B^{-1} + C'A^{-1}C)^{-1}C'A^{-1},$$

provided all the inverses appearing above exist. For a proof, multiply the right-hand side by $A + CBC'$ and show that the product is the identity.

Another useful formula provides the inverse of the partitioned matrix

$$M = \begin{bmatrix} A & B \\ C & D \end{bmatrix}.$$

There holds

$$M^{-1} = \begin{bmatrix} Q & -QBD^{-1} \\ -D^{-1}CQ & D^{-1} + D^{-1}CQBD^{-1} \end{bmatrix},$$

where

$$Q = (A - BD^{-1}C)^{-1},$$

provided all the inverses appearing above exist. For a proof, multiply M with the given expression for M^{-1} and verify that the product is the identity.

Definition A.6: The *characteristic polynomial* ϕ of an $n \times n$ matrix A is defined by $\phi(\lambda) = \det(\lambda I - A)$, where I is the identity matrix of the same size as A . The n (possibly repeated and complex) roots of ϕ are called the *eigenvalues* of A . A vector x (with possibly complex coordinates) such that $Ax = \lambda x$, where λ is an eigenvalue of A , is called an *eigenvector* of A associated with λ .

Proposition A.12: Let A be a square matrix.

- (a) A complex number λ is an eigenvalue of A if and only if there exists a nonzero eigenvector associated with λ .
- (b) A is singular if and only if it has an eigenvalue that is equal to zero.

Note that the only use of complex numbers in this book is in relation to eigenvalues and eigenvectors. All other matrices or vectors are implicitly assumed to have real components.

Proposition A.13: Let A be an $n \times n$ matrix.

- (a) The eigenvalues of a triangular matrix are equal to its diagonal entries.
- (b) If S is a nonsingular matrix and $B = SAS^{-1}$, then the eigenvalues of A and B coincide.
- (c) The eigenvalues of $cI + A$ are equal to $c + \lambda_1, \dots, c + \lambda_n$, where $\lambda_1, \dots, \lambda_n$ are the eigenvalues of A .
- (d) The eigenvalues of A^k are equal to $\lambda_1^k, \dots, \lambda_n^k$, where $\lambda_1, \dots, \lambda_n$ are the eigenvalues of A .
- (e) If A is nonsingular, then the eigenvalues of A^{-1} are the reciprocals of the eigenvalues of A .
- (f) The eigenvalues of A and A' coincide.

Definition A.7: The *spectral radius* $\rho(A)$ of a square matrix A is defined as the maximum of the magnitudes of the eigenvalues of A .

It can be shown that the roots of a polynomial depend continuously on the coefficients of the polynomial. For this reason, the eigenvalues of a square matrix A depend continuously on A , and we obtain the following.

Proposition A.14: The eigenvalues of a square matrix A depend continuously on the elements of A . In particular, $\rho(A)$ is a continuous function of A .

The next two propositions are fundamental for the convergence theory of linear iterative methods.

Proposition A.15: For any induced matrix norm $\|\cdot\|$ and any square matrix A we have

$$\lim_{k \rightarrow \infty} \|A^k\|^{1/k} = \rho(A) \leq \|A\|.$$

Furthermore, given any $\epsilon > 0$, there exists an induced matrix norm $\|\cdot\|$ such that

$$\|A\| = \rho(A) + \epsilon.$$

Proposition A.16: Let A be a square matrix. We have $\lim_{k \rightarrow \infty} A^k = 0$ if and only if $\rho(A) < 1$.

A corollary of the above proposition is that the iteration $x^{k+1} = Ax^k$ converges to 0 for every initial condition x^0 if and only if $\rho(A) < 1$.

A.4 SYMMETRIC AND POSITIVE DEFINITE MATRICES

Symmetric matrices have several special properties, particularly with respect to their eigenvalues and eigenvectors. In this section, $\|\cdot\|$ denotes the Euclidean norm throughout.

Proposition A.17: Let A be a symmetric $n \times n$ matrix. Then:

- (a) The eigenvalues of A are real.
- (b) The matrix A has a set of n mutually orthogonal, real, and nonzero eigenvectors x_1, \dots, x_n .
- (c) Suppose that the eigenvectors in part (b) have been normalized so that $\|x_i\| = 1$ for each i . Then

$$A = \sum_{i=1}^n \lambda_i x_i x_i',$$

where λ_i is the eigenvalue corresponding to x_i .

Proposition A.18: Let A be a symmetric $n \times n$ matrix, let $\lambda_1 \leq \dots \leq \lambda_n$ be its (real) eigenvalues, and let x_1, \dots, x_n be associated orthogonal eigenvectors, normalized so that $\|x_i\| = 1$ for all i . Then:

- (a) $\|A\| = \rho(A) = \max\{|\lambda_1|, |\lambda_n|\}$, where $\|\cdot\|$ is the matrix norm induced by the Euclidean norm.
- (b) $\lambda_1 \|y\|^2 \leq y' Ay \leq \lambda_n \|y\|^2$ for all $y \in \mathbb{R}^n$.
- (c) (*Courant-Fisher Minimax Principle*) For all i , and for all i -dimensional subspaces \bar{S}_i and all $(n-i+1)$ -dimensional subspaces \underline{S}_i , there holds

$$\min_{\|y\|=1, y \in \underline{S}_i} y' Ay \leq \lambda_i \leq \max_{\|y\|=1, y \in \bar{S}_i} y' Ay.$$

Furthermore, equality on the left (right) side above is attained if \underline{S}_i is the subspace spanned by x_i, \dots, x_n (\bar{S}_i is the subspace spanned by x_1, \dots, x_i , respectively).

- (d) (*Interlocking Eigenvalues Lemma*) Let $\tilde{\lambda}_1 \leq \tilde{\lambda}_2 \leq \dots \leq \tilde{\lambda}_n$ be the eigenvalues of $A + bb'$, where b is a vector in \mathbb{R}^n . Then,

$$\lambda_1 \leq \tilde{\lambda}_1 \leq \lambda_2 \leq \tilde{\lambda}_2 \leq \dots \leq \lambda_n \leq \tilde{\lambda}_n.$$

Proof: (a) We already know that $\|A\| \geq \rho(A)$ (Prop. A.15) and we need to show the reverse inequality. We express an arbitrary vector $y \in \mathbb{R}^n$ in the form $y = \sum_{i=1}^n \xi_i x_i$, where each ξ_i is a suitable scalar. Using the orthogonality of the vectors x_i and the Pythagorean theorem (Prop. A.1), we obtain $\|y\|^2 = \sum_{i=1}^n |\xi_i|^2 \cdot \|x_i\|^2$. Using the Pythagorean theorem again, we obtain

$$\|Ay\|^2 = \left\| \sum_{i=1}^n \lambda_i \xi_i x_i \right\|^2 = \sum_{i=1}^n |\lambda_i|^2 \cdot |\xi_i|^2 \cdot \|x_i\|^2 \leq \rho^2(A) \|y\|^2.$$

Since this is true for every y , we obtain $\|A\| \leq \rho(A)$ and the desired result follows.

(b) As in part (a), we express the generic $y \in \mathbb{R}^n$ as $y = \sum_{i=1}^n \xi_i x_i$. We have, using the orthogonality of the vectors x_i , $i = 1, \dots, n$, and the fact $\|x_i\| = 1$,

$$y' Ay = \sum_{i=1}^n \lambda_i |\xi_i|^2 \|x_i\|^2 = \sum_{i=1}^n \lambda_i |\xi_i|^2$$

and

$$\|y\|^2 = \sum_{i=1}^n |\xi_i|^2 \|x_i\|^2 = \sum_{i=1}^n |\xi_i|^2.$$

These two relations prove the desired result.

(c) Let \underline{X}_i be the subspace spanned by x_1, \dots, x_i . The subspaces \underline{X}_i and \underline{S}_i must have a common vector x_0 with $\|x_0\| = 1$, since the sum of their dimensions is $n+1$ [if there was no common nonzero vector, we could take sets of basis vectors for \underline{X}_i and \underline{S}_i (a total of $n+1$ in number), which would have to be linearly independent, yielding a contradiction]. The vector x_0 can be expressed as a linear combination $x_0 = \sum_{j=1}^i \xi_j x_j$, and since $\|x_0\| = 1$ and $\|x_i\| = 1$ for all $i = 1, \dots, n$, we must have

$$\sum_{j=1}^i \xi_j^2 = 1.$$

We also have using the expression $A = \sum_{j=1}^n \lambda_j x_j x_j'$ [cf. Prop. A.17(c)],

$$x_0' A x_0 = \sum_{j=1}^i \lambda_j \xi_j^2 \leq \lambda_i \left(\sum_{j=1}^i \xi_j^2 \right).$$

Combining the last two relations, we obtain $x_0' A x_0 \leq \lambda_i$, which proves the left-hand side of the desired inequality. The right-hand side is proved similarly. Furthermore, we have $x_i' A x_i = \lambda_i$, so equality is attained as in the final assertion.

(d) From part (c) we have

$$\lambda_i = \max_{\underline{S}_i} \min_{\|y\|=1, y \in \underline{S}_i} y' Ay \leq \max_{\underline{S}_i} \min_{\|y\|=1, y \in \underline{S}_i} y' (A + bb') y \leq \tilde{\lambda}_i,$$

so that $\lambda_i \leq \tilde{\lambda}_i$ for all i . Furthermore, from part (c), for some $(n-i+1)$ -dimensional subspace $\tilde{\underline{S}}_i$ we have

$$\tilde{\lambda}_i = \min_{\|y\|=1, y \in \tilde{\underline{S}}_i} y' (A + bb') y.$$

Using this relation and the left-hand side of the inequality of part (c), applied to the subspace $\{y \mid y \in \tilde{\underline{S}}_i, b'y = 0\}$, whose dimension is at least $(n-i)$, we obtain

$$\tilde{\lambda}_i \leq \min_{\|y\|=1, y \in \tilde{\underline{S}}_i, b'y=0} y' (A + bb') y = \min_{\|y\|=1, y \in \tilde{\underline{S}}_i, b'y=0} y' Ay \leq \lambda_{i+1},$$

and the proof is complete. **Q.E.D.**

Proposition A.19: Let A be a square matrix, and let $\|\cdot\|$ be the matrix norm induced by the Euclidean norm. Then:

- (a) If A is symmetric, then $\|A^k\| = \|A\|^k$ for any positive integer k .
- (b) $\|A\|^2 = \|A'A\| = \|AA'\|$.
- (c) If A is symmetric and nonsingular, then $\|A^{-1}\|$ is equal to the reciprocal of the smallest of the absolute values of the eigenvalues of A .

Proof: (a) If A is symmetric then A^k is symmetric. Using Prop. A.18(a), we have $\|A^k\| = \rho(A^k)$. Using Prop. A.13(d), we obtain $\rho(A^k) = \rho(A)^k$, which is equal to $\|A\|^k$ by Prop. A.18(a).

(b) For any vector x such that $\|x\| = 1$, we have, using the Schwartz inequality (Prop. A.2),

$$\|Ax\|^2 = x'A'A x \leq \|x\| \cdot \|A'A x\| \leq \|x\| \cdot \|A'A\| \cdot \|x\| = \|A'A\|.$$

Thus, $\|A\|^2 \leq \|A'A\|$. On the other hand,

$$\|A'A\| = \max_{\|y\|=\|x\|=1} |y'A'A x| \leq \max_{\|y\|=\|x\|=1} \|Ay\| \cdot \|Ax\| = \|A\|^2.$$

Therefore, $\|A\|^2 = \|A'A\|$. The equality $\|A\|^2 = \|A'A\|$ is obtained by replacing A by A' and using Eq. (A.2).

(c) This follows by combining Prop. A.13(e) with Prop. A.18(a). **Q.E.D.**

Definition A.8: A symmetric $n \times n$ matrix A is called *positive definite* if $x'A x > 0$ for all $x \in \mathbb{R}^n$, $x \neq 0$. It is called *nonnegative definite* or *positive semidefinite* if $x'A x \geq 0$ for all $x \in \mathbb{R}^n$.

Throughout this book, the notion of positive and negative definiteness applies exclusively to symmetric matrices. Thus *whenever we say that a matrix is positive or negative (semi)definite, we implicitly assume that the matrix is symmetric.*

Proposition A.20:

- (a) For any $m \times n$ matrix A , the matrix $A'A$ is symmetric and nonnegative definite. $A'A$ is positive definite if and only if A has rank n . In particular, if $m = n$, $A'A$ is positive definite if and only if A is nonsingular.

- (b) A square symmetric matrix is nonnegative definite (respectively, positive definite) if and only if all of its eigenvalues are nonnegative (respectively, positive).
- (c) The inverse of a symmetric positive definite matrix is symmetric and positive definite.

Proof: (a) Symmetry is obvious. For any vector $x \in \mathbb{R}^n$, we have $x'A'A x = \|Ax\|^2 \geq 0$, which establishes nonnegative definiteness. Positive definiteness is obtained if and only if the inequality is strict for every $x \neq 0$, which is the case if and only if $Ax \neq 0$ for every $x \neq 0$. This is equivalent to A having rank n .

(b) Let λ , $x \neq 0$, be an eigenvalue and a corresponding real eigenvector of a symmetric nonnegative definite matrix A . Then $0 \leq x'A x = \lambda x'x = \lambda \|x\|^2$, which proves that $\lambda \geq 0$. For the converse result, let y be an arbitrary vector in \mathbb{R}^n . Let $\lambda_1, \dots, \lambda_n$ be the eigenvalues of A , assumed to be nonnegative, and let x_1, \dots, x_n be a corresponding set of nonzero, real, and orthogonal eigenvectors. Let us express y in the form $y = \sum_{i=1}^n \xi_i x_i$. Then $y'A y = (\sum_{i=1}^n \xi_i x_i)' (\sum_{i=1}^n \xi_i \lambda_i x_i)$. From the orthogonality of the eigenvectors, the latter expression is equal to $\sum_{i=1}^n \xi_i^2 \lambda_i \|x_i\|^2 \geq 0$, which proves that A is nonnegative definite. The proof for the case of positive definite matrices is similar.

(c) The eigenvalues of A^{-1} are the reciprocal of the eigenvalues of A [Prop. A.13(e)], so the result follows using part (b). **Q.E.D.**

Proposition A.21: Let A be a square symmetric nonnegative definite matrix.

- (a) There exists a symmetric matrix Q with the property $Q^2 = A$. Such a matrix is called a *symmetric square root* of A and is denoted by $A^{1/2}$.
- (b) A symmetric square root $A^{1/2}$ is invertible if and only if A is invertible. Its inverse is denoted by $A^{-1/2}$.
- (c) There holds $A^{-1/2} A^{-1/2} = A^{-1}$.
- (d) There holds $AA^{1/2} = A^{1/2}A$.

Proof: (a) Let $\lambda_1, \dots, \lambda_n$ be the eigenvalues of A and let x_1, \dots, x_n be corresponding nonzero, real, and orthogonal eigenvectors normalized so

that $\|x_k\| = 1$ for each k . We let

$$A^{1/2} = \sum_{k=1}^n \lambda_k^{1/2} x_k x_k',$$

where $\lambda_k^{1/2}$ is the nonnegative square root of λ_k . We then have

$$A^{1/2} A^{1/2} = \sum_{i=1}^n \sum_{k=1}^n \lambda_i^{1/2} \lambda_k^{1/2} x_i x_i' x_k x_k' = \sum_{k=1}^n \lambda_k x_k x_k' = A.$$

Here the second equality follows from the orthogonality of distinct eigenvectors; the last equality follows from Prop. A.17(c). We now notice that each one of the matrices $x_k x_k'$ is symmetric, so $A^{1/2}$ is also symmetric.

(b) This follows from the fact that the eigenvalues of A are the squares of the eigenvalues of $A^{1/2}$ [Prop. A.13(d)].

(c) We have $(A^{-1/2} A^{-1/2}) A = A^{-1/2} (A^{-1/2} A^{1/2}) A^{1/2} = A^{-1/2} I A^{1/2} = I$.

(d) We have $AA^{1/2} = A^{1/2} A^{1/2} A^{1/2} = A^{1/2} A$. **Q.E.D.**

A symmetric square root of A is not unique. For example, let $A^{1/2}$ be as in the proof of Prop. A.21(a) and notice that the matrix $-A^{1/2}$ also has the property $(-A^{1/2})(-A^{1/2}) = A$. However, if A is positive definite, it can be shown that the matrix $A^{1/2}$ we have constructed is the only symmetric and positive definite square root of A .

A.5 DERIVATIVES

Let $f : \mathbb{R}^n \mapsto \mathbb{R}$ be some function, fix some $x \in \mathbb{R}^n$, and consider the expression

$$\lim_{\alpha \rightarrow 0} \frac{f(x + \alpha e_i) - f(x)}{\alpha},$$

where e_i is the i th unit vector (all components are 0 except for the i th component which is 1). If the above limit exists, it is called the i th *partial derivative* of f at the point x and it is denoted by $(\partial f / \partial x_i)(x)$ or $\partial f(x) / \partial x_i$ (x_i in this section will denote the i th coordinate of the vector x). Assuming all of these partial derivatives exist, the *gradient* of f at x is defined as the column vector

$$\nabla f(x) = \begin{pmatrix} \frac{\partial f(x)}{\partial x_1} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{pmatrix}.$$

For any $y \in \mathbb{R}^n$, we define the one-sided *directional derivative* of f in the direction y , to be

$$f'(x; y) = \lim_{\alpha \downarrow 0} \frac{f(x + \alpha y) - f(x)}{\alpha},$$

provided that the limit exists. We note from the definitions that

$$f'(x; e_i) = -f'(x; -e_i) \quad \Rightarrow \quad f'(x; e_i) = (\partial f / \partial x_i)(x).$$

If the directional derivative of f at a vector x exists in all directions y and $f'(x; y)$ is a linear function of y , we say that f is *differentiable* at x . This type of differentiability is also called *Gateaux differentiability*. It is seen that f is differentiable at x if and only if the gradient $\nabla f(x)$ exists and satisfies $\nabla f(x)' y = f'(x; y)$ for every $y \in \mathbb{R}^n$. The function f is called *differentiable over a given subset S of \mathbb{R}^n* if it is differentiable at every $x \in S$. The function f is called differentiable (without qualification) if it is differentiable at all $x \in \mathbb{R}^n$.

If f is differentiable over a set S and the gradient $\nabla f(x)$ is continuous at all $x \in S$, f is said to be *continuously differentiable over S* . Such a function is also continuous over S and has the property

$$\lim_{y \rightarrow 0} \frac{f(x + y) - f(x) - \nabla f(x)' y}{\|y\|} = 0, \quad \forall x \in S, \quad (\text{A.3})$$

where $\|\cdot\|$ is an arbitrary vector norm. The above equation can also be used as an alternative definition of differentiability. In particular, f is called *Frechet differentiable* at x if there exists a vector g satisfying Eq. (A.3) with $\nabla f(x)$ replaced by g . If such a vector g exists, it can be seen that all the partial derivatives $(\partial f / \partial x_i)(x)$ exist and that $g = \nabla f(x)$. Frechet differentiability implies (Gateaux) differentiability but not conversely (see for example [OrR70] for a detailed discussion). In this book, when dealing with a differentiable function f , we will always assume that f is continuously differentiable over a given set $[\nabla f(x)$ is a continuous function of x over that set], in which case f is both Gateaux and Frechet differentiable, and the distinctions made above are of no consequence.

Note that the definitions concerning differentiability of f at a point x only involve the values of f in a neighborhood of x . Thus, these definitions can be used for functions f that are not defined on all of \mathbb{R}^n , but are defined instead in a neighborhood of the point at which the derivative is computed.

If $f : \mathbb{R}^n \mapsto \mathbb{R}^m$ is a vector-valued function, it is called differentiable (respectively, continuously differentiable) if each component f_i of f is differentiable (respectively, continuously differentiable). The *gradient matrix* of f , denoted $\nabla f(x)$, is the $n \times m$ matrix whose i th column is the gradient $\nabla f_i(x)$ of f_i . Thus,

$$\nabla f(x) = [\nabla f_1(x) \cdots \nabla f_m(x)].$$

The transpose of ∇f is called the *Jacobian* of f and is a matrix whose ij th entry is equal to the partial derivative $\partial f_i / \partial x_j$.

Now suppose that each one of the partial derivatives of a function $f : \mathbb{R}^n \mapsto \mathbb{R}$ is a continuously differentiable function of x . We use the notation $(\partial^2 f / \partial x_i \partial x_j)(x)$ to indicate the ij th partial derivative of $\partial f / \partial x_j$ at a point $x \in \mathbb{R}^n$. The *Hessian* of f is the matrix whose ij th entry is equal to $(\partial^2 f / \partial x_i \partial x_j)(x)$, and is denoted by $\nabla^2 f(x)$. We have $(\partial^2 f / \partial x_i \partial x_j)(x) = (\partial^2 f / \partial x_j \partial x_i)(x)$ for every x , which implies that $\nabla^2 f(x)$ is symmetric.

If $f : \mathbb{R}^{m+n} \mapsto \mathbb{R}$ is function of (x, y) , where $x = (x_1, \dots, x_m) \in \mathbb{R}^m$ and $y = (y_1, \dots, y_n) \in \mathbb{R}^n$, we write

$$\nabla_x f(x, y) = \begin{pmatrix} \frac{\partial f(x, y)}{\partial x_1} \\ \vdots \\ \frac{\partial f(x, y)}{\partial x_m} \end{pmatrix}, \quad \nabla_y f(x, y) = \begin{pmatrix} \frac{\partial f(x, y)}{\partial y_1} \\ \vdots \\ \frac{\partial f(x, y)}{\partial y_n} \end{pmatrix},$$

$$\nabla_{xx}^2 f(x, y) = \left(\frac{\partial^2 f(x, y)}{\partial x_i \partial x_j} \right), \quad \nabla_{xy}^2 f(x, y) = \left(\frac{\partial^2 f(x, y)}{\partial x_i \partial y_j} \right),$$

$$\nabla_{yy}^2 f(x, y) = \left(\frac{\partial^2 f(x, y)}{\partial y_i \partial y_j} \right).$$

If $f : \mathbb{R}^{m+n} \mapsto \mathbb{R}^r$, $f = (f_1, f_2, \dots, f_r)$, we write

$$\nabla_x f(x, y) = (\nabla_x f_1(x, y) \cdots \nabla_x f_r(x, y)),$$

$$\nabla_y f(x, y) = (\nabla_y f_1(x, y) \cdots \nabla_y f_r(x, y)).$$

Let $f : \mathbb{R}^k \mapsto \mathbb{R}^m$ and $g : \mathbb{R}^m \mapsto \mathbb{R}^n$ be continuously differentiable functions, and let h be their composition, i.e.,

$$h(x) = g(f(x)).$$

Then, the *chain rule* for differentiation states that

$$\nabla h(x) = \nabla f(x) \nabla g(f(x)), \quad \forall x \in \mathbb{R}^k.$$

Some examples of useful relations that follow from the chain rule are:

$$\nabla(f(Ax)) = A' \nabla f(Ax), \quad \nabla^2(f(Ax)) = A' \nabla^2 f(Ax) A,$$

where A is a matrix,

$$\nabla_x(f(h(x), y)) = \nabla h(x) \nabla_h f(h(x), y),$$

$$\nabla_x(f(h(x), g(x))) = \nabla h(x) \nabla_h f(h(x), g(x)) + \nabla g(x) \nabla_g f(h(x), g(x)).$$

We now state the principal theorems relating to differentiable functions that will be useful for our purposes.

Proposition A.22: (Mean Value Theorem) If $f : \mathbb{R} \mapsto \mathbb{R}$ is continuously differentiable over an open interval I , then for every $x, y \in I$, there exists some $\xi \in [x, y]$ such that

$$f(y) - f(x) = \nabla f(\xi)(y - x).$$

Proposition A.23: (Second Order Expansions) Let $f : \mathbb{R}^n \mapsto \mathbb{R}$ be twice continuously differentiable over an open sphere S centered at a vector x .

(a) For all y such that $x + y \in S$,

$$f(x + y) = f(x) + y' \nabla f(x) + \frac{1}{2} y' \left(\int_0^1 \left(\int_0^t \nabla^2 f(x + \tau y) d\tau \right) dt \right) y.$$

(b) For all y such that $x + y \in S$, there exists an $\alpha \in [0, 1]$ such that

$$f(x + y) = f(x) + y' \nabla f(x) + \frac{1}{2} y' \nabla^2 f(x + \alpha y) y.$$

(c) For all y such that $x + y \in S$ there holds

$$f(x + y) = f(x) + y' \nabla f(x) + \frac{1}{2} y' \nabla^2 f(x) y + o(\|y\|^2).$$

Proposition A.24: (Descent Lemma) Let $f : \mathbb{R}^n \mapsto \mathbb{R}$ be continuously differentiable, and let x and y be two vectors in \mathbb{R}^n . Suppose that

$$\|\nabla f(x + ty) - \nabla f(x)\| \leq Lt\|y\|, \quad \forall t \in [0, 1],$$

where L is some scalar. Then

$$f(x + y) \leq f(x) + y' \nabla f(x) + \frac{L}{2} \|y\|^2.$$

Proof: Let t be a scalar parameter and let $g(t) = f(x + ty)$. The chain rule yields $(dg/dt)(t) = y' \nabla f(x + ty)$. Now

$$\begin{aligned} f(x+y) - f(x) &= g(1) - g(0) = \int_0^1 \frac{dg}{dt}(t) dt = \int_0^1 y' \nabla f(x + ty) dt \\ &\leq \int_0^1 y' \nabla f(x) dt + \left| \int_0^1 y' (\nabla f(x + ty) - \nabla f(x)) dt \right| \\ &\leq \int_0^1 y' \nabla f(x) dt + \int_0^1 \|y\| \cdot \|\nabla f(x + ty) - \nabla f(x)\| dt \\ &\leq y' \nabla f(x) + \|y\| \int_0^1 L t \|y\| dt \\ &= y' \nabla f(x) + \frac{L}{2} \|y\|^2. \end{aligned}$$

Q.E.D.

Proposition A.25: (Implicit Function Theorem) Let $f : \mathbb{R}^{n+m} \mapsto \mathbb{R}^m$ be a function of $x \in \mathbb{R}^n$ and $y \in \mathbb{R}^m$ such that:

- (1) $f(\bar{x}, \bar{y}) = 0$.
- (2) f is continuous, and has a continuous and nonsingular gradient matrix $\nabla_y f(x, y)$ in an open set containing (\bar{x}, \bar{y}) .

Then there exist open sets $S_{\bar{x}} \subset \mathbb{R}^n$ and $S_{\bar{y}} \subset \mathbb{R}^m$ containing \bar{x} and \bar{y} , respectively, and a continuous function $\phi : S_{\bar{x}} \mapsto S_{\bar{y}}$ such that $\bar{y} = \phi(\bar{x})$ and $f(x, \phi(x)) = 0$ for all $x \in S_{\bar{x}}$. The function ϕ is unique in the sense that if $x \in S_{\bar{x}}$, $y \in S_{\bar{y}}$, and $f(x, y) = 0$, then $y = \phi(x)$. Furthermore, if for some $p > 0$, f is p times continuously differentiable the same is true for ϕ , and we have

$$\nabla \phi(x) = -\nabla_x f(x, \phi(x)) (\nabla_y f(x, \phi(x)))^{-1}, \quad \forall x \in S_{\bar{x}}.$$

As a final word of caution to the reader, let us mention that one can easily get confused with gradient notation and its use in various formulas, such as for example the order of multiplication of various gradients in the chain rule and the implicit function theorem. Perhaps the safest guideline to minimize errors is to remember our conventions:

- (a) A vector is viewed as a column vector (an $n \times 1$ matrix).
- (b) The gradient ∇f of a scalar function $f : \mathbb{R}^n \mapsto \mathbb{R}$ is also viewed as a column vector.

- (c) The gradient matrix ∇f of a vector function $f : \mathbb{R}^n \mapsto \mathbb{R}^m$ with components f_1, \dots, f_m is the $n \times m$ matrix whose columns are the (column) vectors $\nabla f_1, \dots, \nabla f_m$.

With these rules in mind one can use “dimension matching” as an effective guide to writing correct formulas quickly.

A.6 CONTRACTION MAPPINGS

Many iterative algorithms can be written as

$$x^{k+1} = g(x^k), \quad k = 0, 1, \dots,$$

where g is a mapping from a subset X of \mathbb{R}^n into itself and has the property

$$\|g(x) - g(y)\| \leq \gamma \|x - y\|, \quad \forall x, y \in X. \quad (\text{A.4})$$

Here $\|\cdot\|$ is some norm, and γ is a scalar with $0 \leq \gamma < 1$. Such a mapping is called a *contraction mapping*, or simply a *contraction*. The scalar γ is called the *contraction modulus* of g . Note that a mapping g may be a contraction for some choice of the norm $\|\cdot\|$ and fail to be a contraction under a different choice of norm.

Let there be given a mapping $g : X \mapsto X$. Any vector $x^* \in X$ satisfying $g(x^*) = x^*$ is called a *fixed point* of g and the iteration $x^{k+1} = g(x^k)$ is an important algorithm for finding such a fixed point. The following is the central result regarding contraction mappings.

Proposition A.26: (Contraction Mapping Theorem) Suppose that $g : X \mapsto X$ is a contraction with modulus $\gamma \in [0, 1)$ and that X is a closed subset of \mathbb{R}^n . Then:

- (a) (*Existence and Uniqueness of Fixed Point*) The mapping g has a unique fixed point $x^* \in X$.
- (b) (*Convergence*) For every initial vector $x^0 \in X$, the sequence $\{x^k\}$ generated by $x^{k+1} = g(x^k)$ converges to x^* . In particular,

$$\|x^k - x^*\| \leq \gamma^k \|x^0 - x^*\|, \quad \forall k \geq 0.$$

Proof: (a) Fix some $x^0 \in X$ and consider the sequence $\{x^k\}$ generated by $x^{k+1} = g(x^k)$. We have, from inequality (A.4),

$$\|x^{k+1} - x^k\| \leq \gamma \|x^k - x^{k-1}\|,$$

for all $k \geq 1$, which implies

$$\|x^{k+1} - x^k\| \leq \gamma^k \|x^1 - x^0\|, \quad \forall k \geq 0.$$

It follows that for every $k \geq 0$ and $m \geq 1$, we have

$$\begin{aligned} \|x^{k+m} - x^k\| &\leq \sum_{i=1}^m \|x^{k+i} - x^{k+i-1}\| \\ &\leq \gamma^k (1 + \gamma + \cdots + \gamma^{m-1}) \|x^1 - x^0\| \\ &\leq \frac{\gamma^k}{1 - \gamma} \|x^1 - x^0\|. \end{aligned}$$

Therefore, $\{x^k\}$ is a Cauchy sequence and must converge to a limit x^* (Prop. A.5). Furthermore, since X is closed, x^* belongs to X . We have for all $k \geq 1$,

$$\|g(x^*) - x^*\| \leq \|g(x^*) - x^k\| + \|x^k - x^*\| \leq \gamma \|x^* - x^{k-1}\| + \|x^k - x^*\|$$

and since x^k converges to x^* , we obtain $g(x^*) = x^*$. Therefore, the limit x^* of x^k is a fixed point of g . It is a unique fixed point because if y^* were another fixed point, we would have

$$\|x^* - y^*\| = \|g(x^*) - g(y^*)\| \leq \gamma \|x^* - y^*\|,$$

which implies that $x^* = y^*$.

(b) We have

$$\|x^{k'} - x^*\| = \|g(x^{k'-1}) - g(x^*)\| \leq \gamma \|x^{k'-1} - x^*\|,$$

for all $k' \geq 1$, so by applying this relation successively for $k' = k, k-1, \dots, 1$, we obtain the desired result. **Q.E.D.**

APPENDIX B:

Convex Analysis

Convexity is a central concept in nonlinear programming. In this appendix, we collect definitions, notational conventions, and several results from the theory of convex sets and functions. A classical and extensive reference on convex analysis is Rockafellar's book [Roc70]. Related and additional material can be found in Stoer and Witzgall [StW70], Ekeland and Teman [EkT76], Rockafellar [Roc84], Hiriart-Urruty and Lemarechal [HiL93], and Rockafellar and Wets [RoW97]. A discussion of generalized notions of convexity, including quasiconvexity and pseudoconvexity, and their applications in optimization can be found in the books by Avriel [Avr76], Bazaraa, Sherali, and Shetty [BSS93], Mangasarian [Man69], and the references quoted therein.

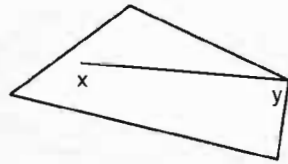
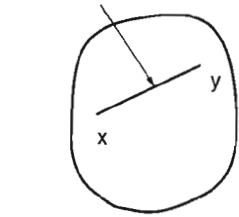
B.1 CONVEX SETS AND FUNCTIONS

The notions of a convex set and a convex function are defined below and are illustrated in Figs. B.1 and B.2, respectively.

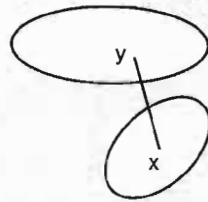
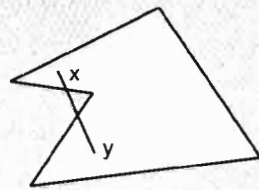
Definition B.1: Let C be a subset of \mathbb{R}^n . We say that C is *convex* if

$$\alpha x + (1 - \alpha)y \in C, \quad \forall x, y \in C, \forall \alpha \in [0, 1]. \quad (\text{B.1})$$

$$\alpha x + (1 - \alpha)y, 0 < \alpha < 1$$



Convex Sets



Nonconvex Sets

Figure B.1. Illustration of the definition of a convex set. For convexity, linear interpolation between two points in the set must yield a point within the set.

Definition B.2: Let C be a convex subset of \mathbb{R}^n . A function $f : C \mapsto \mathbb{R}$ is called *convex* if

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y), \quad \forall x, y \in C, \forall \alpha \in [0, 1]. \quad (\text{B.2})$$

The function f is called *concave* if $-f$ is convex. The function f is called *strictly convex* if the above inequality is strict for all $x, y \in C$ with $x \neq y$, and all $\alpha \in (0, 1)$. For a function $f : \mathbb{R}^n \mapsto \mathbb{R}$, we also say that f is *convex over the convex set C* if Eq. (B.2) holds.

The following proposition provides some means for verifying convexity of a set.

Proposition B.1:

- (a) For any collection $\{C_i \mid i \in I\}$ of convex sets, the set intersection $\cap_{i \in I} C_i$ is convex.
- (b) The vector sum $\{x_1 + x_2 \mid x_1 \in C_1, x_2 \in C_2\}$ of two convex sets C_1 and C_2 is convex.

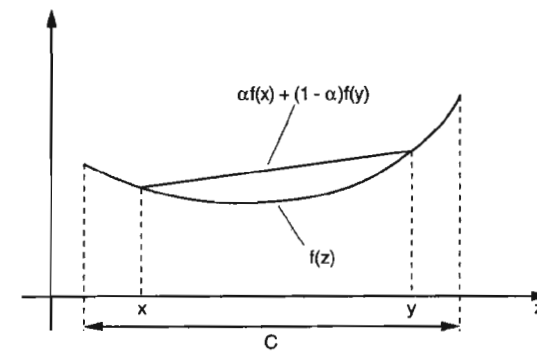


Figure B.2. Illustration of the definition of a convex function. The linear interpolation $\alpha f(x) + (1 - \alpha)f(y)$ overestimates the function value $f(\alpha x + (1 - \alpha)y)$. Note that the domain of the function must be a convex set.

- (c) The image of a convex set under a linear transformation is convex.
- (d) If C is a convex set and $f : C \mapsto \mathbb{R}$ is a convex function, the level sets $\{x \in C \mid f(x) \leq \alpha\}$ and $\{x \in C \mid f(x) < \alpha\}$ are convex for all scalars α .

Proof: The proof is straightforward using the definitions (B.1) and (B.2). For example, to prove part (a), we take two points x and y from $\cap_{i \in I} C_i$, and we use the convexity of C_i to argue that the line segment connecting x and y belongs to all the sets C_i , and hence, to their intersection. The proofs of parts (b)-(d) are similar and are left as exercises for the reader. **Q.E.D.**

We occasionally deal with convex functions that can take the value of infinity. A function $f : C \mapsto (-\infty, \infty]$, where C is a convex subset of \mathbb{R}^n , is also called convex if condition (B.2) holds. (Here the rules of arithmetic are extended to include $\infty + \infty = \infty$, $0 \cdot \infty = 0$, and $\alpha \cdot \infty = \infty$ if $\alpha > 0$.) The *effective domain* of f is the convex set

$$\text{dom}(f) = \{x \in C \mid f(x) < \infty\}.$$

By restricting the definition of a convex function to its effective domain we can avoid calculations with ∞ , and we will often do this. However, in some analyses it is more economical to use convex functions that can take the value of infinity.

The *epigraph* of a function $f : C \mapsto (-\infty, \infty]$, where C is a convex subset of \mathbb{R}^n , is the subset of \mathbb{R}^{n+1} given by

$$\text{epi}(f) = \{(x, w) \mid x \in C, w \in \mathbb{R}, f(x) \leq w\}.$$

It can be seen that f is convex if and only if $\text{epi}(f)$ is a convex set. This is a useful property, since it allows us to translate results about convex sets into results about convex functions. Another useful property, obtained by repeated application of inequality (B.2), is that if $x_1, \dots, x_m \in C$, $\alpha_1, \dots, \alpha_m \geq 0$, and $\sum_{i=1}^m \alpha_i = 1$, then

$$f\left(\sum_{i=1}^m \alpha_i x_i\right) \leq \sum_{i=1}^m \alpha_i f(x_i). \quad (\text{B.3})$$

This is a special case of *Jensen's inequality* and can be used to prove a number of interesting inequalities in applied mathematics and probability theory.

The following proposition provides some means for recognizing convex functions.

Proposition B.2:

- (a) A linear function is convex.
- (b) Any vector norm is convex.
- (c) The weighted sum of convex functions, with positive weights, is convex.
- (d) If I is an index set, $C \subset \mathbb{R}^n$ is a convex set, and $f_i : C \mapsto \mathbb{R}$ is convex for each $i \in I$, then the function $h : C \mapsto (-\infty, \infty]$ defined by

$$h(x) = \sup_{i \in I} f_i(x)$$

is also convex.

Proof: Parts (a) and (c) are immediate consequences of the definition of convexity.

Let $\|\cdot\|$ be a vector norm. For any $x, y \in \mathbb{R}^n$ and any $\alpha \in [0, 1]$, we have

$$\|\alpha x + (1 - \alpha)y\| \leq \|\alpha x\| + \|(1 - \alpha)y\| = \alpha\|x\| + (1 - \alpha)\|y\|,$$

which proves part (b).

For part (d), let us fix some $x, y \in C$, $\alpha \in [0, 1]$, and let $z = \alpha x + (1 - \alpha)y$. For every $i \in I$, we have

$$f_i(z) \leq \alpha f_i(x) + (1 - \alpha)f_i(y) \leq \alpha h(x) + (1 - \alpha)h(y).$$

Taking the supremum over all $i \in I$, we conclude that $h(z) \leq \alpha h(x) + (1 - \alpha)h(y)$, so h is convex. **Q.E.D.**

Characterizations of Differentiable Convex Functions

For differentiable functions, there is an alternative characterization of convexity, given in the following proposition and illustrated in Fig. B.3.

Proposition B.3: Let $C \subset \mathbb{R}^n$ be a convex set and let $f : \mathbb{R}^n \mapsto \mathbb{R}$ be differentiable over C .

- (a) The function f is convex over C if and only if

$$f(z) \geq f(x) + (z - x)' \nabla f(x), \quad \forall x, z \in C. \quad (\text{B.4})$$

- (b) If the inequality (B.4) is strict whenever $x \neq z$, then f is strictly convex over C .

Proof: (a) Suppose that f is convex. Let $x \in C$ and $z \in C$. By the convexity of C , we obtain $x + \alpha(z - x) \in C$ for every $\alpha \in [0, 1]$. Furthermore,

$$\lim_{\alpha \downarrow 0} \frac{f(x + \alpha(z - x)) - f(x)}{\alpha} = (z - x)' \nabla f(x). \quad (\text{B.5})$$

Using the convexity of f , we have

$$f(x + \alpha(z - x)) \leq \alpha f(z) + (1 - \alpha)f(x), \quad \forall \alpha \in [0, 1],$$

from which

$$\frac{f(x + \alpha(z - x)) - f(x)}{\alpha} \leq f(z) - f(x), \quad \forall \alpha \in [0, 1].$$

Taking the limit as $\alpha \downarrow 0$ and using Eq. (B.5), we obtain Eq. (B.4).

For the proof of the converse, suppose that inequality (B.4) is true. We fix some $x, y \in C$ and some $\alpha \in [0, 1]$. Let $z = \alpha x + (1 - \alpha)y$. Using inequality (B.4) twice, we obtain

$$f(x) \geq f(z) + (x - z)' \nabla f(z),$$

$$f(y) \geq f(z) + (y - z)' \nabla f(z).$$

We multiply the first inequality by α , the second by $(1 - \alpha)$, and add them to obtain

$$\alpha f(x) + (1 - \alpha)f(y) \geq f(z) + (\alpha x + (1 - \alpha)y - z)' \nabla f(z) = f(z),$$

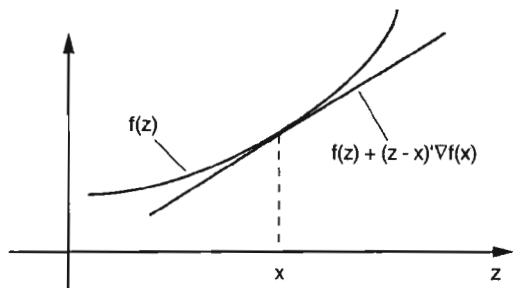


Figure B.3. Characterization of convexity in terms of first derivatives. The condition $f(z) \geq f(x) + (z-x)' \nabla f(x)$ states that a linear approximation, based on the first order Taylor series expansion, underestimates a convex function.

which proves that f is convex.

(b) The proof for the strictly convex case is almost identical to the proof of the corresponding statement of part (a) and is left for the reader. **Q.E.D.**

For twice differentiable convex functions, there is another characterization of convexity as shown by the following proposition.

Proposition B.4: Let $C \subset \mathbb{R}^n$ be a convex set, let $f : \mathbb{R}^n \mapsto \mathbb{R}$ be twice continuously differentiable over C , and let Q be a real symmetric $n \times n$ matrix.

- (a) If $\nabla^2 f(x)$ is positive semidefinite for all $x \in C$, then f is convex over C .
- (b) If $\nabla^2 f(x)$ is positive definite for every $x \in C$, then f is strictly convex over C .
- (c) If $C = \mathbb{R}^n$ and f is convex, then $\nabla^2 f(x)$ is positive semidefinite for all $x \in C$.
- (d) The quadratic function $f(x) = x'Qx$, where Q is a symmetric matrix, is convex if and only if Q is positive semidefinite. Furthermore, f is strictly convex if and only if Q is positive definite.

Proof: (a) By Prop. A.23(b) of Appendix A, for all $x, y \in C$ we have

$$f(y) = f(x) + (y-x)' \nabla f(x) + \frac{1}{2} (y-x)' \nabla^2 f(x + \alpha(y-x)) (y-x)$$

for some $\alpha \in [0, 1]$. Therefore, using the positive semidefiniteness of $\nabla^2 f$, we obtain

$$f(y) \geq f(x) + (y-x)' \nabla f(x), \quad \forall x, y \in C.$$

From Prop. B.3(a), we conclude that f is convex.

(b) Similar to the proof of part (a), we have $f(y) > f(x) + (y-x)' \nabla f(x)$ for all $x, y \in C$ with $x \neq y$, and the result follows from Prop. B.3(b).

(c) Suppose that $f : \mathbb{R}^n \mapsto \mathbb{R}$ is convex and suppose, to derive a contradiction, that there exist some $x \in \mathbb{R}^n$ and some $z \in \mathbb{R}^n$ such that $z' \nabla^2 f(x) z < 0$. Using the continuity of $\nabla^2 f$, we see that we can choose the magnitude of z to be small enough so that $z' \nabla^2 f(x + \alpha z) z < 0$ for every $\alpha \in [0, 1]$. Then, using again Prop. A.23(b) of Appendix A, we obtain $f(x+z) < f(x) + z' \nabla f(x)$, which, in view of Prop. B.3(a), contradicts the convexity of f .

(d) An easy calculation shows that $\nabla^2 f(x) = 2Q$ for all $x \in \mathbb{R}^n$. Hence, from parts (a) and (c), we obtain that f is convex if and only if Q is positive semidefinite.

If Q is positive definite, then strict convexity of f follows from part (b). For the converse, suppose that f is strictly convex. Then part (c) implies that Q is positive semidefinite and it remains to show that Q is actually positive definite. In view of Prop. A.20(b) of Appendix A, it suffices to show that zero is not an eigenvalue of Q . Suppose the contrary. Then there exists some $x \neq 0$ such that $Qx = 0$. It follows that

$$\frac{1}{2} (f(x) + f(-x)) = 0 = f(0),$$

which contradicts the strict convexity of f . **Q.E.D.**

The conclusion of Prop. B.4(c) also holds if C is assumed to have nonempty interior instead of being equal to \mathbb{R}^n ; see Exercise B.1.2. The following proposition considers a strengthened form of strict convexity characterized by the following equation:

$$(\nabla f(x) - \nabla f(y))' (x - y) \geq \alpha \|x - y\|^2, \quad \forall x, y \in \mathbb{R}^n, \quad (\text{B.6})$$

Convex functions with this property are called *strongly convex*.

Proposition B.5: (Strong Convexity) Let $f : \mathbb{R}^n \mapsto \mathbb{R}$ be continuously differentiable and let α be a positive scalar. If f is strongly convex then f is strictly convex. Furthermore, if f is twice continuously differentiable, then strong convexity of f is equivalent to the positive semidefiniteness of $\nabla^2 f(x) - \alpha I$ for every $x \in \mathbb{R}^n$, where I is the identity matrix.

Proof: Fix some $x, y \in \mathbb{R}^n$ such that $x \neq y$, and define the function $h : \mathbb{R} \mapsto \mathbb{R}$ by $h(t) = f(x + t(y-x))$. Consider some $t, t' \in \mathbb{R}$ such that

$t < t'$. Using the chain rule and Eq. (B.6), we have

$$\begin{aligned} & \left(\frac{dh}{dt}(t') - \frac{dh}{dt}(t) \right) (t' - t) \\ &= \left(\nabla f(x + t'(y - x)) - \nabla f(x + t(y - x)) \right)' (y - x) (t' - t) \\ &\geq \alpha (t' - t)^2 \|x - y\|^2 > 0. \end{aligned}$$

Thus, dh/dt is strictly increasing and for any $t \in (0, 1)$, we have

$$\frac{h(t) - h(0)}{t} = \frac{1}{t} \int_0^t \frac{dh}{d\tau}(\tau) d\tau < \frac{1}{1-t} \int_t^1 \frac{dh}{d\tau}(\tau) d\tau = \frac{h(1) - h(t)}{1-t}.$$

Equivalently, $th(1) + (1-t)h(0) > h(t)$. The definition of h yields $tf(y) + (1-t)f(x) > f(ty + (1-t)x)$. Since this inequality has been proved for arbitrary $t \in (0, 1)$ and $x \neq y$, we conclude that f is strictly convex.

Suppose now that f is twice continuously differentiable and Eq. (B.6) holds. Let c be a scalar. We use Prop. A.23(b) of Appendix A twice to obtain

$$f(x + cy) = f(x) + cy' \nabla f(x) + \frac{c^2}{2} y' \nabla^2 f(x + tcy) y,$$

and

$$f(x) = f(x + cy) - cy' \nabla f(x + cy) + \frac{c^2}{2} y' \nabla^2 f(x + scy) y,$$

for some t and s belonging to $[0, 1]$. Adding these two equations and using Eq. (B.6), we obtain

$$\frac{c^2}{2} y' (\nabla^2 f(x + scy) + \nabla^2 f(x + tcy)) y = (\nabla f(x + cy) - \nabla f(x))' (cy) \geq \alpha c^2 \|y\|^2.$$

We divide both sides by c^2 and then take the limit as $c \rightarrow 0$ to conclude that $y' \nabla^2 f(x) y \geq \alpha \|y\|^2$. Since this inequality is valid for every $y \in \mathbb{R}^n$, it follows that $\nabla^2 f(x) - \alpha I$ is positive semidefinite.

For the converse, assume that $\nabla^2 f(x) - \alpha I$ is positive semidefinite for all $x \in \mathbb{R}^n$. Consider the function $g : \mathbb{R} \mapsto \mathbb{R}$ defined by

$$g(t) = \nabla f(tx + (1-t)y)' (x - y).$$

Using the mean value theorem (Prop. A.22 in Appendix A), we have $(\nabla f(x) - \nabla f(y))' (x - y) = g(1) - g(0) = (dg/dt)(t)$ for some $t \in [0, 1]$. The result follows because

$$\frac{dg}{dt}(t) = (x - y)' \nabla^2 f(tx + (1-t)y) (x - y) \geq \alpha \|x - y\|^2,$$

where the last inequality is a consequence of the positive semidefiniteness of $\nabla^2 f(tx + (1-t)y) - \alpha I$. **Q.E.D.**

Convex and Affine Hulls

Let X be a subset of \mathbb{R}^n . A *convex combination* of elements of X , is a vector of the form $\sum_{i=1}^m \alpha_i x_i$, where x_1, \dots, x_m belong to X and $\alpha_1, \dots, \alpha_m$ are scalars such that

$$\alpha_i \geq 0, \quad i = 1, \dots, m, \quad \sum_{i=1}^m \alpha_i = 1.$$

The *convex hull* of X , denoted $\text{conv}(X)$, is the set of all convex combinations of elements of X . In particular, if X consists of a finite number of vectors x_1, \dots, x_m , its convex hull is

$$\text{conv}(\{x_1, \dots, x_m\}) = \left\{ \sum_{i=1}^m \alpha_i x_i \mid \alpha_i \geq 0, i = 1, \dots, m, \sum_{i=1}^m \alpha_i = 1 \right\}.$$

It is straightforward to verify that $\text{conv}(X)$ is a convex set, and using this, to assert that $\text{conv}(X)$ is the intersection of all convex sets containing X .

We recall that a linear manifold M is a set of the form $x + S = \{z \mid z - x \in S\}$, where S is a subspace, called the *subspace parallel to M* . If S is a subset of \mathbb{R}^n , the *affine hull* of S , denoted $\text{aff}(S)$, is the intersection of all linear manifolds containing S . Note that $\text{aff}(S)$ is itself a linear manifold and that it contains $\text{conv}(S)$. It can be seen that the affine hull of S and the affine hull of $\text{conv}(S)$ coincide.

The following is a fundamental characterization of convex sets.

Proposition B.6: (Caratheodory's Theorem) Let X be a subset of \mathbb{R}^n . Every element of $\text{conv}(X)$ can be represented as a convex combination of no more than $n + 1$ elements of X .

Proof: Let $x \in \text{conv}(X)$. Then, we can represent x as $\sum_{i=1}^m \alpha_i x_i$ for some vectors $x_i \in X$ and scalars $\alpha_i \geq 0$ with $\sum_{i=1}^m \alpha_i = 1$. Let us assume that m is the minimal number of vectors for which such a representation of x is possible; in particular, this implies that $\alpha_i > 0$ for all i . Suppose, in order to arrive at a contradiction, that $m > n + 1$, and let S be the subspace parallel to $\text{aff}(X)$. The $m - 1$ vectors $x_2 - x_1, \dots, x_m - x_1$ belong to S , and since $m - 1 > n$, they must be linearly dependent. Therefore, there exist scalars $\lambda_2, \dots, \lambda_m$ at least one of which is positive, such that

$$\sum_{i=2}^m \lambda_i (x_i - x_1) = 0.$$

Letting $\mu_i = \lambda_i$ for $i = 2, \dots, m$ and $\mu_1 = -\sum_{i=2}^m \lambda_i$, we see that

$$\sum_{i=1}^m \mu_i x_i = 0, \quad \sum_{i=1}^m \mu_i = 0,$$

while at least one of the scalars μ_2, \dots, μ_m is positive. Define

$$\bar{\alpha}_i = \alpha_i - \bar{\gamma} \mu_i, \quad i = 1, \dots, m,$$

where $\bar{\gamma} > 0$ is the largest γ such that $\alpha_i - \gamma \mu_i \geq 0$ for all i . Then, since $\sum_{i=1}^m \mu_i x_i = 0$, we see that x is also represented as $\sum_{i=1}^m \bar{\alpha}_i x_i$. Furthermore, in view of the fact $\sum_{i=1}^m \mu_i = 0$ and the choice of $\bar{\gamma}$, the coefficients $\bar{\alpha}_i$ are nonnegative, sum to one, and at least one of them is zero. Thus, x can be represented as a convex combination of fewer than m vectors of X , contradicting our earlier assumption. **Q.E.D.**

Closure and Continuity Properties

We now explore some generic topological properties of convex sets and functions.

Let C be a convex subset of \mathbb{R}^n . We say that x is a *relative interior point* of C , if $x \in C$ and there exists a neighborhood N of x such that $N \cap \text{aff}(C) \subset C$, that is, if x is an interior point of C relative to $\text{aff}(C)$. The *relative interior* of C , denoted $\text{ri}(C)$, is the set of all relative interior points of C . For example, if C is a line segment connecting two distinct points in the plane, then $\text{ri}(C)$ consists of all points of C except for the end points.

Proposition B.7:

- (a) (*Nonemptiness of Relative Interior*) If C is a nonempty convex set, $\text{ri}(C)$ is nonempty and has the same affine hull as C .
- (b) (*Line Segment Principle*) If C is a convex set, $x \in \text{ri}(C)$ and $\bar{x} \in C$, then all points on the line segment connecting x and \bar{x} , except possibly \bar{x} , belong to $\text{ri}(C)$, i.e., $\alpha x + (1 - \alpha)\bar{x} \in \text{ri}(C)$ for all $\alpha \in (0, 1]$.

Proof: (a) By using a transformation argument if necessary, we assume without loss of generality that $0 \in C$. Then, the affine hull of C , $\text{aff}(C)$, is a subspace with dimension denoted by m . If $m = 0$, then C and $\text{aff}(C)$ consist of a single point, which satisfies the definition of a relative interior point. If $m > 0$, we can find m linearly independent vectors x_1, \dots, x_m from C ; otherwise there would exist a set of $r < m$ linearly independent

vectors from C , whose span contains C , contradicting the fact that the dimension of $\text{aff}(C)$ is m . Thus x_1, \dots, x_m form a basis for the subspace $\text{aff}(C)$. It can be seen that the set

$$S = \left\{ x \mid x = \sum_{i=1}^m \alpha_i x_i, \sum_{i=1}^m \alpha_i < 1, \alpha_i > 0, i = 1, \dots, m \right\}$$

is open relative to $\text{aff}(C)$; that is, if $x \in S$, there exists an open set N_x such that $x \in N_x$ and $N_x \cap \text{aff}(C) \subset S$. [To see this, note that S is the image of the open subset of \mathbb{R}^m

$$\left\{ (\alpha_1, \dots, \alpha_m) \mid \sum_{i=1}^m \alpha_i < 1, \alpha_i > 0, i = 1, \dots, m \right\}$$

under the invertible linear transformation from \mathbb{R}^m onto $\text{aff}(C)$ that maps $(\alpha_1, \dots, \alpha_m)$ into $\sum_{i=1}^m \alpha_i x_i$; openness of sets is preserved by invertible linear transformations.] Since $S \subset C$, it follows that all points of S are relative interior points of C .

(b) See Fig. B.4. **Q.E.D.**

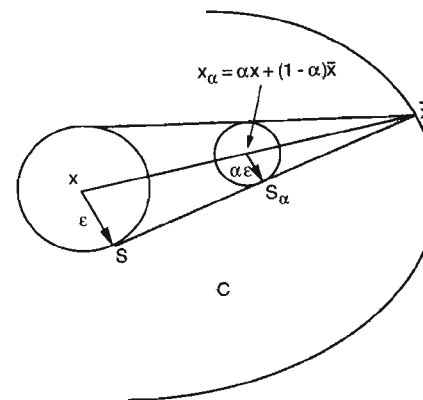


Figure B.4. Proof of the line segment principle. Since $x \in \text{ri}(C)$, there exists a sphere $S = \{z \mid \|z - x\| < \epsilon\}$ such that $S \cap \text{aff}(C) \subset C$. For all $\alpha \in (0, 1]$, let $x_\alpha = \alpha x + (1 - \alpha)\bar{x}$ and let $S_\alpha = \{z \mid \|z - x_\alpha\| < \alpha\epsilon\}$. It can be seen that each point of $S_\alpha \cap \text{aff}(C)$ is a convex combination of \bar{x} and some point of $S \cap \text{aff}(C)$. Therefore, $S_\alpha \cap \text{aff}(C) \subset C$, implying that $x_\alpha \in \text{ri}(C)$.

The closure of a set $X \subset \mathbb{R}^n$, denoted $\text{cl}(X)$, is the set of all limit points of sequences from X . It is not generally true that the closedness of

convex sets is preserved by taking vector sums, applying linear transformations, or forming convex hulls (for examples, see the subsequent Fig. B.8). We have, however, the following:

Proposition B.8:

- (a) The closure $\text{cl}(C)$ and the relative interior $\text{ri}(C)$ of a convex set C are convex sets.
- (b) The vector sum of two closed convex sets at least one of which is compact, is a closed convex set.
- (c) The image of a convex and compact set under a linear transformation is a convex and compact set.
- (d) The convex hull of a compact set is compact.

Proof: (a) Let $S_\epsilon = \{x + z \mid x \in C, \|x - z\| \leq \epsilon\}$. Then $\text{cl}(C) = \bigcap_{\epsilon > 0} S_\epsilon$, and since each set S_ϵ can be seen to be convex, the same is true of $\text{cl}(C)$. The convexity of $\text{ri}(C)$ follows from the line segment principle [Prop. B.7(b)].

(b) Let C_1 and C_2 be closed convex sets and suppose that C_2 is compact. Their vector sum $C = \{x_1 + x_2 \mid x_1 \in C_1, x_2 \in C_2\}$ is convex by Prop. B.1(b). To show that C is also closed, consider a convergent sequence $\{x_1^k + x_2^k\} \subset C$ with $\{x_1^k\} \subset C_1$ and $\{x_2^k\} \subset C_2$. Then $\{x_2^k\}$ is bounded, since C_2 is compact, and since $\{x_1^k + x_2^k\}$ converges, it follows that $\{x_1^k\}$ is also bounded. Thus, $\{(x_1^k, x_2^k)\}$ is bounded, and $\{(x_1^k, x_2^k)\}$ has a limit point $(\tilde{x}_1, \tilde{x}_2)$ with $\tilde{x}_1 \in C_1$ and $\tilde{x}_2 \in C_2$, since C_1 and C_2 are closed. The vector $\tilde{x}_1 + \tilde{x}_2$, which is the limit of $\{x_1^k + x_2^k\}$, must therefore belong to C , proving that C is closed.

(c) Let C be a convex and compact set, A be a matrix, and $\{Ax_k\}$ be a sequence with $\{x_k\} \subset C$. Then, $\{x_k\}$ has a convergent subsequence $\{x_k\}_K$ and the subsequence $\{Ax_k\}_K$ is also convergent. Therefore, the image of C under A is compact. It is also convex by Prop. B.1(c).

(d) Let X be a compact subset of \mathbb{R}^n . By Caratheodory's theorem (Prop. B.6), a sequence in $\text{conv}(X)$ can be expressed as $\{\sum_{i=1}^{n+1} \alpha_i^k x_i^k\}$, where for all k and i , $\alpha_i^k \geq 0$, $x_i^k \in X$, and $\sum_{i=1}^{n+1} \alpha_i^k = 1$. Since the sequence

$$\{(\alpha_1^k, \dots, \alpha_{n+1}^k, x_1^k, \dots, x_{n+1}^k)\}$$

belongs to a compact set, it has a limit point $\{(\alpha_1, \dots, \alpha_{n+1}, x_1, \dots, x_{n+1})\}$ such that $\sum_{i=1}^{n+1} \alpha_i = 1$, and for all i , $\alpha_i \geq 0$, and $x_i \in X$. Thus, the vector $\sum_{i=1}^{n+1} \alpha_i x_i$, which belongs to $\text{conv}(X)$, is a limit point of the sequence $\{\sum_{i=1}^{n+1} \alpha_i^k x_i^k\}$, showing that $\text{conv}(X)$ is compact. **Q.E.D.**

The following result will also be very useful to us.

Proposition B.9:

- (a) If $f : \mathbb{R}^n \mapsto \mathbb{R}$ is convex, then it is continuous. More generally, if $C \subset \mathbb{R}^n$ is convex and $f : C \mapsto \mathbb{R}$ is convex, then f is continuous in the relative interior of C .
- (b) Let X^* be the set of minimizing points of a convex function $f : \mathbb{R}^n \mapsto \mathbb{R}$ over a closed convex set X , and assume that X^* is nonempty and bounded. Then the level set

$$L_a = \{x \in X \mid f(x) \leq a\}$$

is compact for each scalar a .

Proof: (a) Restricting attention to the affine hull of C and using a transformation argument if necessary, we assume without loss of generality, that the origin is an interior point of C and that the unit cube $S = \{z \mid \|z\|_\infty \leq 1\}$ is contained in C . Let $e_i, i = 1, \dots, 2^n$, be the corners of S , that is, each e_i is a vector whose entries belong to $\{-1, 1\}$. It is not difficult to see that any $x \in S$ can be expressed in the form $x = \sum_{i=1}^{2^n} a_i e_i$, where each a_i is a nonnegative scalar and $\sum_{i=1}^{2^n} a_i = 1$. Let $A = \max_i f(e_i)$. From Jensen's inequality [Eq. (B.3)], it follows that $f(x) \leq A$ for every $x \in S$.

Let $\{x_k\}$ be a sequence in \mathbb{R}^n that converges to zero. For the purpose of proving continuity at zero, we can assume that $x_k \in S$ for all k . Using the definition of a convex function [Eq. (B.2)], we have

$$f(x_k) \leq (1 - \|x_k\|_\infty) f(0) + \|x_k\|_\infty f\left(\frac{x_k}{\|x_k\|_\infty}\right).$$

Letting k tend to infinity, $\|x_k\|_\infty$ goes to zero and we obtain

$$\limsup_{k \rightarrow \infty} f(x_k) \leq f(0) + A \limsup_{k \rightarrow \infty} \|x_k\|_\infty = f(0).$$

Inequality (B.2) also implies that

$$f(0) \leq \frac{\|x_k\|_\infty}{\|x_k\|_\infty + 1} f\left(\frac{-x_k}{\|x_k\|_\infty}\right) + \frac{1}{\|x_k\|_\infty + 1} f(x_k)$$

and letting k tend to infinity, we obtain $f(0) \leq \liminf_{k \rightarrow \infty} f(x_k)$. Thus, $\lim_{k \rightarrow \infty} f(x_k) = f(0)$ and f is continuous at zero.

(b) If X is bounded, then using also the continuity of f , which was proved in part (a), it follows that L_a is compact. We thus assume that X is

unbounded. Fix some $x^* \in X^*$ and let $b \in \mathbb{R}$ be such that $x \notin X^*$ for all x with $\|x - x^*\| = b$ (there exists such a b because X^* is bounded). Denote

$$S_b = \{x \in X \mid \|x - x^*\| = b\}, \quad \tilde{f} = \inf_{x \in S_b} f(x).$$

Since X is unbounded, closed, and convex, it is seen that S_b is nonempty and compact, and since f is continuous, it follows from Weierstrass' theorem (Prop. A.8 in Appendix A) that the infimum above is attained at some point of S_b and we have

$$\tilde{f} > f(x^*).$$

For each $x \in X$ with $\|x - x^*\| > b$, let

$$\hat{\alpha} = \frac{b}{\|x - x^*\|}, \quad \hat{x} = (1 - \hat{\alpha})x^* + \hat{\alpha}x.$$

By convexity of X , we have $\hat{x} \in X$, and by convexity of f , we have

$$(1 - \hat{\alpha})f(x^*) + \hat{\alpha}f(x) \geq f(\hat{x}).$$

Since $\|\hat{x} - x^*\| = \hat{\alpha}\|x - x^*\| = b$, we also have $\hat{x} \in S_b$, so that

$$f(\hat{x}) \geq \tilde{f}.$$

Combining these two relations and using the definition of $\hat{\alpha}$, we obtain

$$\begin{aligned} f(x) &\geq f(x^*) + \frac{f(\hat{x}) - f(x^*)}{\hat{\alpha}} \\ &\geq f(x^*) + \frac{\tilde{f} - f(x^*)}{\hat{\alpha}} \\ &= f(x^*) + \frac{\tilde{f} - f(x^*)}{b} \|x - x^*\|. \end{aligned}$$

Since $\tilde{f} > f(x^*)$, we see that if $x \in X$ and $f(x) \leq a$, then

$$\|x - x^*\| \leq \max \left\{ b, \frac{b(a - f(x^*))}{\tilde{f} - f(x^*)} \right\}.$$

Hence the level set L_a is bounded and it is also closed by continuity of f . **Q.E.D.**

Another way to phrase Prop. B.9(b) is that *if one level set of a convex function $f : \mathbb{R}^n \mapsto \mathbb{R}$ is compact, all level sets are compact.*

Local and Global Minima

Let $X \subset \mathbb{R}^n$ and let $f : X \mapsto \mathbb{R}$ be a function. A vector $x \in X$ is called a *local minimum* of f if there exists some $\epsilon > 0$ such that $f(x) \leq f(y)$ for every $y \in X$ satisfying $\|x - y\| \leq \epsilon$, where $\|\cdot\|$ is some vector norm. A vector $x \in X$ is called a *global minimum* of f if $f(x) \leq f(y)$ for every $y \in X$. A local or global maximum is defined similarly (compare also with Section 1.1).

Under convexity assumptions, the distinction between local and global minima is unnecessary as shown by the following proposition.

Proposition B.10: If $C \subset \mathbb{R}^n$ is a convex set and $f : C \mapsto \mathbb{R}$ is a convex function, then a local minimum of f is also a global minimum. If in addition f is strictly convex, then there exists at most one global minimum of f .

Proof: Suppose that x is a local minimum of f but not a global minimum. Then there exists some $y \neq x$ such that $f(y) < f(x)$. Using inequality (B.2), we conclude that $f(\alpha x + (1 - \alpha)y) < f(x)$ for every $\alpha \in [0, 1]$. This contradicts the assumption that x is a local minimum.

Suppose that f is strictly convex, and two distinct global minima x and y exist. Then their average $(x + y)/2$ must belong to C , since C is convex, and the value of f must be smaller at the average than at x and y by the strict convexity of f . Since x and y are global minima, we obtain a contradiction. **Q.E.D.**

The Projection Theorem

We close this section with a basic result of analysis and optimization, which will also be used later in this appendix.

Proposition B.11: (Projection Theorem) Let C be a closed convex set and let $\|\cdot\|$ be the Euclidean norm.

- (a) For every $x \in \mathbb{R}^n$, there exists a unique vector $z \in C$ that minimizes $\|z - x\|$ over all $z \in C$. This vector is called the *projection of x on C* , and is denoted by $[x]^+$, i.e.,

$$[x]^+ = \arg \min_{z \in C} \|z - x\|.$$

(b) Given some $x \in \mathbb{R}^n$, a vector $z \in C$ is equal to $[x]^+$ if and only if

$$(y - z)'(x - z) \leq 0, \quad \forall y \in C.$$

(c) The mapping $f : \mathbb{R}^n \mapsto C$ defined by $f(x) = [x]^+$ is continuous and nonexpansive, i.e.,

$$\|[x]^+ - [y]^+\| \leq \|x - y\|, \quad \forall x, y \in \mathbb{R}^n.$$

Proof: (a) Fix x and let w be some element of C . Minimizing $\|x - z\|$ over all $z \in C$ is equivalent to minimizing the same function over all $z \in C$ such that $\|x - z\| \leq \|x - w\|$, which is a compact set. Furthermore, the function g defined by $g(z) = \|z - x\|^2$ is continuous. Existence of a minimizing vector follows by Weierstrass' theorem (Prop. A.8 in Appendix A).

To prove uniqueness, notice that the square of the Euclidean norm is a strictly convex function of its argument [Prop. B.4(d)]. Therefore, g is strictly convex and it follows that its minimum is attained at a unique point (Prop. B.10).

(b) For all y and z in C we have

$$\|y - x\|^2 = \|y - z\|^2 + \|z - x\|^2 - 2(y - z)'(x - z) \geq \|z - x\|^2 - 2(y - z)'(x - z).$$

Therefore, if z is such that $(y - z)'(x - z) \leq 0$ for all $y \in C$, we have $\|y - x\|^2 \geq \|z - x\|^2$ for all $y \in C$, implying that $z = [x]^+$.

Conversely, let $z = [x]^+$, consider any $y \in C$, and for $\alpha > 0$, define $y_\alpha = \alpha y + (1 - \alpha)z$. We have

$$\begin{aligned} \|x - y_\alpha\|^2 &= \|(1 - \alpha)(x - z) + \alpha(x - y)\|^2 \\ &= (1 - \alpha)^2\|x - z\|^2 + \alpha^2\|x - y\|^2 + 2(1 - \alpha)\alpha(x - z)'(x - y). \end{aligned}$$

Viewing $\|x - y_\alpha\|^2$ as a function of α , we have

$$\frac{\partial}{\partial \alpha} \{\|x - y_\alpha\|^2\} \Big|_{\alpha=0} = -2\|x - z\|^2 + 2(x - z)'(x - y) = -2(y - z)'(x - z).$$

Therefore, if $(y - z)'(x - z) > 0$ for some $y \in C$, then

$$\frac{\partial}{\partial \alpha} \{\|x - y_\alpha\|^2\} \Big|_{\alpha=0} < 0$$

and for positive but small enough α , we obtain $\|x - y_\alpha\| < \|x - z\|$. This contradicts the fact $z = [x]^+$ and shows that $(y - z)'(x - z) \leq 0$ for all $y \in C$.

(c) Let x and y be elements of \mathbb{R}^n . From part (b), we have $(w - [x]^+)'(x - [x]^+) \leq 0$ for all $w \in C$. Since $[y]^+ \in C$, we obtain

$$([y]^+ - [x]^+)'(x - [x]^+) \leq 0.$$

Similarly,

$$([x]^+ - [y]^+)'(y - [y]^+) \leq 0.$$

Adding these two inequalities, we obtain

$$([y]^+ - [x]^+)'(x - [x]^+ - y + [y]^+) \leq 0.$$

By rearranging and by using the Schwartz inequality, we have

$$\|[y]^+ - [x]^+\|^2 \leq ([y]^+ - [x]^+)'(y - x) \leq \|[y]^+ - [x]^+\| \cdot \|y - x\|,$$

showing that $[\cdot]^+$ is nonexpansive and *a fortiori* continuous. **Q.E.D.**

Figure B.5 illustrates the necessary and sufficient condition of part (b) of the projection theorem.

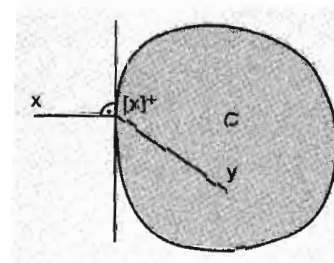


Figure B.5. Illustration of the condition satisfied by the projection $[x]^+$. For each vector $y \in C$, the vectors $x - [x]^+$ and $y - [x]^+$ form an angle larger than or equal to 90 degrees or, equivalently, $(y - [x]^+)'(x - [x]^+) \leq 0$.

EXERCISES

B.1.1

Let g be a convex, monotonically nondecreasing function of a single variable [i.e., $g(y) \leq g(\bar{y})$ for $y < \bar{y}$], and let f be a convex function defined on a convex set $C \subset \mathbb{R}^n$. Show that the function h defined by

$$h(x) = g(f(x))$$

is convex on C . Use this fact to show that the function $h(x) = e^{\beta x' Q x}$, where β is a positive scalar and Q is a positive semidefinite symmetric matrix, is convex over \mathbb{R}^n .

B.1.2

Use the line segment principle and the method of proof of Prop. B.4(c) to show that if C is a convex set with nonempty interior, and $f : \mathbb{R}^n \mapsto \mathbb{R}$ is twice continuously differentiable over C with $\nabla^2 f(x)$ positive semidefinite for all $x \in C$, then f is convex over C .

B.1.3 (Arithmetic-Geometric Mean Inequality)

Show that if $\alpha_1, \dots, \alpha_n$ are positive scalars with $\sum_{i=1}^n \alpha_i = 1$, then for every set of positive scalars x_1, \dots, x_n , we have

$$x_1^{\alpha_1} x_2^{\alpha_2} \cdots x_n^{\alpha_n} \leq \alpha_1 x_1 + \alpha_2 x_2 + \cdots + \alpha_n x_n,$$

with equality if and only if $x_1 = x_2 = \cdots = x_n$. *Hint:* Show that $-\ln x$ is a strictly convex decreasing function on $(0, \infty)$.

B.1.4

Use the result of Exercise B.1.3 to verify Young's inequality

$$xy \leq \frac{x^p}{p} + \frac{y^q}{q},$$

where $p > 0$, $q > 0$, $1/p + 1/q = 1$, $x \geq 0$, and $y \geq 0$. Then, use Young's inequality to verify Holder's inequality

$$\sum_{i=1}^n |x_i y_i| \leq \left(\sum_{i=1}^n |x_i|^p \right)^{1/p} \left(\sum_{i=1}^n |y_i|^q \right)^{1/q}.$$

B.1.5

Let $f : \mathbb{R}^{n+m} \mapsto \mathbb{R}$ be a convex function. Consider the function $F : \mathbb{R}^n \mapsto \mathbb{R}$ given by

$$F(x) = \inf_{u \in U} f(x, u),$$

where U be any nonempty and convex subset of \mathbb{R}^m such that $F(x) > -\infty$ for all $x \in \mathbb{R}^n$. Show that F is convex. *Hint:* There cannot exist $\alpha \in [0, 1]$, $x_1, x_2, u_1 \in U, u_2 \in U$ such that $F(\alpha x_1 + (1 - \alpha)x_2) > \alpha f(x_1, u_1) + (1 - \alpha)f(x_2, u_2)$.

B.1.6

Let $f : \mathbb{R}^n \mapsto \mathbb{R}$ be a differentiable function. Show that f is convex over a convex set C if and only if

$$(\nabla f(x) - \nabla f(y))'(x - y) \geq 0, \quad \forall x, y \in C.$$

Hint: The condition above says that the function f , restricted on the line segment connecting x and y , has monotonically nondecreasing gradient; see also the proof of Prop. B.5.

B.1.7 (Caratheodory's Theorem for Cones)

Let X be the cone generated by a subset of vectors $S \subset \mathbb{R}^n$, i.e., the set of vectors x of the form

$$x = \sum_{i \in I} a_i x_i,$$

where I is a finite index set, and for all $i \in I$, $x_i \in S$ and a_i is a nonnegative scalar. Show that any nonzero vector from X can be represented as a positive combination of no more than n vectors from S . Furthermore, these vectors can be chosen to be linearly independent. *Hint:* Let x be a nonzero vector from X , and let m be the smallest integer such that x has the form $\sum_{i=1}^m a_i x_i$, where $a_i > 0$ and $x_i \in S$ for all $i = 1, \dots, m$. If the vectors x_i were linearly dependent, there would exist scalars $\lambda_1, \dots, \lambda_m$, at least one of which is positive, and such that $\sum_{i=1}^m \lambda_i x_i = 0$. Consider the linear combination $\sum_{i=1}^m (a_i - \bar{\gamma} \lambda_i) x_i$, where $\bar{\gamma}$ is the largest γ such that $a_i - \gamma \lambda_i \geq 0$ for all i , to arrive at a contradiction.

B.1.8 (Properties of Relative Interiors) [Roc70]

(a) If C is a convex set in \mathbb{R}^n , then:

$$(i) \text{ cl}(C) = \text{cl}(\text{ri}(C)).$$

$$(ii) \text{ ri}(C) = \text{ri}(\text{cl}(C)).$$

$$(iii) \text{ ri}(A \cdot C) = A \cdot \text{ri}(C) \text{ for all } m \times n \text{ matrices } A.$$

(b) If C_1 and C_2 are convex sets in \mathbb{R}^n , then:

$$(i) \text{ ri}(C_1 + C_2) = \text{ri}(C_1) + \text{ri}(C_2).$$

$$(ii) \text{ ri}(C_1 \cap C_2) = \text{ri}(C_1) \cap \text{ri}(C_2), \text{ provided the sets } \text{ri}(C_1) \text{ and } \text{ri}(C_2) \text{ have a nonempty intersection.}$$

(c) If C_1 and C_2 are convex subsets of \mathbb{R}^n and \mathbb{R}^m , respectively, then

$$\text{ri}(C_1 \times C_2) = \text{ri}(C_1) \times \text{ri}(C_2).$$

B.2 SEPARATING HYPERPLANES

A *hyperplane* is a set of the form $\{x \mid a'x = b\}$, where $a \in \mathbb{R}^n$, $a \neq 0$, and $b \in \mathbb{R}$, as illustrated in Fig. B.6. An equivalent definition is that a hyperplane in \mathbb{R}^n is a linear manifold of dimension $n - 1$. The vector a called the *normal* vector of the hyperplane (it is orthogonal to the difference $x - y$ of any two vectors x and y of the hyperplane). The two sets

$$\{x \mid a'x \geq b\}, \quad \{x \mid a'x \leq b\},$$

are called the *halfspaces* associated with the hyperplane (also referred to as the *positive* and *negative halfspaces*, respectively). We have the following

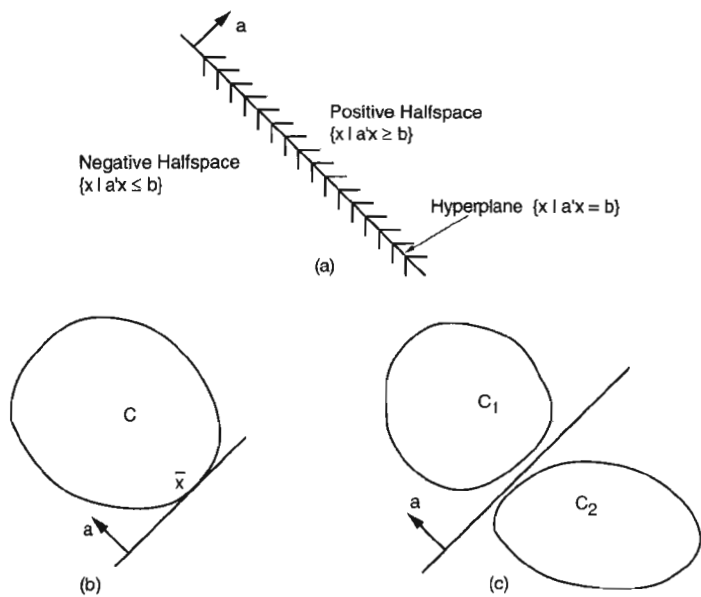


Figure B.6. (a) A hyperplane $\{x \mid a'x = b\}$ divides the space in two halfspaces as illustrated. (b) Geometric interpretation of the supporting hyperplane theorem. (c) Geometric interpretation of the separating hyperplane theorem.

result, which is also illustrated in Fig. B.6. The proof is based on the projection theorem and is illustrated in Fig. B.7.

Proposition B.12: (Supporting Hyperplane Theorem) If $C \subset \mathbb{R}^n$ is a convex set and \bar{x} is a point that does not belong to the interior of C , there exists a vector $a \neq 0$ such that

$$a'x \geq a'\bar{x}, \quad \forall x \in C. \quad (\text{B.7})$$

Proof: Denote by \bar{C} the closure of C , which is a convex set by Prop. B.8. Let $\{x_k\}$ be a sequence of vectors not belonging to \bar{C} , which converges to \bar{x} ; such a sequence exists because \bar{x} does not belong to the interior of C . If \hat{x}_k is the projection of x_k on \bar{C} , we have by part (b) of the projection theorem (Prop. B.11)

$$(\hat{x}_k - x_k)'(x - \hat{x}_k) \geq 0, \quad \forall x \in \bar{C}.$$

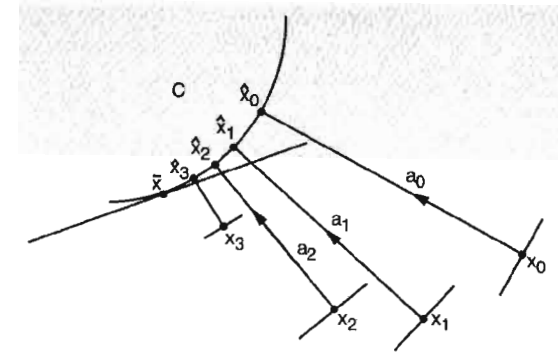


Figure B.7. Illustration of the proof of the supporting hyperplane theorem for the case where the vector \bar{x} belongs to the closure of C . We choose a sequence $\{x_k\}$ of vectors not belonging to the closure of C which converges to \bar{x} , and we project x_k on the closure of C . We then consider, for each k , the hyperplane that is orthogonal to the line segment connecting x_k and its projection, and passes through x_k . These hyperplanes “converge” to a hyperplane that supports C at \bar{x} .

Hence we obtain for all k and $x \in \bar{C}$,

$$(\hat{x}_k - x_k)'x \geq (\hat{x}_k - x_k)'\hat{x}_k = (\hat{x}_k - x_k)'(\hat{x}_k - x_k) + (\hat{x}_k - x_k)'x_k \geq (\hat{x}_k - x_k)'x_k.$$

We can write this inequality as

$$a'_k x \geq a'_k x_k, \quad \forall x \in \bar{C}, \quad k = 0, 1, \dots, \quad (\text{B.8})$$

where

$$a_k = \frac{\hat{x}_k - x_k}{\|\hat{x}_k - x_k\|}.$$

We have $\|a_k\| = 1$ for all k , and hence the sequence $\{a_k\}$ has a subsequence that converges to a nonzero limit a . By considering Eq. (B.8) for all a_k belonging to this subsequence and by taking the limit as $k \rightarrow \infty$, we obtain Eq. (B.7). **Q.E.D.**

Proposition B.13: (Separating Hyperplane Theorem) If C_1 and C_2 are two nonempty and disjoint convex subsets of \mathbb{R}^n , there exists a hyperplane that separates them, i.e., a vector $a \neq 0$ such that

$$a'x_1 \leq a'x_2, \quad \forall x_1 \in C_1, \quad x_2 \in C_2. \quad (\text{B.9})$$

Proof: Consider the convex set

$$C = \{x \mid x = x_2 - x_1, x_1 \in C_1, x_2 \in C_2\}.$$

Since C_1 and C_2 are disjoint, the origin does not belong to C , so by the supporting hyperplane theorem there exists a vector $a \neq 0$ such that

$$0 \leq a'x, \quad \forall x \in C,$$

which is equivalent to Eq. (B.9). **Q.E.D.**

Proposition B.14: (Strict Separation Theorem) If C_1 and C_2 are two nonempty and disjoint convex sets such that C_1 is closed and C_2 is compact, there exists a hyperplane that strictly separates them, i.e., a vector $a \neq 0$ and a scalar b such that

$$a'x_1 < b < a'x_2, \quad \forall x_1 \in C_1, x_2 \in C_2. \quad (\text{B.10})$$

Proof: Consider the problem

$$\begin{aligned} &\text{minimize } \|x_1 - x_2\| \\ &\text{subject to } x_1 \in C_1, x_2 \in C_2. \end{aligned} \quad (\text{B.11})$$

The set

$$C = \{x_1 - x_2 \mid x_1 \in C_1, x_2 \in C_2\}$$

is convex and closed by Prop. B.8(b). Since problem (B.11) is the problem of projecting the origin on C , we conclude using Prop. B.11(a), that problem (B.11) has at least one solution (\bar{x}_1, \bar{x}_2) . Let

$$a = \frac{\bar{x}_2 - \bar{x}_1}{2}, \quad \bar{x} = \frac{\bar{x}_1 + \bar{x}_2}{2}, \quad b = a'\bar{x}.$$

Then, $a \neq 0$, since $\bar{x}_1 \in C_1$, $\bar{x}_2 \in C_2$, and C_1 and C_2 are disjoint. The hyperplane

$$\{x \mid a'x = b\}$$

contains \bar{x} , and it can be seen from problem (B.11) that \bar{x}_1 is the projection of \bar{x} on C_1 , and \bar{x}_2 is the projection of \bar{x} on C_2 (see Fig. B.8). By Prop. B.11(b), we have

$$(\bar{x} - \bar{x}_1)'(x_1 - \bar{x}_1) \leq 0, \quad \forall x_1 \in C_1$$

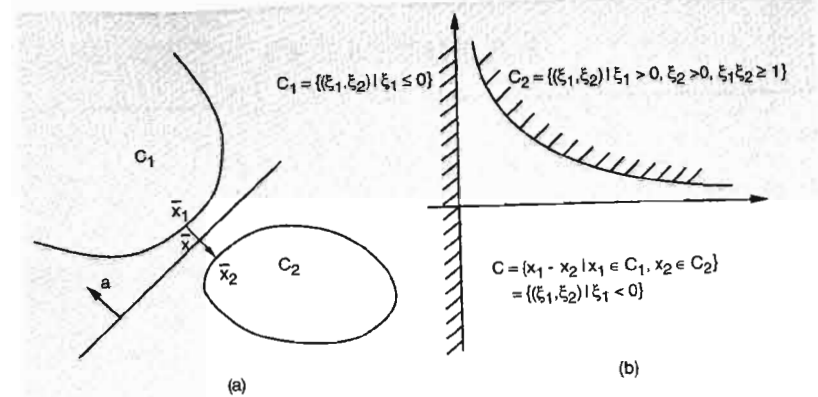


Figure B.8. (a) Illustration of the construction of a strictly separating hyperplane of two disjoint closed convex sets C_1 and C_2 one of which is also bounded (cf. Prop. B.14). (b) An example showing that if none of the two sets is compact, there may not exist a strictly separating hyperplane. This is due to the fact that the set $C = \{x_1 - x_2 \mid x_1 \in C_1, x_2 \in C_2\}$ is equal to $\{(\xi_1, \xi_2) \mid \xi_1 < 0\}$ and is not closed, even though C_1 and C_2 are closed. This is also an example where vector sum as well as linear transformation of closed convex sets does not preserve closure.

or equivalently, since $\bar{x} - \bar{x}_1 = a$,

$$a'x_1 \leq a'\bar{x}_1 = a'\bar{x} + a'(\bar{x}_1 - \bar{x}) = b - \|a\|^2 < b, \quad \forall x_1 \in C_1.$$

Thus, the left-hand side of Eq. (B.10) is proved. The right-hand side is proved similarly. **Q.E.D.**

The preceding proposition may be used to provide a fundamental characterization of closed convex sets.

Proposition B.15: Every closed convex subset of \mathbb{R}^n is the intersection of the halfspaces that contain it.

Proof: Let C be the set at issue. Clearly C is contained in the intersection of the halfspaces that contain C , so we focus on proving the reverse inclusion. Let $x \notin C$. Applying the strict separation theorem (Prop. B.14) to the sets C and $\{x\}$, we see that there exists a halfspace containing C but not containing x . Hence, if $x \notin C$, then x cannot belong to the intersection of the halfspaces containing C , proving the result. **Q.E.D.**

EXERCISES

B.2.1

Let C_1 and C_2 be two nonempty, convex sets, which are at positive Euclidean distance from each other, that is,

$$\inf_{x_1 \in C_1, x_2 \in C_2} \|x_1 - x_2\| > 0.$$

Show that there exists a hyperplane that strictly separates them. *Hint:* Adapt the proof of Prop. B.14.

B.3 CONES AND POLYHEDRAL CONVEXITY

We now develop some basic results regarding cones and also discuss the geometry of polyhedral sets. A set $C \subset \mathbb{R}^n$ is said to be a *cone* if $ax \in C$ for all $a \geq 0$ and $x \in C$. We introduce three important types of cones.

Given a cone C , the cone given by

$$C^\perp = \{y \mid y'x \leq 0, \forall x \in C\},$$

is called the *polar cone* of C .

A cone C is said to be *finitely generated*, if it has the form

$$C = \left\{ x \mid x = \sum_{j=1}^r \mu_j a_j, \mu_j \geq 0, j = 1, \dots, r \right\},$$

where a_1, \dots, a_r are some vectors.

A cone C is said to be *polyhedral*, if it has the form

$$C = \{x \mid a'_j x \leq 0, j = 1, \dots, r\},$$

where a_1, \dots, a_r are some vectors.

Figure B.9 illustrates the above definitions. It is straightforward to show that the polar cone of any cone, as well as all finitely generated and polyhedral cones are convex, by verifying the definition of convexity of Eq. (B.1). Furthermore, polar and polyhedral cones are closed, since they are intersections of closed halfspaces. Finitely generated cones are also closed

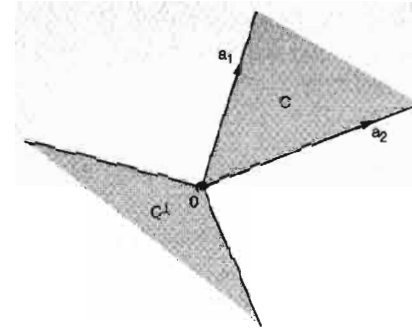


Figure B.9. Illustration of a cone and its polar in \mathbb{R}^2 . Here, a_1 and a_2 are given vectors, $C = \{x \mid x = \mu_1 a_1 + \mu_2 a_2, \mu_1 \geq 0, \mu_2 \geq 0\}$, which is a finitely generated cone, and $C^\perp = \{y \mid y'a_1 \leq 0, y'a_2 \leq 0\}$, which is a polyhedral cone.

as shown in part (b) of the following proposition, which also provides some additional important results.

Proposition B.16:

- (a) (*Polar Cone Theorem*) For any nonempty closed convex cone C , we have $(C^\perp)^\perp = C$.
- (b) Let a_1, \dots, a_r be vectors of \mathbb{R}^n . Then the finitely generated cone

$$C = \left\{ x \mid x = \sum_{j=1}^r \mu_j a_j, \mu_j \geq 0, j = 1, \dots, r \right\} \quad (\text{B.12})$$

is closed and its polar cone is the polyhedral cone given by

$$C^\perp = \{x \mid x'a_j \leq 0, j = 1, \dots, r\}. \quad (\text{B.13})$$

- (c) (*Minkowski - Weyl Theorem*) A cone is polyhedral if and only if it is finitely generated.
- (d) (*Farkas' Lemma*) Let x, e_1, \dots, e_m , and a_1, \dots, a_r be vectors of \mathbb{R}^n . We have $x'y \leq 0$ for all vectors $y \in \mathbb{R}^n$ such that

$$y'e_i = 0, \quad \forall i = 1, \dots, m, \quad y'a_j \leq 0, \quad \forall j = 1, \dots, r,$$

if and only if x can be expressed as

$$x = \sum_{i=1}^m \lambda_i e_i + \sum_{j=1}^r \mu_j a_j,$$

where λ_i and μ_j are some scalars with $\mu_j \geq 0$ for all j .

Proof: (a) See Fig. B.10.

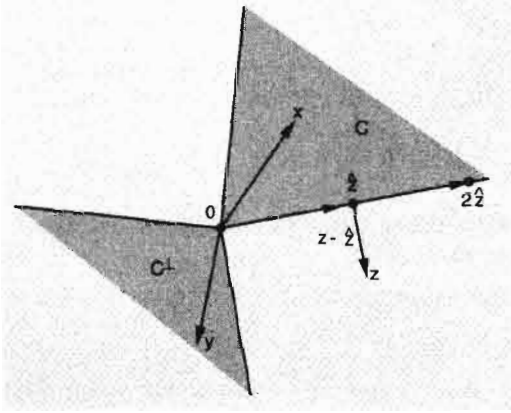


Figure B.10. Proof of the polar cone theorem. If $x \in C$, then for all $y \in C^\perp$, we have $x'y \leq 0$, which implies that $x \in (C^\perp)^\perp$. Hence, $C \subset (C^\perp)^\perp$. To prove the reverse inclusion, take $z \in (C^\perp)^\perp$, and let \hat{z} be the unique projection of z on C , as shown in the figure. Since C is closed, the projection exists by the projection theorem (Prop. B.11), which also implies that

$$(z - \hat{z})'(x - \hat{z}) \leq 0, \quad \forall x \in C.$$

By taking $x = 0$ and $x = 2\hat{z}$ in the preceding relation, it is seen that

$$(z - \hat{z})'\hat{z} = 0.$$

Combining the last two relations, we obtain $(z - \hat{z})'x \leq 0$ for all $x \in C$. Therefore, $(z - \hat{z}) \in C^\perp$, and since $z \in (C^\perp)^\perp$, we obtain $(z - \hat{z})'z \leq 0$, which when added to $(z - \hat{z})'\hat{z} = 0$ yields $\|z - \hat{z}\|^2 \leq 0$. Therefore, $z = \hat{z}$ and $z \in C$. It follows that $(C^\perp)^\perp \subset C$.

(b) We first show that the polar cone of C has the desired form (B.13). If y satisfies $y'a_j \leq 0$ for all j , then $y'x \leq 0$ for all $x \in C$, so the set in the right-hand side of Eq. (B.13) is a subset of C^\perp . Conversely, if $y \in C^\perp$, that is, if $y'x \leq 0$ for all $x \in C$, then (since a_j belong to C) we have $y'a_j \leq 0$, for all j . Thus, C^\perp is a subset of the set in the right-hand side of Eq. (B.13).

To show that C is closed, it will suffice to show that C is polyhedral; this will also prove half of the Minkowski-Weyl theorem [part (c)]. Our proof, due to [Wet90], is constructive and uses induction on the number of

vectors r . We will also give an alternative proof, which is simpler than the first but does not show simultaneously half of part (c).

To start the induction, we assume without loss of generality that $a_1 = 0$. Then, for $r = 1$, we have $C = \{0\}$, which is polyhedral, since it can be expressed as

$$\{x \mid u_i'x \leq 0, -u_i'x \leq 0, i = 1, \dots, n\},$$

where u_i is the i th unit coordinate vector.

Assume that for some $r \geq 2$, the set

$$C_{r-1} = \left\{ x \mid x = \sum_{j=1}^{r-1} \mu_j a_j, \mu_j \geq 0 \right\}$$

has a polyhedral representation

$$P_{r-1} = \{x \mid b_j'x \leq 0, j = 1, \dots, m\}.$$

Let

$$\beta_j = a_r' b_j, \quad j = 1, \dots, m,$$

and define the index sets

$$J^- = \{j \mid \beta_j < 0\}, \quad J^0 = \{j \mid \beta_j = 0\}, \quad J^+ = \{j \mid \beta_j > 0\}.$$

Let also

$$b_{l,k} = b_l - \frac{\beta_l}{\beta_k} b_k, \quad \forall l \in J^+, k \in J^-.$$

We will show that the set

$$C_r = \left\{ x \mid x = \sum_{j=1}^r \mu_j a_j, \mu_j \geq 0 \right\}$$

has the polyhedral representation

$$P_r = \{x \mid b_j'x \leq 0, j \in J^- \cup J^0, b_{l,k}'x \leq 0, l \in J^+, k \in J^-\},$$

thus completing the induction.

We have $C_r \subset P_r$ because by construction, all the vectors a_1, \dots, a_r satisfy the inequalities defining P_r . To show the reverse inclusion, we fix a vector $x \in P_r$ and we verify that there exists $\mu_r \geq 0$ such that

$$x - \mu_r a_r \in P_{r-1},$$

which is equivalent to

$$\gamma \leq \mu_r \leq \delta,$$

where

$$\gamma = \max \left\{ 0, \max_{j \in J^+} \frac{b'_j x}{\beta_j} \right\}, \quad \delta = \min_{j \in J^-} \frac{b'_j x}{\beta_j}.$$

Since $x \in P_r$, we have

$$0 \leq \frac{b'_k x}{\beta_k}, \quad \forall k \in J^-, \quad (\text{B.14})$$

and also $b'_{l,k} x \leq 0$ for all $l \in J^+$, $k \in J^-$, or equivalently

$$\frac{b'_l x}{\beta_l} \leq \frac{b'_k x}{\beta_k}, \quad \forall l \in J^+, k \in J^-. \quad (\text{B.15})$$

Equations (B.14) and (B.15) imply that $\gamma \leq \delta$, thereby completing the proof.

We now give an alternative proof that C is closed, which is based again on induction on the number of vectors r . When $r = 1$, C is either $\{0\}$ (if $a_1 = 0$) or a halfline, and is therefore closed. Suppose, for some $r \geq 1$, all cones of the form

$$\left\{ x \mid x = \sum_{j=1}^r \mu_j a_j, \mu_j \geq 0 \right\}$$

are closed. Then, we will show that a cone of the form

$$C_{r+1} = \left\{ x \mid x = \sum_{j=1}^{r+1} \mu_j a_j, \mu_j \geq 0 \right\}$$

is also closed. Without loss of generality, assume that $\|a_j\| = 1$ for all j . There are two cases: (i) The vectors $-a_1, \dots, -a_{r+1}$ belong to C_{r+1} , in which case C_{r+1} is the subspace spanned by a_1, \dots, a_{r+1} and is therefore closed, and (ii) The negative of one of the vectors, say $-a_{r+1}$, does not belong to C_{r+1} . In this case, consider the cone

$$C_r = \left\{ x \mid x = \sum_{j=1}^r \mu_j a_j, \mu_j \geq 0 \right\},$$

which is closed by the induction hypothesis. Let

$$m = \min_{x \in C_r, \|x\|=1} a'_{r+1} x.$$

Since, the set $\{x \in C_r \mid \|x\| = 1\}$ is nonempty and compact, the minimum above is attained at some x^* by Weierstrass' theorem. We have, using the Schwartz inequality,

$$m = a'_{r+1} x^* \geq -\|a_{r+1}\| \cdot \|x^*\| = -1,$$

with equality if and only if $x^* = -a_{r+1}$. It follows that

$$m > -1,$$

since otherwise we would have $x^* = -a_{r+1}$, which violates the hypothesis $(-a_{r+1}) \notin C_r$. Let $\{x_k\}$ be a convergent sequence in C_{r+1} . We will prove that its limit belongs to C_{r+1} , thereby showing that C_{r+1} is closed. Indeed, for all k , we have $x_k = \xi_k a_{r+1} + y_k$, where $\xi_k \geq 0$ and $y_k \in C_r$. Using the fact $\|a_{r+1}\| = 1$, we obtain

$$\begin{aligned} \|x_k\|^2 &= \xi_k^2 + \|y_k\|^2 + 2\xi_k a'_{r+1} y_k \\ &\geq \xi_k^2 + \|y_k\|^2 + 2m\xi_k \|y_k\| \\ &= (\xi_k - \|y_k\|)^2 + 2(1+m)\xi_k \|y_k\|. \end{aligned}$$

Since $\{x_k\}$ converges, $\xi_k \geq 0$, and $1+m > 0$, it follows that the sequences $\{\xi_k\}$ and $\{y_k\}$ are bounded and hence, they have limit points denoted by ξ and y , respectively. The limit of $\{x_k\}$ is

$$\lim_{k \rightarrow \infty} (\xi_k a_{r+1} + y_k) = \xi a_{r+1} + y,$$

which belongs to C_{r+1} , since $\xi \geq 0$ and $y \in C_r$ (by the closure hypothesis on C_r). We conclude that C_{r+1} is closed, completing the proof.

(c) We have already shown in the proof of part (b) that a finitely generated cone is polyhedral. To show the reverse, we use parts (a) and (b) to conclude that the polar of any polyhedral cone [cf. Eq. (B.13)] is finitely generated [cf. Eq. (B.12)]. The finitely generated cone (B.12) has already been shown to be polyhedral, so its polar, which is the "typical" polyhedral cone (B.13), is finitely generated. This completes the proof.

(d) Define $a_{r+i} = e_i$ and $a_{r+m+i} = -e_i$, $i = 1, \dots, m$. The result to be shown translates to

$$x \in C \iff x \in P^\perp,$$

where

$$\begin{aligned} C &= \left\{ x \mid x = \sum_{j=1}^{r+2m} \mu_j a_j, \mu_j \geq 0 \right\}, \\ P &= \{y \mid y' a_j \leq 0, j = 1, \dots, r+2m\}. \end{aligned}$$

Since by part (b), $P = C^\perp$ and C is closed, we have by part (a), $P^\perp = (C^\perp)^\perp = C$. **Q.E.D.**

Polyhedral Sets

A nonempty subset of \mathbb{R}^n is said to be a *polyhedral set* (or *polyhedron*) if it is of the form

$$P = \{x \mid a'_j x \leq b_j, j = 1, \dots, r\},$$

where a_j are some vectors and b_j are some scalars.

The following is a fundamental result, showing that a polyhedral set can be represented as the sum of the convex hull of a finite set of points and a finitely generated cone. The proof is based on an interesting construction that can be used to translate results about polyhedral cones to results about polyhedral sets.

Proposition B.17: A set P is polyhedral if and only if there exist a nonempty and finite set of vectors $\{v_1, \dots, v_m\}$, and a finitely generated cone C such that

$$P = \left\{ x \mid x = y + \sum_{j=1}^m \mu_j v_j, y \in C, \sum_{j=1}^m \mu_j = 1, \mu_j \geq 0, j = 1, \dots, m \right\}.$$

Proof: Assume that P is polyhedral. Then, it has the form

$$P = \{x \mid a'_j x \leq b_j, j = 1, \dots, r\},$$

for some vectors a_j and some scalars b_j . Consider the polyhedral cone of \mathbb{R}^{n+1}

$$\hat{P} = \{(x, w) \mid 0 \leq w, a'_j x \leq b_j w, j = 1, \dots, r\}$$

and note that

$$P = \{x \mid (x, 1) \in \hat{P}\}.$$

By the Minkowski – Weyl theorem [Prop. B.16(c)], \hat{P} is finitely generated, so it has the form

$$\hat{P} = \left\{ (x, w) \mid x = \sum_{j=1}^m \mu_j v_j, w = \sum_{j=1}^m \mu_j d_j, \mu_j \geq 0, j = 1, \dots, m \right\},$$

for some vectors v_j and scalars d_j . Since $w \geq 0$ for all vectors $(x, w) \in \hat{P}$, we see that $d_j \geq 0$ for all j . Let

$$J^+ = \{j \mid d_j > 0\}, \quad J^0 = \{j \mid d_j = 0\}.$$

By replacing μ_j by μ_j/d_j for all $j \in J^+$, we obtain the equivalent description

$$\hat{P} = \left\{ (x, w) \mid x = \sum_{j=1}^m \mu_j v_j, w = \sum_{j \in J^+} \mu_j, \mu_j \geq 0, j = 1, \dots, m \right\}.$$

Since $P = \{x \mid (x, 1) \in \hat{P}\}$, we obtain

$$P = \left\{ x \mid x = \sum_{j \in J^+} \mu_j v_j + \sum_{j \in J^0} \mu_j v_j, \sum_{j \in J^+} \mu_j = 1, \mu_j \geq 0, j = 1, \dots, m \right\}.$$

Thus, P is the vector sum of the convex hull of the vectors $v_j, j \in J^+$, plus the finitely generated cone

$$\left\{ \sum_{j \in J^0} \mu_j v_j \mid \mu_j \geq 0, j \in J^0 \right\}.$$

To prove that the vector sum of the convex hull of a finite set of points with a finitely generated cone is a polyhedral set, we use a reverse argument; we pass to a finitely generated cone description, we use the Minkowski – Weyl theorem to assert that this cone is polyhedral, and we finally construct a polyhedral set description. The details are left as an exercise for the reader. **Q.E.D.**

B.4 EXTREME POINTS

A vector x is said to be an *extreme point* of a convex set C if x belongs to C and there do not exist vectors $y \in C$ and $z \in C$, with $y \neq x$ and $z \neq x$, and a scalar $\alpha \in (0, 1)$ such that $x = \alpha y + (1 - \alpha)z$. An equivalent definition is that x cannot be expressed as a convex combination of some vectors of C , all of which are different from x .

An important fact that forms the basis for the simplex method of linear programming, is that if a linear function f attains a minimum over a polyhedral set C having at least one extreme point, then f attains a minimum at some extreme point of C (as well as possibly at some other nonextreme points). We will prove this fact after considering the more general case where f is concave and C is closed and convex. We first show a preliminary result.

Proposition B.18: Let C be a nonempty, closed, convex set in \mathbb{R}^n .

- (a) If H is a hyperplane that passes through a boundary point of C and contains C in one of its halfspaces, then every extreme point of $T = C \cap H$ is also an extreme point of C .
- (b) C has at least one extreme point if and only if it does not contain a line, that is, a set L of the form $L = \{x + \alpha d \mid \alpha \in \mathbb{R}\}$ with $d \neq 0$.

Proof: (a) Let \bar{x} be an element of T which is not an extreme point of C . Then we have $\bar{x} = \alpha y + (1 - \alpha)z$ for some $\alpha \in (0, 1)$, and some $y \in C$ and $z \in C$, with $y \neq \bar{x}$ and $z \neq \bar{x}$. Since $\bar{x} \in H$, \bar{x} is a boundary point of C , and the halfspace containing C is of the form $\{x \mid a'x \geq a'\bar{x}\}$, where $a \neq 0$. Then $a'y \geq a'\bar{x}$ and $a'z \geq a'\bar{x}$, which in view of $\bar{x} = \alpha y + (1 - \alpha)z$, implies that $a'y = a'\bar{x}$ and $a'z = a'\bar{x}$. Therefore, $y \in T$ and $z \in T$, showing that \bar{x} cannot be an extreme point of T .

(b) Assume that C has an extreme point x and contains a line $L = \{\bar{x} + \alpha d \mid \alpha \in \mathbb{R}\}$, where $d \neq 0$. We will arrive at a contradiction. For each integer $n > 0$, the vector

$$x_n = \left(1 - \frac{1}{n}\right)x + \frac{1}{n}(\bar{x} + nd) = x + d + \frac{1}{n}(\bar{x} - x)$$

lies in the line segment connecting x and $\bar{x} + nd$, so it belongs to C . Since C is closed, $x + d = \lim_{n \rightarrow \infty} x_n$ must also belong to C . Similarly, we show that $x - d$ must belong to C . Thus $x - d$, x , and $x + d$ all belong to C , contradicting the hypothesis that x is an extreme point.

Conversely, we use induction on the dimension of the space to show that if C does not contain a line, it must have an extreme point. This is true in the real line \mathbb{R}^1 , so assume it is true in \mathbb{R}^{n-1} . If a nonempty, closed, convex subset C of \mathbb{R}^n contains no line, it must have some boundary point \bar{x} . Take any hyperplane H passing through \bar{x} and containing C in one of its halfspaces. Then, since H is an $(n - 1)$ -dimensional manifold, the set $C \cap H$ lies in an $(n - 1)$ -dimensional space and contains no line, so by the induction hypothesis, it must have an extreme point. By part (a), this extreme point must also be an extreme point of C . **Q.E.D.**

We say that a set $C \subset \mathbb{R}^n$ is *bounded from below* if there exists a vector $b \in \mathbb{R}^n$ such that $x \geq b$ for all $x \in C$.

Proposition B.19: Let C be a closed convex set which is bounded from below and let $f : C \rightarrow \mathbb{R}$ be a concave function. Then if f attains a minimum over C , it attains a minimum at some extreme point of C .

Proof: We first show that f attains a minimum at some boundary point of C . Let x^* be a vector where f attains a minimum over C . If x^* is a boundary point we are done, so assume that x^* is an interior point of C . Let

$$L = \{x \mid x = x^* + \lambda d, \lambda \in \mathbb{R}\}$$

be a line passing through x^* , where d is a vector with strictly positive coordinates. Then, using the boundedness from below, convexity, and closure of C , we see that the set $C \cap L$ contains a set of the form

$$\{x^* + \lambda d \mid \lambda_1 \leq \lambda \leq \lambda_2\}$$

for some $\lambda_2 > 0$ and some $\lambda_1 < 0$ for which the vector

$$\bar{x} = x^* + \lambda_1 d$$

is a boundary point of C . If $f(\bar{x}) > f(x^*)$, we have by concavity of f ,

$$\begin{aligned} f(x^*) &\geq \frac{\lambda_2}{\lambda_2 - \lambda_1} f(\bar{x}) + \left(1 - \frac{\lambda_2}{\lambda_2 - \lambda_1}\right) f(x^* + \lambda_2 d) \\ &> \frac{\lambda_2}{\lambda_2 - \lambda_1} f(x^*) + \left(1 - \frac{\lambda_2}{\lambda_2 - \lambda_1}\right) f(x^* + \lambda_2 d). \end{aligned}$$

It follows that $f(x^*) > f(x^* + \lambda_2 d)$. This contradicts the optimality of x^* , proving that $f(\bar{x}) = f(x^*)$.

We have shown that the minimum of f is attained at some boundary point \bar{x} of C . If \bar{x} is an extreme point of C , we are done. If it is not an extreme point, consider a hyperplane H passing through \bar{x} and containing C in one of its halfspaces. The intersection $T_1 = C \cap H$ is closed, convex, bounded from below, and lies in a linear manifold M_1 of dimension $n - 1$. Furthermore, f attains its minimum over T_1 at \bar{x} . Thus, by the preceding argument, it also attains its minimum at some boundary point x_1 of T_1 . If x_1 is an extreme point of T_1 , then by Prop. B.18, it is also an extreme point of C and the result follows. If x_1 is not an extreme point of T_1 , then we view M_1 as a space of dimension $n - 1$ and we form T_2 , the intersection of T_1 with a hyperplane in M_1 that passes through x_1 and contains T_1 in one of its halfspaces. This hyperplane will be of dimension $n - 2$. We can continue this process for at most n times, when a set T_n consisting of a single point is obtained. This point is an extreme point of T_n and, by repeated application of Prop. B.18, an extreme point of C . **Q.E.D.**

As a corollary we have the following:

Proposition B.20: Let C be a closed convex set and let $f : C \mapsto \mathbb{R}$ be a concave function. Assume that for some invertible $n \times n$ matrix A and some $b \in \mathbb{R}^n$ we have

$$Ax \geq b, \quad \forall x \in C.$$

Then if f attains a minimum over C , it attains a minimum at some extreme point of C .

Proof: Consider the transformation $x = A^{-1}y$ and the problem of minimizing

$$h(y) = f(A^{-1}y)$$

over $Y = \{y \mid A^{-1}y \in C\}$. The function h is concave over the closed convex set Y . Furthermore, $y \geq b$ for all $y \in Y$ and hence Y is bounded from below. By Prop. B.19, h attains a minimum at some extreme point y^* of Y . Then f attains its minimum over C at $x^* = A^{-1}y^*$, while x^* is an extreme point of C , since it can be verified that invertible transformations of sets map extreme points to extreme points. **Q.E.D.**

Extreme Points of Polyhedral Sets

We now consider a polyhedral set P and we characterize the set of its extreme points (also called *vertices*). By Prop. B.17, P can be represented as

$$P = C + \hat{P},$$

where C is a finitely generated cone C and \hat{P} is the convex hull of some vectors v_1, \dots, v_m :

$$\hat{P} = \left\{ x \mid x = \sum_{j=1}^m \mu_j v_j, \sum_{j=1}^m \mu_j = 1, \mu_j \geq 0, j = 1, \dots, m \right\}.$$

We note that an extreme point \bar{x} of P cannot be of the form $\bar{x} = c + \hat{x}$, where $c \neq 0$, $c \in C$, and $\hat{x} \in \hat{P}$, since in this case \bar{x} would be the midpoint of the line segment connecting the distinct vectors \hat{x} and $2c + \hat{x}$. Therefore, an extreme point of P must belong to \hat{P} , and since $\hat{P} \subset P$, it must also be an extreme point of \hat{P} . An extreme point of \hat{P} must be one of the vectors v_1, \dots, v_m , since otherwise this point would be expressible as a convex combination of v_1, \dots, v_m . Thus the set of extreme points of P is either empty or finite. Using Prop. B.18(b), it follows that *the set of extreme points of P is nonempty and finite if and only if P contains no line.*

If P is bounded, then we must have $P = \hat{P}$, and it can be shown that P is equal to the convex hull of its extreme points (not just the convex hull of the vectors v_1, \dots, v_m). The proof is sketched in Exercise B.4.1.

The following proposition gives another and more specific characterization of extreme points of polyhedral sets, and is central in the theory of linear programming.

Proposition B.21: Let P be a polyhedral set in \mathbb{R}^n .

(a) If P has the form

$$P = \{x \mid a'_j x \leq b_j, j = 1, \dots, r\},$$

where a_j and b_j are given vectors and scalars, respectively, then a vector $v \in P$ is an extreme point of P if and only if the set

$$A_v = \{a_j \mid a'_j v = b_j, j = 1, \dots, r\}$$

contains n linearly independent vectors.

(b) If P has the form

$$P = \{x \mid Ax = b, x \geq 0\},$$

where A is a given $m \times n$ matrix and b is a given vector, then a vector $v \in P$ is an extreme point of P if and only if the columns of A corresponding to the nonzero coordinates of v are linearly independent.

(c) (*Fundamental Theorem of Linear Programming*) Assume that P has at least one extreme point. Then if a linear function attains a minimum over P , it attains a minimum at some extreme point of P .

Proof: (a) If the set A_v contains fewer than n linearly independent vectors, then the system of equations

$$a'_j w = 0, \quad \forall a_j \in A_v$$

has a nonzero solution \bar{w} . For sufficiently small $\gamma > 0$, we have $v + \gamma \bar{w} \in P$ and $v - \gamma \bar{w} \in P$, thus showing that v is not an extreme point. Thus, if v is an extreme point, A_v must contain n linearly independent vectors.

Conversely, suppose that A_v contains a subset \bar{A}_v consisting of n linearly independent vectors. Suppose that for some $y \in P$, $z \in P$, and $\alpha \in (0, 1)$, we have $v = \alpha y + (1 - \alpha)z$. Then for all $a_j \in \bar{A}_v$, we have

$$b_j = a'_j v = \alpha a'_j y + (1 - \alpha) a'_j z \leq \alpha b_j + (1 - \alpha) b_j = b_j.$$

Thus v , y , and z are all solutions of the system of n linearly independent equations

$$a'_j w = b_j, \quad \forall a_j \in \bar{A}_v.$$

Hence $v = y = z$, implying that v is an extreme point.

(b) Let k be the number of zero coordinates of v , and consider the matrix \bar{A} , which is the same as A except that the columns corresponding to the zero coordinates of v are set to zero. We write P in the form

$$P = \{x \mid Ax \leq b, -Ax \leq -b, -x \leq 0\},$$

and apply the result of part (a). We obtain that v is an extreme point if and only if \bar{A} contains $n - k$ linearly independent rows, which is equivalent to the $n - k$ nonzero columns of \bar{A} (corresponding to the nonzero coordinates of v) being linearly independent.

(c) Since P is polyhedral, it has a representation

$$P = \{x \mid Ax \geq b\},$$

for some $m \times n$ matrix A and some $b \in \mathbb{R}^m$. If A had rank less than n , then its nullspace would contain some nonzero vector \bar{x} , so P would contain a line parallel to \bar{x} , contradicting the existence of an extreme point [cf. Prop. B.18(b)]. Thus A has rank n and hence it must contain n linearly independent rows that constitute an $n \times n$ invertible submatrix \hat{A} . If \hat{b} is the corresponding subvector of b , we see that every $x \in P$ satisfies $\hat{A}x \geq \hat{b}$. The result then follows using Prop. B.20. **Q.E.D.**

EXERCISES

B.4.1

Show that a polyhedron of the form

$$P = \left\{ x \mid x = \sum_{j=1}^m \mu_j v_j, \sum_{j=1}^m \mu_j = 1, \mu_j \geq 0, j = 1, \dots, m \right\}. \quad (\text{B.16})$$

is the convex hull of its extreme points. *Hint:* Use induction on the dimension of the space. Suppose that all bounded polyhedra of $(n-1)$ -dimensional spaces have a representation of the form (B.16), but there is a bounded polyhedron $P \subset \mathbb{R}^n$

and a vector $x \in P$, which is not in the convex hull P_E of the extreme points of P . Let \hat{x} be the projection of x on P_E and let \bar{x} be a solution of the problem

$$\begin{aligned} &\text{maximize } (x - \hat{x})'z \\ &\text{subject to } z \in P. \end{aligned}$$

The polyhedron

$$\hat{P} = P \cap \{z \mid (x - \hat{x})'z = (x - \hat{x})'\bar{x}\}$$

is equal to the convex hull of its extreme points by the induction hypothesis. Show that $P_E \cap \hat{P} = \emptyset$, while, by Prop. B.18(a), each of the extreme points of \hat{P} is also an extreme point of P , arriving at a contradiction.

B.5 DIFFERENTIABILITY ISSUES

Convex functions have interesting differentiability properties, which we discuss in this section. We first consider convex functions of a single variable.

Let I be an interval of real numbers, and let $f : I \rightarrow \mathbb{R}$ be convex. If $x, y, z \in I$ and $x < y < z$, then we can show the relation

$$\frac{f(y) - f(x)}{y - x} \leq \frac{f(z) - f(x)}{z - x} \leq \frac{f(z) - f(y)}{z - y}, \quad (\text{B.17})$$

which is illustrated in Fig. B.11. For a formal proof, note that, using the definition of a convex function [cf. Eq. (B.2)], we obtain

$$f(y) \leq \left(\frac{y - x}{z - x} \right) f(z) + \left(\frac{z - y}{z - x} \right) f(x)$$

and either of the desired inequalities follows by appropriately rearranging terms.

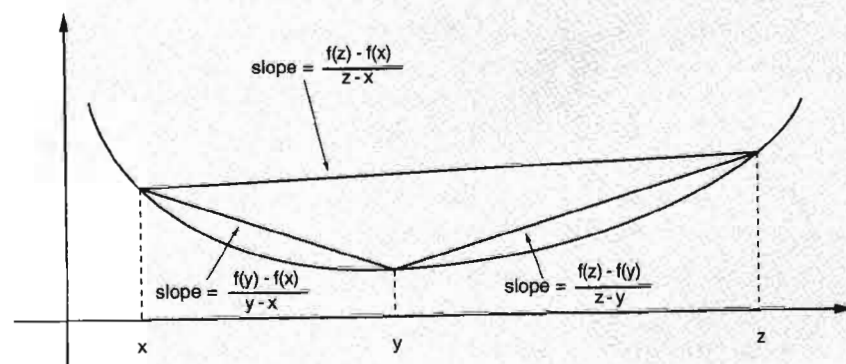


Figure B.11. Illustration of the inequalities (B.17). The rate of change of the function f is nondecreasing with its argument.

Let a and b be the infimum and the supremum, respectively, of I , also referred to as the *end points* of I . For any $x \in I$, $x \neq b$, and for any $\alpha > 0$ such that $x + \alpha \in I$, we define

$$s^+(x, \alpha) = \frac{f(x + \alpha) - f(x)}{\alpha}.$$

Let $0 < \alpha \leq \alpha'$. We use the first inequality in Eq. (B.17) with $y = x + \alpha$ and $z = x + \alpha'$ to obtain $s^+(x, \alpha) \leq s^+(x, \alpha')$. Therefore, $s^+(x, \alpha)$ is a nondecreasing function of α and, as α decreases to zero, it converges either to a finite number or to $-\infty$. Let $f^+(x)$ be the value of the limit, which we call the *right derivative* of f at the point x . Similarly, if $x \in I$, $x \neq a$, $\alpha > 0$, and $x - \alpha \in I$, we define

$$s^-(x, \alpha) = \frac{f(x) - f(x - \alpha)}{\alpha},$$

which is, by a symmetrical argument, a nonincreasing function of α . Its limit as α decreases to zero, denoted by $f^-(x)$, is called the *left derivative* of f at the point x , and is either finite or equal to ∞ .

In the case where the end points a and b belong to the domain I of f , we define for completeness $f^-(a) = -\infty$ and $f^+(b) = \infty$.

Proposition B.22: Let $I \subset \mathbb{R}$ be a convex interval and let $f : I \mapsto \mathbb{R}$ be a convex function. Let a and b be the end points of I .

- (a) We have $f^-(y) \leq f^+(y)$ for every $y \in I$.
- (b) If x belongs to the interior of I , then $f^+(x)$ and $f^-(x)$ are finite.
- (c) If $x, z \in I$ and $x < z$, then $f^+(x) \leq f^-(z)$.
- (d) The functions $f^-, f^+ : I \mapsto [-\infty, +\infty]$ are nondecreasing.
- (e) The function f^+ (respectively, f^-) is right- (respectively, left-) continuous at every interior point of I . Also, if $a \in I$ (respectively, $b \in I$) and f is continuous at a (respectively, b), then f^+ (respectively, f^-) is right- (respectively, left-) continuous at a (respectively, b).
- (f) If f is differentiable at a point x belonging to the interior of I , then $f^+(x) = f^-(x) = (df/dx)(x)$.
- (g) For any $x, z \in I$ and any d satisfying $f^-(x) \leq d \leq f^+(x)$, we have

$$f(z) \geq f(x) + d(z - x).$$

- (h) The function $f^+ : I \mapsto (-\infty, \infty]$ [respectively, $f^- : I \mapsto [-\infty, \infty)$] is upper (respectively, lower) semicontinuous at every $x \in I$.

Proof: (a) If y is an end point of I , the result is trivial because $f^-(a) = -\infty$ and $f^+(b) = \infty$. We assume that y is an interior point, we let $\alpha > 0$, and use Eq. (B.17), with $x = y - \alpha$ and $z = y + \alpha$, to obtain $s^-(y, \alpha) \leq s^+(y, \alpha)$. Taking the limit as α decreases to zero, we obtain $f^-(y) \leq f^+(y)$.

(b) Let x belong to the interior of I and let $\alpha > 0$ be such that $x - \alpha \in I$. Then $f^-(x) \geq s^-(x, \alpha) > -\infty$. For similar reasons, we obtain $f^+(x) < \infty$. Part (a) then implies that $f^-(x) < \infty$ and $f^+(x) > -\infty$.

(c) We use Eq. (B.17), with $y = (z + x)/2$, to obtain $s^+(x, (z - x)/2) \leq s^-(z, (z - x)/2)$. The result then follows because $f^+(x) \leq s^+(x, (z - x)/2)$ and $s^-(z, (z - x)/2) \leq f^-(z)$.

(d) This follows by combining parts (a) and (c).

(e) Fix some $x \in I$, $x \neq b$, and some positive δ and α such that $x + \delta + \alpha < b$. We allow x to be equal to a , in which case f is assumed to be continuous at a . We have $f^+(x + \delta) \leq s^+(x + \delta, \alpha)$. We take the limit, as δ decreases to zero, to obtain $\lim_{\delta \downarrow 0} f^+(x + \delta) \leq s^+(x, \alpha)$. We have used here the fact that $s^+(x, \alpha)$ is a continuous function of x , which is a consequence of the continuity of f (Prop. B.9). We now let α decrease to zero to obtain $\lim_{\alpha \downarrow 0} \lim_{\delta \downarrow 0} f^+(x + \delta) \leq f^+(x)$. The reverse inequality is also true because f^+ is nondecreasing and this proves the right-continuity of f^+ . The proof for f^- is similar.

(f) This is immediate from the definition of f^+ and f^- .

(g) Fix some $x, z \in I$. The result is trivially true for $x = z$. We only consider the case $x < z$; the proof for the case $x > z$ is similar. Since $s^+(x, \alpha)$ is nondecreasing in α , we have $(f(z) - f(x))/(z - x) \geq s^+(x, \alpha)$ for α belonging to $(0, z - x)$. Letting α decrease to zero, we obtain $(f(z) - f(x))/(z - x) \geq f^+(x) \geq d$ and the result follows.

(h) This follows from parts (a), (d), (e), and the definition of semicontinuity (Definition A.4 in Appendix A). **Q.E.D.**

We now consider the *directional derivative* $f'(x; y)$ of a convex function $f : \mathbb{R}^n \mapsto \mathbb{R}$ at a vector $x \in \mathbb{R}^n$ in the direction $y \in \mathbb{R}^n$. This derivative is equal to the right derivative $F_y^+(0)$ of the convex scalar function $F_y(\alpha) = f(x + \alpha y)$ at $\alpha = 0$, i.e.,

$$f'(x; y) = \lim_{\alpha \downarrow 0} \frac{f(x + \alpha y) - f(x)}{\alpha} = \lim_{\alpha \downarrow 0} \frac{F_y(\alpha) - F_y(0)}{\alpha} = F_y^+(0), \quad (\text{B.18})$$

and the limit in the above equation is guaranteed to exist. Similarly, the left derivative $F_y^-(0)$ of F_y is equal to $-f'(x; -y)$ and, by using Prop. B.22(a), we obtain $F_y^-(0) \leq F_y^+(0)$, or equivalently,

$$-f'(x; -y) \leq f'(x; y), \quad \forall y \in \mathbb{R}^n. \quad (\text{B.19})$$

The directional derivative can be used to provide a necessary and sufficient condition for optimality in the problem of minimizing a convex function $f : \mathbb{R}^n \mapsto \mathbb{R}$ over a convex set $X \subset \mathbb{R}^n$. In particular, x^* is a global minimum of f over X if and only if

$$f'(x^*; x - x^*) \geq 0, \quad \forall x \in X.$$

This follows from the definition (B.18) of directional derivative, and from the fact that the difference quotient

$$\frac{f(x^* + \alpha(x - x^*)) - f(x^*)}{\alpha}$$

is a monotonically nondecreasing function of α .

The following proposition generalizes the upper semicontinuity property of right derivatives of scalar convex functions [Prop. B.22(h)], and shows that if f is differentiable, then its gradient is continuous.

Proposition B.23: Let $f : \mathbb{R}^n \mapsto \mathbb{R}$ be convex, and let $\{f_k\}$ be a sequence of convex functions $f_k : \mathbb{R}^n \mapsto \mathbb{R}$ with the property that $\lim_{k \rightarrow \infty} f_k(x_k) = f(x)$ for every $x \in \mathbb{R}^n$ and every sequence $\{x_k\}$ that converges to x . Then for any $x \in \mathbb{R}^n$ and $y \in \mathbb{R}^n$, and any sequences $\{x_k\}$ and $\{y_k\}$ converging to x and y , respectively, we have

$$\limsup_{k \rightarrow \infty} f'_k(x_k; y_k) \leq f'(x; y). \quad (\text{B.20})$$

Furthermore, if f is differentiable at all $x \in \mathbb{R}^n$, then its gradient $\nabla f(x)$ is a continuous function of x .

Proof: For any $\mu > f'(x; y)$, there exists an $\bar{\alpha} > 0$ such that

$$\frac{f(x + \alpha y) - f(x)}{\alpha} < \mu, \quad \forall \alpha \leq \bar{\alpha}.$$

Hence, for $\alpha \leq \bar{\alpha}$, we have

$$\frac{f_k(x_k + \alpha y_k) - f_k(x_k)}{\alpha} < \mu$$

for all sufficiently large k , and using Eq. (B.18), we obtain

$$\limsup_{k \rightarrow \infty} f'_k(x_k; y_k) < \mu.$$

Since this is true for all $\mu > f'(x; y)$, inequality (B.20) follows.

If f is differentiable at all $x \in \mathbb{R}^n$, then using the continuity of f and the part of the proposition just proved, we have for every sequence $\{x_k\}$ converging to x and every $y \in \mathbb{R}^n$,

$$\limsup_{k \rightarrow \infty} \nabla f(x_k)'y = \limsup_{k \rightarrow \infty} f'(x_k; y) \leq f'(x; y) = \nabla f(x)'y.$$

By replacing y by $-y$ in the preceding argument, we obtain

$$-\liminf_{k \rightarrow \infty} \nabla f(x_k)'y = \limsup_{k \rightarrow \infty} (-\nabla f(x_k)'y) \leq -\nabla f(x)'y.$$

Therefore, we have $\nabla f(x_k)'y \rightarrow \nabla f(x)'y$ for every y , which implies that $\nabla f(x_k) \rightarrow \nabla f(x)$. Hence, the gradient is continuous. **Q.E.D.**

Subgradients and Subdifferentials

Given a convex function $f : \mathbb{R}^n \mapsto \mathbb{R}$, we say that a vector $d \in \mathbb{R}^n$ is a *subgradient* of f at a point $x \in \mathbb{R}^n$ if

$$f(z) \geq f(x) + (z - x)'d, \quad \forall z \in \mathbb{R}^n. \quad (\text{B.21})$$

If instead f is a concave function, we say that d is a subgradient of f at x if $-d$ is a subgradient of the convex function $-f$ at x . The set of all subgradients of a convex (or concave) function f at $x \in \mathbb{R}^n$ is called the *subdifferential* of f at x , and is denoted by $\partial f(x)$. Figure B.12 provides some examples of subdifferentials.

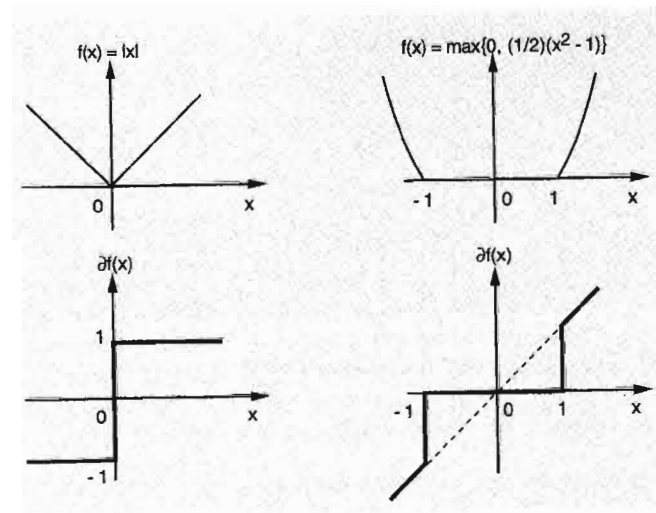


Figure B.12. The subdifferential of some scalar convex functions as a function of the argument x .

We next provide the relationship between the directional derivative and the subdifferential, and prove some basic properties of subgradients.

Proposition B.24: Let $f : \mathbb{R}^n \mapsto \mathbb{R}$ be convex. For every $x \in \mathbb{R}^n$, the following hold:

(a) A vector d is a subgradient of f at x if and only if

$$f'(x; y) \geq y'd, \quad \forall y \in \mathbb{R}^n.$$

(b) The subdifferential $\partial f(x)$ is a nonempty, convex, and compact set, and there holds

$$f'(x; y) = \max_{d \in \partial f(x)} y'd, \quad \forall y \in \mathbb{R}^n. \quad (\text{B.22})$$

In particular, f is differentiable at x with gradient $\nabla f(x)$, if and only if it has $\nabla f(x)$ as its unique subgradient at x . Furthermore, if X is a bounded set, the set $\cup_{x \in X} \partial f(x)$ is bounded.

(c) If a sequence $\{x_k\}$ converges to x and $d_k \in \partial f(x_k)$ for all k , the sequence $\{d_k\}$ is bounded and each of its limit points is a subgradient of f at x .

(d) If f is equal to the sum $f_1 + \dots + f_m$ of convex functions $f_j : \mathbb{R}^n \mapsto \mathbb{R}$, $j = 1, \dots, m$, then $\partial f(x)$ is equal to the vector sum $\partial f_1(x) + \dots + \partial f_m(x)$.

(e) If f is equal to the composition of a convex function $h : \mathbb{R}^m \mapsto \mathbb{R}$ and an $m \times n$ matrix A [$f(x) = h(Ax)$], then $\partial f(x)$ is equal to $A'\partial h(Ax) = \{A'g \mid g \in \partial h(Ax)\}$.

(f) x minimizes f over a convex set $X \subset \mathbb{R}^n$ if and only if there exists a subgradient $d \in \partial f(x)$ such that

$$d'(z - x) \geq 0, \quad \forall z \in X.$$

Proof: (a) The subgradient inequality (B.21) is equivalent to

$$\frac{f(x + \alpha y) - f(x)}{\alpha} \geq y'd, \quad \forall y \in \mathbb{R}^n, \alpha > 0.$$

Since the quotient on the left above decreases monotonically to $f'(x; y)$ as $\alpha \downarrow 0$ [Eq. (B.17)], we conclude that the subgradient inequality (B.21) is equivalent to $f'(x; y) \geq y'd$ for all $y \in \mathbb{R}^n$. Therefore we obtain

$$d \in \partial f(x) \iff f'(x; y) \geq y'd, \quad \forall y \in \mathbb{R}^n. \quad (\text{B.23})$$

(b) From Eq. (B.23), we see that $\partial f(x)$ is the intersection of the closed halfspaces $\{d \mid y'd \leq f'(x; y)\}$, where y ranges over the nonzero vectors of \mathbb{R}^n . It follows that $\partial f(x)$ is closed and convex. It is also bounded, since otherwise, for some $y \in \mathbb{R}^n$, $y'd$ could be made unbounded by proper choice of $d \in \partial f(x)$, contradicting Eq. (B.23). Since $\partial f(x)$ is both closed and bounded, it is compact.

To show that $\partial f(x)$ is nonempty and that Eq. (B.22) holds, we first observe that Eq. (B.23) implies that $f'(x; y) \geq \max_{d \in \partial f(x)} y'd$ [where the maximum is $-\infty$ if $\partial f(x)$ is empty]. To show the reverse inequality, take any x and y in \mathbb{R}^n , and consider the subset of \mathbb{R}^{n+1}

$$C_1 = \{(\mu, z) \mid \mu > f(z)\},$$

and the half-line

$$C_2 = \{(\mu, z) \mid \mu = f(x) + \alpha f'(x; y), z = x + \alpha y, \alpha \geq 0\};$$

see Fig. B.13. Using the definition of directional derivative and the convexity of f , it follows that these two sets are nonempty, convex, and disjoint. By applying the separating hyperplane theorem (Prop. B.13), we see that there exists a nonzero vector $(\gamma, w) \in \mathbb{R}^{n+1}$ such that

$$\gamma\mu + w'z \leq \gamma(f(x) + \alpha f'(x; y)) + w'(x + \alpha y), \quad \forall \alpha \geq 0, z \in \mathbb{R}^n, \mu > f(z). \quad (\text{B.24})$$

We cannot have $\gamma > 0$ since then the left-hand side above could be made arbitrarily large by choosing μ sufficiently large. Also if $\gamma = 0$, then Eq. (B.24) implies that $w = 0$, which is a contradiction. Therefore, $\gamma < 0$ and by dividing with γ in Eq. (B.24), we obtain

$$\mu + (z - x)'(w/\gamma) \geq f(x) + \alpha f'(x; y) + \alpha y'(w/\gamma), \quad \forall \alpha \geq 0, z \in \mathbb{R}^n, \mu > f(z). \quad (\text{B.25})$$

By taking the limit in the above relation as $\alpha \downarrow 0$ and $\mu \downarrow f(z)$, we obtain $f(z) \geq f(x) + (z - x)'(-w/\gamma)$ for all $z \in \mathbb{R}^n$, implying that $(-w/\gamma) \in \partial f(x)$. By taking $z = x$ and $\alpha = 1$ in Eq. (B.25), and by taking the limit as $\mu \downarrow f(x)$, we obtain $y'(-w/\gamma) \geq f'(x; y)$, which implies that $\max_{d \in \partial f(x)} y'd \geq f'(x; y)$. The proof of Eq. (B.22) is complete.

From the definition of directional derivative, we see that f is differentiable at x with gradient $\nabla f(x)$ if and only if the directional derivative $f'(x; y)$ is a linear function of the form $f'(x; y) = \nabla f(x)'y$. Thus, from Eq. (B.22), f is differentiable at x with gradient $\nabla f(x)$, if and only if it has $\nabla f(x)$ as its unique subgradient at x .

Finally, let X be a bounded set. To show that $\cup_{x \in X} \partial f(x)$ is bounded, we assume the contrary, i.e. that there exists a sequence $\{x_k\} \subset X$, and a sequence $\{d_k\}$ with $d_k \in \partial f(x_k)$ for all k and $\|d_k\| \rightarrow \infty$. Without loss of generality, we assume that $d_k \neq 0$ for all k , and we denote $y_k = d_k/\|d_k\|$.

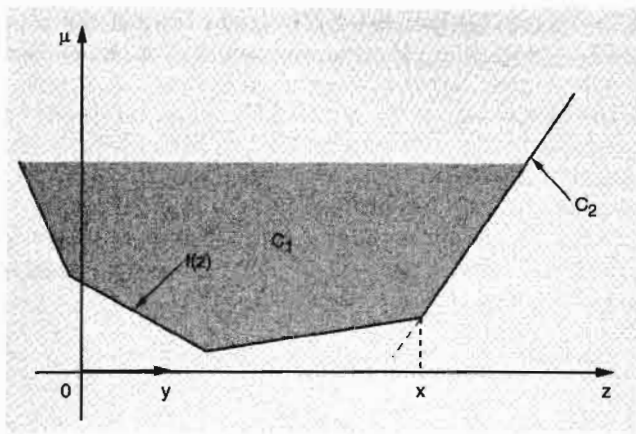


Figure B.13. Illustration of the sets C_1 and C_2 used in the hyperplane separation argument of the proof of Prop. B.24.

Since both $\{x_k\}$ and $\{y_k\}$ are bounded, they must contain convergent subsequences. We assume without loss of generality that x_k converges to some x and y_k converges to some y with $\|y\| = 1$. By Eq. (B.22), we have

$$f'(x_k; y_k) \geq d'_k y_k = \|d_k\|,$$

so it follows that $f'(x_k; y_k) \rightarrow \infty$. This contradicts, however, Eq. (B.20), which implies that $\limsup_{k \rightarrow \infty} f'(x_k; y_k) \leq f'(x; y)$.

(c) By part (b), the sequence $\{d_k\}$ is bounded, and by part (a), we have

$$y'd_k \leq f'(x_k; y), \quad \forall y \in \mathbb{R}^n.$$

If d is a limit point of $\{d_k\}$, we have by taking limit in the above relation and by using Prop. B.23

$$y'd \leq \limsup_{k \rightarrow \infty} f'(x_k; y) \leq f'(x; y), \quad \forall y \in \mathbb{R}^n.$$

Therefore, by part (a), we have $d \in \partial f(x)$.

(d) It will suffice to prove the result for the case where $f = f_1 + f_2$. If $d_1 \in \partial f_1(x)$ and $d_2 \in \partial f_2(x)$, then from the subgradient inequality (B.21), we have

$$f_1(z) \geq f_1(x) + (z - x)'d_1, \quad \forall z \in \mathbb{R}^n,$$

$$f_2(z) \geq f_2(x) + (z - x)'d_2, \quad \forall z \in \mathbb{R}^n,$$

so by adding, we obtain

$$f(z) \geq f(x) + (z - x)'(d_1 + d_2), \quad \forall z \in \mathbb{R}^n.$$

Hence $d_1 + d_2 \in \partial f(x)$, implying that $\partial f_1(x) + \partial f_2(x) \subset \partial f(x)$.

To prove the reverse inclusion, suppose to come to a contradiction, that there exists a $d \in \partial f(x)$ such that $d \notin \partial f_1(x) + \partial f_2(x)$. Since by part (b), the sets $\partial f_1(x)$ and $\partial f_2(x)$ are compact, the set $\partial f_1(x) + \partial f_2(x)$ is compact (cf. Prop. B.8), and by Prop. B.14, there exists a hyperplane strictly separating $\{d\}$ from $\partial f_1(x) + \partial f_2(x)$, i.e., a vector y and a scalar b such that

$$y'(d_1 + d_2) < b < y'd, \quad \forall d_1 \in \partial f_1(x), d_2 \in \partial f_2(x).$$

From this we obtain

$$\max_{d_1 \in \partial f_1(x)} y'd_1 + \max_{d_2 \in \partial f_2(x)} y'd_2 < y'd,$$

or using part (b),

$$f'_1(x; y) + f'_2(x; y) < y'd.$$

By using the definition of directional derivative, $f'_1(x; y) + f'_2(x; y) = f'(x; y)$, so we have

$$f'(x; y) < y'd,$$

which is a contradiction in view of part (a).

(e) It is seen using the definition of directional derivative that

$$f'(x; y) = h'(Ax; Ay), \quad \forall y \in \mathbb{R}^n.$$

Let $g \in \partial h(Ax)$ and $d = A'g$. Then by part (a), we have

$$g'z \leq h'(Ax; z) \quad \forall z \in \mathbb{R}^m,$$

and in particular,

$$g'Ay \leq h'(Ax; Ay) \quad \forall y \in \mathbb{R}^n,$$

or

$$(A'g)'y \leq f(x; y), \quad \forall y \in \mathbb{R}^n.$$

Hence, by part (a), we have $A'g \in \partial f(x)$, so that $A'\partial h(Ax) \subset \partial f(x)$.

To prove the reverse inclusion, suppose to come to a contradiction, that there exists a $d \in \partial f(x)$ such that $d \notin A'\partial h(Ax)$. Since by part (b), the set $\partial h(Ax)$ is compact, the set $A'\partial h(Ax)$ is also compact (cf. Prop. B.8), and by Prop. B.14, there exists a hyperplane strictly separating $\{d\}$ from $A'\partial h(Ax)$, i.e., a vector y and a scalar b such that

$$y'(A'g) < b < y'd, \quad \forall g \in \partial h(Ax).$$

From this we obtain

$$\max_{g \in \partial h(Ax)} (Ay)'g < y'd,$$

or using part (b),

$$h'(Ax; Ay) < y'd.$$

Since $h'(Ax; Ay) = f'(x; y)$, it follows that

$$f'(x; y) < y'd,$$

which is a contradiction in view of part (a).

(f) Suppose that for some $d \in \partial f(x)$ and all $z \in X$, we have $d'(z - x) \geq 0$. Then, since from the definition of a subgradient we have $f(z) - f(x) \geq d'(z - x)$ for all $z \in X$, we obtain $f(z) - f(x) \geq 0$ for all $z \in X$, so x minimizes f over X .

Conversely, suppose that x minimizes f over X . Consider the set of feasible directions of X at x

$$W = \{w \neq 0 \mid x + \alpha w \in X \text{ for some } \alpha > 0\},$$

and the cone

$$\hat{W} = \{d \mid d'w \geq 0, \forall w \in W\}$$

(this is equal to $-W^\perp$, the set of all d such that $-d$ belongs to the polar cone W^\perp). If $\partial f(x)$ and \hat{W} have a point in common, we are done, so to arrive at a contradiction, assume the opposite, i.e., $\partial f(x) \cap \hat{W} = \emptyset$. Since $\partial f(x)$ is compact and \hat{W} is closed, by Prop. B.14 there exists a hyperplane strictly separating $\partial f(x)$ and \hat{W} , i.e., a vector y and a scalar c such that

$$g'y < c < d'y, \quad \forall g \in \partial f(x), \forall d \in \hat{W}.$$

Using the fact that \hat{W} is a closed cone, it follows that

$$c < 0 \leq d'y, \quad \forall d \in \hat{W}, \quad (\text{B.26})$$

which when combined with the preceding inequality, also yields

$$\max_{g \in \partial f(x)} g'y < c < 0.$$

Thus, using part (b), we have $f'(x; y) < 0$, while from Eq. (B.26), we see that y belongs to the polar cone of W^\perp , which by the polar cone theorem [Prop. B.16(a)], implies that y is in the closure of the set of feasible directions W . Hence for a sequence y^k of feasible directions converging to y we have $f'(x; y^k) < 0$, and this contradicts the optimality of x . **Q.E.D.**

Note that Prop. B.24(f) generalizes the optimality condition of Chapter 2 for the case where f is differentiable: $\nabla f(x)'(z - x) \geq 0$ for all $z \in X$. In the special case where $X = \mathbb{R}^n$, we obtain a basic necessary and sufficient condition for unconstrained optimality of x :

$$0 \in \partial f(x).$$

This optimality condition is also evident from the subgradient inequality (B.21).

Danskin's Min-Max Theorem

We next consider the directional derivative and the subdifferential of the function $f(x) = \max_{z \in Z} \phi(x, z)$.

Proposition B.25: (Danskin's Theorem) Let $Z \subset \mathbb{R}^m$ be a compact set, and let $\phi : \mathbb{R}^n \times Z \mapsto \mathbb{R}$ be continuous and such that $\phi(\cdot, z) : \mathbb{R}^n \mapsto \mathbb{R}$ is convex for each $z \in Z$.

(a) The function $f : \mathbb{R}^n \mapsto \mathbb{R}$ given by

$$f(x) = \max_{z \in Z} \phi(x, z) \quad (\text{B.27})$$

is convex and has directional derivative given by

$$f'(x; y) = \max_{z \in Z(x)} \phi'(x, z; y), \quad (\text{B.28})$$

where $\phi'(x, z; y)$ is the directional derivative of the function $\phi(\cdot, z)$ at x in the direction y , and $Z(x)$ is the set of maximizing points in Eq. (B.27)

$$Z(x) = \left\{ \bar{z} \mid \phi(x, \bar{z}) = \max_{z \in Z} \phi(x, z) \right\}.$$

In particular, if $Z(x)$ consists of a unique point \bar{z} and $\phi(\cdot, \bar{z})$ is differentiable at x , then f is differentiable at x , and $\nabla f(x) = \nabla_x \phi(x, \bar{z})$, where $\nabla_x \phi(x, \bar{z})$ is the vector with coordinates

$$\frac{\partial \phi(x, \bar{z})}{\partial x_i}, \quad i = 1, \dots, n.$$

(b) If $\phi(\cdot, z)$ is differentiable for all $z \in Z$ and $\nabla_x \phi(x, \cdot)$ is continuous on Z for each x , then

$$\partial f(x) = \text{conv} \{ \nabla_x \phi(x, z) \mid z \in Z(x) \}, \quad \forall x \in \mathbb{R}^n.$$

In particular, if ϕ is linear in x for all $z \in Z$, i.e.,

$$\phi(x, z) = a'_z x + b_z, \quad \forall z \in Z,$$

then

$$\partial f(x) = \text{conv} \{ a_z \mid z \in Z(x) \}.$$

Proof: (a) The convexity of f has been established in Prop. B.2(d). We note that since ϕ is continuous and Z is compact, the set $Z(x)$ is nonempty by Weierstrass' theorem (Prop. A.8 in Appendix A) and f is finite. For any $z \in Z(x)$, $y \in \mathbb{R}^n$, and $\alpha > 0$, we use the definition of f to obtain

$$\frac{f(x + \alpha y) - f(x)}{\alpha} \geq \frac{\phi(x + \alpha y, z) - \phi(x, z)}{\alpha}.$$

Taking the limit as α decreases to zero, we obtain $f'(x; y) \geq \phi'(x, z; y)$. Since this is true for every $z \in Z(x)$, we conclude that

$$f'(x; y) \geq \sup_{z \in Z(x)} \phi'(x, z; y), \quad \forall y \in \mathbb{R}^n. \quad (\text{B.29})$$

To prove the reverse inequality and that the supremum in the right-hand side of the above inequality is attained, consider a sequence $\{\alpha_k\}$ of positive scalars that converges to zero and let $x_k = x + \alpha_k y$. For each k , let z_k be a vector in $Z(x_k)$. Since $\{z_k\}$ belongs to the compact set Z , it has a subsequence converging to some $\bar{z} \in Z$. Without loss of generality, we assume that the entire sequence $\{z_k\}$ converges to \bar{z} . We have

$$\phi(x_k, z_k) \geq \phi(x_k, \bar{z}), \quad \forall z \in Z,$$

so by taking the limit as $k \rightarrow \infty$ and by using the continuity of ϕ , we obtain

$$\phi(x, \bar{z}) \geq \phi(x, z), \quad \forall z \in Z.$$

Therefore, $\bar{z} \in Z(x)$. We now have

$$\begin{aligned} f'(x; y) &\leq \frac{f(x + \alpha_k y) - f(x)}{\alpha_k} \\ &= \frac{\phi(x + \alpha_k y, z_k) - \phi(x, \bar{z})}{\alpha_k} \\ &\leq \frac{\phi(x + \alpha_k y, z_k) - \phi(x, z_k)}{\alpha_k} \\ &\leq -\phi'(x + \alpha_k y, z_k; -y) \\ &\leq \phi'(x + \alpha_k y, z_k; y), \end{aligned} \quad (\text{B.30})$$

where the last inequality follows from inequality (B.19). We apply Prop. B.23 to the functions f_k defined by $f_k(\cdot) = \phi(\cdot, z_k)$, and with $x_k = x + \alpha_k y$, to obtain

$$\limsup_{k \rightarrow \infty} \phi'(x + \alpha_k y, z_k; y) \leq \phi'(x, \bar{z}; y). \quad (\text{B.31})$$

We take the limit in inequality (B.30) as $k \rightarrow \infty$, and we use inequality (B.31) to conclude that

$$f'(x; y) \leq \phi'(x, \bar{z}; y).$$

This relation together with inequality (B.29) proves Eq. (B.28).

For the last statement of part (a), if $Z(x)$ consists of the unique point \bar{z} , Eq. (B.28) and the differentiability assumption on ϕ yield

$$f'(x; y) = \phi'(x, \bar{z}; y) = y' \nabla_x \phi(x, \bar{z}), \quad \forall y \in \mathbb{R}^n,$$

which implies that $\nabla f(x) = \nabla_x \phi(x, \bar{z})$.

(b) By part (a), we have

$$f'(x; y) = \max_{z \in Z(x)} \nabla_x \phi(x, z)' y,$$

while by Prop. B.24, we have

$$f'(x; y) = \max_{z \in \partial f(x)} d' y.$$

For all $\bar{z} \in Z(x)$ and $y \in \mathbb{R}^n$, we have

$$\begin{aligned} f(y) &= \max_{z \in Z} \phi(y, z) \\ &\geq \phi(y, \bar{z}) \\ &\geq \phi(x, \bar{z}) + \nabla_x \phi(x, \bar{z})'(y - x) \\ &= f(x) + \nabla_x \phi(x, \bar{z})'(y - x). \end{aligned}$$

Therefore, $\nabla_x \phi(x, \bar{z})$ is a subgradient of f at x , implying that

$$\text{conv}\{\nabla_x \phi(x, z) \mid z \in Z(x)\} \subset \partial f(x).$$

To prove the reverse inclusion, we use a hyperplane separation argument. By the continuity of $\nabla_x \phi(x, \cdot)$ and the compactness of Z , we see that $Z(x)$ is compact, and therefore also the set $\{\nabla_x \phi(x, z) \mid z \in Z(x)\}$ is compact. By Prop. B.8(d), it follows that $\text{conv}\{\nabla_x \phi(x, z) \mid z \in Z(x)\}$ is compact. If $d \in \partial f(x)$ while $d \notin \text{conv}\{\nabla_x \phi(x, z) \mid z \in Z(x)\}$, by the strict separation theorem (Prop. B.14), there exists $y \neq 0$, and $\gamma \in \mathbb{R}$, such that

$$d' y > \gamma > \nabla_x \phi(x, z)' y, \quad \forall z \in Z(x).$$

Therefore, we have

$$d' y > \max_{z \in Z(x)} \nabla_x \phi(x, z)' y = f'(x; y),$$

contradicting Prop. B.24. Therefore, $\partial f(x) \subset \text{conv}\{\nabla_x \phi(x, z) \mid z \in Z(x)\}$ and the proof is complete. **Q.E.D.**

Subgradients of Extended-Real Valued Convex Functions

In this book the major emphasis is on real-valued convex functions $f : \mathbb{R}^n \mapsto \mathbb{R}$, which are defined over the entire space \mathbb{R}^n and are convex over \mathbb{R}^n . There are, however, important cases, prominently arising in the context of duality, where we must deal with functions $g : D \mapsto \mathbb{R}$ that are defined over a convex subset D of \mathbb{R}^n , and are convex over D . This type of function may also be specified as the extended real-valued function $f : \mathbb{R}^n \mapsto (-\infty, \infty]$ given by

$$f(x) = \begin{cases} g(x) & \text{if } x \in D, \\ \infty & \text{otherwise,} \end{cases}$$

with D referred to as the *effective domain* of f .

The notion of a subdifferential and a subgradient of such a function can be developed along the lines of the present section. In particular, given a convex function $f : \mathbb{R}^n \mapsto (-\infty, \infty]$, a vector d is a subgradient of f at a vector x such that $f(x) < \infty$ if the subgradient inequality holds, i.e.,

$$f(z) \geq f(x) + (z - x)'d, \quad \forall z \in \mathbb{R}^n.$$

If $g : D \mapsto \mathbb{R}$ is a concave function (that is, $-g$ is a convex function over the convex set D), it can also be represented as the extended real-valued function $f : \mathbb{R}^n \mapsto [-\infty, \infty)$, where

$$f(x) = \begin{cases} g(x) & \text{if } x \in D, \\ -\infty & \text{otherwise.} \end{cases}$$

As earlier, we say that d is a subgradient of f at an $x \in D$ if $-d$ is a subgradient of the convex function $-g$ at x .

The subdifferential $\partial f(x)$ is the set of all subgradients of the convex (or concave) function f . By convention, $\partial f(x)$ is considered empty for all x with $f(x) = \infty$. Note that contrary to the case of real-valued functions, $\partial f(x)$ may be empty, or closed but unbounded. For example, the extended real-valued convex function given by

$$f(x) = \begin{cases} -\sqrt{x} & \text{if } 0 \leq x \leq 1, \\ \infty & \text{otherwise,} \end{cases}$$

has the subdifferential

$$\partial f(x) = \begin{cases} -\frac{1}{2\sqrt{x}} & \text{if } 0 < x < 1, \\ [-1/2, \infty) & \text{if } x = 1, \\ \emptyset & \text{if } x \leq 0 \text{ or } 1 < x. \end{cases}$$

Thus, $\partial f(x)$ can be empty and can be unbounded at points x that belong to the effective domain of f (as in the cases $x = 0$ and $x = 1$, respectively, of the above example). However, it can be shown that $\partial f(x)$ is nonempty

and compact at points x that are *interior* points of the effective domain of f , as also illustrated by the above example.

One can provide generalized versions of the results of Props. B.24 and B.25 within the context of extended real-valued convex functions, but with appropriate adjustments and additional assumptions to deal with cases where $\partial f(x)$ may be empty or noncompact. The reader will find a detailed account of the corresponding theory in the book by Rockafellar [Roc70].

APPENDIX C:

Line Search Methods

In this appendix we describe algorithms for one-dimensional minimization. These are iterative algorithms, used to implement (approximately) the line minimization stepsize rules.

We briefly present three practical methods. The first two use polynomial interpolation, one requiring derivatives, the second only function values. The third, the Golden Section method, also requires just function values. By contrast with the interpolation methods, it does not depend on the existence of derivatives of the minimized function and may be applied even to discontinuous functions. Its validity depends, however, on a certain unimodality assumption.

In our presentation of the interpolation methods, we consider minimization of the function

$$g(\alpha) = f(x + \alpha d),$$

where f is continuously differentiable. By the chain rule, we have

$$g'(\alpha) = \frac{dg(\alpha)}{d\alpha} = \nabla f(x + \alpha d)'d.$$

We assume that $g'(0) = \nabla f(x)'d < 0$, that is, d is a descent direction at x . We give no convergence or rate of convergence results, but under some fairly natural assumptions, it can be shown that the interpolation methods converge superlinearly.

C.1 CUBIC INTERPOLATION

The cubic interpolation method successively determines at each iteration an appropriate interval $[a, b]$ within which a local minimum of g is guaranteed to exist. It then fits a cubic polynomial to the values $g(a)$, $g(b)$, $g'(a)$,

$g'(b)$. The minimizing point $\bar{\alpha}$ of this cubic polynomial lies within $[a, b]$ and replaces one of the two points a or b for the next iteration.

Cubic Interpolation

Step 1: (Determination of the Initial Interval) Let $s > 0$ be some scalar. (Note: If d "approximates well" the Newton direction, then we take $s = 1$.) Evaluate $g(\alpha)$ and $g'(\alpha)$ at the points $\alpha = 0, s, 2s, 4s, 8s, \dots$, until two successive points a and b are found such that either $g'(b) \geq 0$ or $g(b) \geq g(a)$. Then, it can be seen that a local minimum of g exists within the interval (a, b) . [Note: If $g(s)$ is "much larger" than $g(0)$, it is advisable to replace s by βs , where $\beta \in (0, 1)$, for example $\beta = \frac{1}{2}$ or $\beta = \frac{1}{5}$, and repeat this step.] One can show that this step can be carried out if $\lim_{\alpha \rightarrow \infty} g(\alpha) > g(0)$.

Step 2: (Updating of the Current Interval) Given the current interval $[a, b]$, a cubic polynomial is fitted to the four values $g(a), g'(a), g(b), g'(b)$. The cubic can be shown to have a unique minimum $\bar{\alpha}$ in the interval $(a, b]$ given by

$$\bar{\alpha} = b - \frac{g'(b) + w - z}{g'(b) - g'(a) + 2w}(b - a),$$

where

$$z = \frac{3(g(a) - g(b))}{b - a} + g'(a) + g'(b),$$

$$w = \sqrt{z^2 - g'(a)g'(b)}.$$

If $g'(\bar{\alpha}) \geq 0$ or $g(\bar{\alpha}) \geq g(a)$ replace b by $\bar{\alpha}$. If $g'(\bar{\alpha}) < 0$ and $g(\bar{\alpha}) < g(a)$ replace a by $\bar{\alpha}$. (Note: In practice the computation is terminated once the length of the current interval becomes smaller than a prespecified tolerance or else we obtain $\bar{\alpha} = b$.)

C.2 QUADRATIC INTERPOLATION

This method uses three points a, b , and c such that $a < b < c$, and $g(a) > g(b)$ and $g(b) < g(c)$. Such a set of points is referred to as a *three-point pattern*. It can be seen that a local minimum of g must lie between the extreme points a and c of a three-point pattern a, b, c . At each iteration, the method fits a quadratic polynomial to the three values $g(a), g(b)$, and $g(c)$, and replaces one of the points a, b , and c by the minimizing point of this quadratic polynomial (see Fig. C.1).

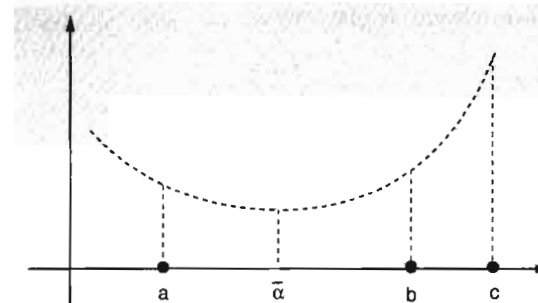


Figure C.1. A three-point pattern and the associated quadratic polynomial. If $\bar{\alpha}$ minimizes the quadratic, a new three point pattern is obtained using $\bar{\alpha}$ and two of the three points a, b , and c (a, b , and $\bar{\alpha}$ in the example of the figure).

Quadratic Interpolation

Step 1: (Determination of Initial Three-Point Pattern) We search along the line as in the cubic interpolation method until we find three successive points a, b , and c with $a < b < c$ such that $g(a) > g(b)$ and $g(b) < g(c)$. As for the cubic interpolation method, we assume that this stage can be carried out, and we can show that this is guaranteed if $\lim_{\alpha \rightarrow \infty} g(\alpha) > g(0)$.

Step 2: (Updating the Current Three-Point Pattern) Given the current three-point pattern a, b, c , we fit a quadratic polynomial to the values $g(a), g(b)$, and $g(c)$, and we determine its unique minimum $\bar{\alpha}$. It can be shown that $\bar{\alpha} \in (a, c)$ and that

$$\bar{\alpha} = \frac{1}{2} \frac{g(a)(c^2 - b^2) + g(b)(a^2 - c^2) + g(c)(b^2 - a^2)}{g(a)(c - b) + g(b)(a - c) + g(c)(b - a)}.$$

Then, we form a new three-point pattern as follows. If $\bar{\alpha} > b$, we replace a or c by $\bar{\alpha}$ depending on whether $g(\bar{\alpha}) < g(b)$ or $g(\bar{\alpha}) > g(b)$, respectively. If $\bar{\alpha} < b$, we replace c or a by $\bar{\alpha}$ depending on whether $g(\bar{\alpha}) < g(b)$ or $g(\bar{\alpha}) > g(b)$, respectively. [Note: If $g(\bar{\alpha}) = g(b)$ then a special local search near $\bar{\alpha}$ should be conducted to replace $\bar{\alpha}$ by a point $\bar{\alpha}'$ with $g(\bar{\alpha}') \neq g(b)$. The computation is terminated when the length of the three-point pattern is smaller than a certain tolerance.]

An alternative possibility for quadratic interpolation is to determine the minimum \bar{a} of the quadratic polynomial that has the same value as g at the points 0 and a , and the same first derivative as g at 0. It can be

verified that this minimum is given by

$$\bar{a} = \frac{g'(0)a^2}{2(g'(0)a + g(0) - g(a))}.$$

C.3 THE GOLDEN SECTION METHOD

Here, we assume that $g(\alpha)$ is *strictly unimodal* in the interval $[0, s]$, as defined in Fig. C.2. The Golden Section method minimizes g over $[0, s]$ by determining at the k th iteration an interval $[\alpha_k, \bar{\alpha}_k]$ containing α^* . These intervals are obtained using the number

$$\tau = \frac{3 - \sqrt{5}}{2},$$

which satisfies $\tau = (1 - \tau)^2$ and is related to the Fibonacci number sequence. The significance of this number will be seen shortly.

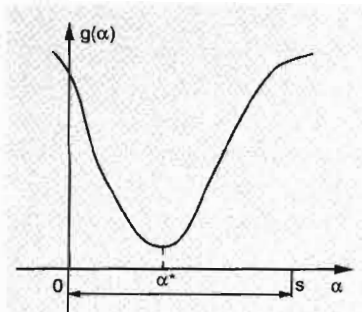


Figure C.2. A strictly unimodal function g over an interval $[0, s]$ is defined as a function that has a unique global minimum α^* in $[0, s]$ and if α_1, α_2 are two points in $[0, s]$ such that $\alpha_1 < \alpha_2 < \alpha^*$ or $\alpha^* < \alpha_1 < \alpha_2$, then $g(\alpha_1) > g(\alpha_2) > g(\alpha^*)$ or $g(\alpha^*) < g(\alpha_1) < g(\alpha_2)$, respectively. An example of a strictly unimodal function, is a function which is strictly convex over $[0, s]$.

Initially, we take

$$[\alpha_0, \bar{\alpha}_0] = [0, s].$$

Given $[\alpha_k, \bar{\alpha}_k]$, we determine $[\alpha_{k+1}, \bar{\alpha}_{k+1}]$ so that $\alpha^* \in [\alpha_{k+1}, \bar{\alpha}_{k+1}]$ as follows. We calculate

$$b_k = \alpha_k + \tau(\bar{\alpha}_k - \alpha_k)$$

$$\bar{b}_k = \bar{\alpha}_k - \tau(\bar{\alpha}_k - \alpha_k)$$

and $g(b_k), g(\bar{b}_k)$. Then:

(1) If $g(b_k) < g(\bar{b}_k)$ we set

$$\alpha_{k+1} = \alpha_k, \quad \bar{\alpha}_{k+1} = b_k \quad \text{if} \quad g(\alpha_k) \leq g(b_k)$$

$$\alpha_{k+1} = \alpha_k, \quad \bar{\alpha}_{k+1} = \bar{b}_k \quad \text{if} \quad g(\alpha_k) > g(b_k).$$

(2) If $g(b_k) > g(\bar{b}_k)$ we set

$$\alpha_{k+1} = \bar{b}_k, \quad \bar{\alpha}_{k+1} = \bar{\alpha}_k \quad \text{if} \quad g(\bar{b}_k) \geq g(\bar{\alpha}_k)$$

$$\alpha_{k+1} = b_k, \quad \bar{\alpha}_{k+1} = \bar{a}_k \quad \text{if} \quad g(\bar{b}_k) < g(\alpha_k).$$

(3) If $g(b_k) = g(\bar{b}_k)$ we set

$$\alpha_{k+1} = b_k, \quad \bar{\alpha}_{k+1} = \bar{b}_k.$$

Based on the definition of a strictly unimodal function it can be shown (see Fig. C.3) that the intervals $[\alpha_k, \bar{\alpha}_k]$ contain α^* and their lengths converge to zero. In practice, the computation is terminated once $(\bar{\alpha}_k - \alpha_k)$ becomes smaller than a prespecified tolerance.

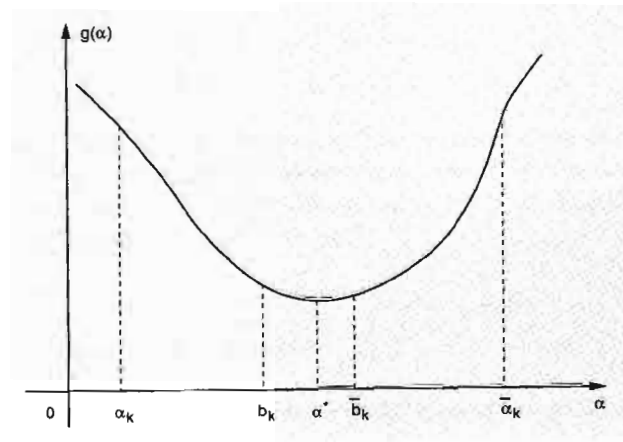


Figure C.3. Golden Section search. Given the interval $[\alpha_k, \bar{\alpha}_k]$ containing the minimum α^* , we calculate

$$b_k = \alpha_k + \tau(\bar{\alpha}_k - \alpha_k)$$

and

$$\bar{b}_k = \bar{\alpha}_k - \tau(\bar{\alpha}_k - \alpha_k).$$

The new interval $[\alpha_{k+1}, \bar{\alpha}_{k+1}]$ has either b_k or \bar{b}_k as one of its endpoints.

An important fact, which rests on the choice of the particular number τ is that

$$[\alpha_{k+1}, \bar{\alpha}_{k+1}] = [\alpha_k, \bar{b}_k] \implies \bar{b}_{k+1} = b_k,$$

$$[\alpha_{k+1}, \bar{\alpha}_{k+1}] = [b_k, \bar{\alpha}_k] \implies b_{k+1} = \bar{b}_k.$$

In other words, a trial point b_k or \bar{b}_k that is not used as the end point of the next interval continues to be a trial point for the next iteration. The reader can verify this, using the property

$$\tau = (1 - \tau)^2.$$

Thus, in either of the above situations, the values \bar{b}_{k+1} , $g(\bar{b}_{k+1})$ or b_{k+1} , $g(b_{k+1})$ are available and need not be recomputed at the next iteration, requiring a single function evaluation instead of two.

APPENDIX D:

Implementation of Newton's Method

In this appendix we describe a globally convergent version of Newton's method based on the modified Cholesky factorization approach discussed in Section 1.4. A computer code implementing the method is available from the author on request.

D.1 CHOLSKY FACTORIZATION

We will give an algorithm for factoring a positive definite symmetric matrix A as

$$A = LL',$$

where L is lower triangular. This is the *Cholesky factorization*. Let a_{ij} be the elements of A and let A_i be the i th leading principal submatrix of A , that is, the submatrix

$$A_i = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1i} \\ a_{21} & a_{22} & \cdots & a_{2i} \\ \vdots & \vdots & \ddots & \vdots \\ a_{i1} & a_{i2} & \cdots & a_{ii} \end{bmatrix}.$$

It is seen that this submatrix is positive definite, since for any $y \in \mathbb{R}_i$, $y \neq 0$, we have by the positive definiteness of A

$$y' A_i y = \begin{bmatrix} y' & 0 \end{bmatrix} A \begin{bmatrix} y \\ 0 \end{bmatrix} > 0.$$

The factorization of A is obtained by successive factorization of A_1, A_2, \dots . Indeed we have $A_1 = L_1 L_1'$, where $L_1 = [\sqrt{a_{11}}]$. Suppose we have the Cholesky factorization of A_{i-1} ,

$$A_{i-1} = L_{i-1} L_{i-1}' \quad (D.1)$$

Let us write

$$A_i = \begin{bmatrix} A_{i-1} & \beta_i \\ \beta_i' & a_{ii} \end{bmatrix}, \quad (D.2)$$

where β_i is the column vector

$$\beta_i = \begin{bmatrix} a_{1i} \\ \vdots \\ a_{i-1,i} \end{bmatrix}. \quad (D.3)$$

Based on Eqs. (D.1)-(D.3), it can be verified that

$$A_i = L_i L_i',$$

where

$$L_i = \begin{bmatrix} L_{i-1} & 0 \\ l_i' & \lambda_{ii} \end{bmatrix}, \quad (D.4)$$

and

$$l_i = L_{i-1}^{-1} \beta_i, \quad \lambda_{ii} = \sqrt{a_{ii} - l_i' l_i}. \quad (D.5)$$

The scalar λ_{ii} is well defined because it can be shown that $a_{ii} - l_i' l_i > 0$. This is seen by defining $b = A_{i-1}^{-1} \beta_i$, and by using the positive definiteness of A_i to write

$$\begin{aligned} 0 &< [b' \quad -1] A_i \begin{bmatrix} b \\ -1 \end{bmatrix} = b' A_{i-1} b - 2b' \beta_i + a_{ii} \\ &= b' \beta_i - 2b' \beta_i + a_{ii} = a_{ii} - b' \beta_i \\ &= a_{ii} - \beta_i' A_{i-1}^{-1} \beta_i = a_{ii} - \beta_i' (L_{i-1} L_{i-1}')^{-1} \beta_i \\ &= a_{ii} - (L_{i-1}^{-1} \beta_i)' (L_{i-1}^{-1} \beta_i) = a_{ii} - l_i' l_i. \end{aligned}$$

The preceding construction can also be used to show that the Cholesky factorization is unique among factorizations involving lower triangular matrices with positive elements along the diagonal. Indeed, A_1 has a unique such factorization, and if A_{i-1} has a unique factorization $A_{i-1} = L_{i-1} L_{i-1}'$, then L_i is uniquely determined from the requirement $A_i = L_i L_i'$ with the diagonal elements of L_i positive, and Eqs. (D.4) and (D.5).

Cholesky Factorization by Columns

In the preceding algorithm, we calculate L by rows, that is, we first calculate the first row of L , then the second row, etc. An alternative and equivalent method is to calculate L by columns, that is, first calculate the first column of L , then the second column, etc. To see how this can be done, we note that the first column of A is equal to the first column of L multiplied with l_{11} , that is,

$$a_{i1} = l_{11} l_{i1}, \quad i = 1, \dots, n,$$

from which we obtain

$$\begin{aligned} l_{11} &= \sqrt{a_{11}}, \\ l_{i1} &= \frac{a_{i1}}{l_{11}}, \quad i = 2, \dots, n. \end{aligned}$$

Similarly, given columns $1, 2, \dots, j-1$ of L , we equate the elements of the j th column of A with the corresponding elements of LL' and we obtain the elements of the j th column of L as follows:

$$\begin{aligned} l_{jj} &= \sqrt{a_{jj} - \sum_{m=1}^{j-1} l_{jm}^2}, \\ l_{ij} &= \frac{a_{ij} - \sum_{m=1}^{j-1} l_{jm} l_{im}}{l_{jj}}, \quad i = j+1, \dots, n. \end{aligned}$$

D.2 APPLICATION TO A MODIFIED NEWTON METHOD

Consider now adding to A a diagonal correction E and simultaneously factoring the matrix

$$F = A + E,$$

where E is such that F is positive definite. The elements of E are introduced sequentially during the factorization process as some diagonal elements of the triangular factor are discovered, which are either negative or are close to zero, indicating that A is either not positive definite or is nearly singular. As discussed in Section 1.4, this is a principal method by which Newton's method is modified to enhance its global convergence properties. The precise mechanization is as follows:

We first fix positive scalars μ_1 and μ_2 , where $\mu_1 < \mu_2$. We calculate the first column of the triangular factor L of F by

$$l_{11} = \begin{cases} \sqrt{a_{11}} & \text{if } \mu_1 < a_{11}, \\ \sqrt{\mu_2} & \text{otherwise,} \end{cases}$$

$$l_{i1} = \frac{a_{i1}}{l_{11}}, \quad i = 2, \dots, n.$$

Similarly, given columns 1, 2, ..., $j-1$ of L , we obtain the elements of the j th column from the equations

$$l_{jj} = \begin{cases} \sqrt{a_{jj} - \sum_{m=1}^{j-1} l_{jm}^2} & \text{if } \mu_1 < a_{11} - \sum_{m=1}^{j-1} l_{jm}^2, \\ \sqrt{\mu_2} & \text{otherwise,} \end{cases}$$

$$l_{ij} = \frac{a_{ij} - \sum_{m=1}^{j-1} l_{jm} l_{im}}{l_{jj}}, \quad i = j+1, \dots, n.$$

In words, if the diagonal element of LL' comes out less than μ_1 , we bring it up to μ_2 .

Note that the j th diagonal element of the correction matrix E is equal to zero if $\mu_1 < a_{jj} - \sum_{m=1}^{j-1} l_{jm}^2$ and is equal to

$$\mu_2 - \left(a_{jj} - \sum_{m=1}^{j-1} l_{jm}^2 \right)$$

otherwise.

The preceding scheme can be used to modify Newton's method, where at the k th iteration, we add a diagonal correction Δ^k to the Hessian $\nabla^2 f(x^k)$ and simultaneously obtain the Cholesky factorization $L^k L^{k'}$ of $\nabla^2 f(x^k) + \Delta^k$ as described above. A modified Newton direction d^k is then obtained by first solving the triangular system

$$L^k y = -\nabla f(x^k),$$

and then solving the triangular system

$$L^{k'} d^k = y.$$

Solving the first system is called *forward elimination* and is accomplished in $O(n^2)$ arithmetic operations using the equations

$$y_1 = -\frac{\partial f(x^k)/\partial x_1}{l_{11}},$$

$$y_i = -\frac{\partial f(x^k)/\partial x_i + \sum_{m=1}^{i-1} l_{im} y_m}{l_{ii}}, \quad i = 2, \dots, n,$$

where l_{im} is the im th element of L^k . Solving the second system is called *back substitution* and is accomplished again in $O(n^2)$ arithmetic operations using the equations

$$d^n = \frac{y^n}{l_{nn}},$$

$$d_i = \frac{y_i - \sum_{m=i+1}^n l_{mi} d_m}{l_{ii}}, \quad i = 1, \dots, n-1.$$

The next point x^{k+1} is obtained from

$$x^{k+1} = x^k + \alpha^k d^k,$$

where α^k is chosen by the Armijo rule with unity initial step whenever the Hessian is not modified ($\Delta^k = 0$) and by means of a line minimization otherwise.

Assuming fixed values of μ_1 and μ_2 , the following may be verified for the modified Newton's method just described:

- The algorithm is globally convergent in the sense that every limit point of $\{x^k\}$ is a stationary point of f . This can be shown using Prop. 1.2.1 in Section 1.2.
- For each local minimum x^* with positive definite Hessian, there exist scalars $\mu > 0$ and $\epsilon > 0$ such that if $\mu_1 \leq \mu$ and $\|x^0 - x^*\| \leq \epsilon$, then $x^k \rightarrow x^*$, $\Delta^k = 0$, and $\alpha^k = 1$ for all k . In other words if μ_1 is not chosen too large, the Hessian will never be modified near x^* , the method will be reduced to the pure form of Newton's method, and the convergence to x^* will be superlinear. The theoretical requirement that μ_1 be sufficiently small can be eliminated by making μ_1 dependent on the norm of the gradient (e.g. $\mu_1 = c\|\nabla f(x^k)\|$, where c is some positive scalar).

Practical Choice of Parameters and Stepsize Selection

We now address some practical issues. As discussed earlier, one should try to choose μ_1 small in order to avoid detrimental modification of the Hessian. Some trial and error with one's particular problem may be required here. As a practical matter, we recommend choosing initially $\mu_1 = 0$ and increasing μ_1 only if difficulties arise due to roundoff error or extremely large norm of calculated direction. (Choosing $\mu_1 = 0$, runs counter to our convergence theory because the generated directions are not guaranteed to be gradient related, but the practical consequences of this are typically insignificant.)

The parameter μ_2 should generally be chosen considerably larger than μ_1 . It can be seen that choosing μ_2 very small can make the modified Hessian matrix $L^k L^{k'}$ nearly singular. On the other hand, choosing μ_2 very large has the effect of making nearly zero the coordinates of d^k that correspond to nonzero diagonal elements of the correction matrix Δ^k . Generally, some trial and error is necessary to determine a proper value of μ_2 . A good guideline is to try a relatively small value of μ_2 and to increase μ_2 if the stepsize generated by the line minimization algorithm is substantially smaller than unity. The idea here is that small values of μ_2 tend to

produce directions d^k with large value of norm and hence small values of stepsize. Thus a small value of stepsize indicates that μ_2 is chosen smaller than appropriate, and suggests that an increase of μ_2 is desirable. It is also possible to construct along these lines an adaptive scheme that changes the values of μ_1 and μ_2 in the course of the algorithm.

The following scheme to set and adjust μ_1 and μ_2 has worked well for the author. At each iteration k , we determine the maximal absolute diagonal element of the Hessian, that is,

$$w^k = \max \left\{ \left| \frac{\partial^2 f(x^k)}{(x_1)^2} \right|, \dots, \left| \frac{\partial^2 f(x^k)}{(x_n)^2} \right| \right\},$$

and we set μ_1 and μ_2 to

$$\mu_1 = r_1 w^k, \quad \mu_2 = r_2 w^k.$$

The scalar r_1 is set at some "small" (or zero) value. The scalar r_2 is changed each time the Hessian is modified; it is multiplied by 5 if the stepsize obtained by the minimization rule is less than 0.2, and it is divided by 5 each time the stepsize is larger than 0.9.

Finally, regarding stepsize selection, any of a large number of possible line minimization algorithms can be used for those iterations where the Hessian is modified (in other iterations the Armijo rule with unity initial stepsize is used). One possibility is to use quadratic interpolation based on function values; see Section C.2 in Appendix C.

It is worth noting that if the cost function is quadratic, then it can be shown that a unity stepsize results in cost reduction for any values of μ_1 and μ_2 . In other words if f is quadratic (not necessarily positive definite), we have

$$f(x^k - (F^k)^{-1} \nabla f(x^k)) \leq f(x^k),$$

where $F^k = \nabla^2 f(x^k) + \Delta^k$ and Δ^k is any positive definite matrix such that F^k is positive definite. As a result, a stepsize near unity is appropriate for initiating the line minimization algorithm. This fact can be used to guide the implementation of the line minimization routine.

References

- [AHR97] Auslender, A., Cominetti, R., and Haddou, M., 1997. "Asymptotic Analysis for Penalty and Barrier Methods in Convex and Linear Programming, Math. Operations Res., Vol. 22, pp. 43-62.
- [AHR93] Anstreicher, K. M., den Hertog, D., Roos, C., and Terlaky, T., 1993. "A Long Step Barrier Method for Convex Quadratic Programming," *Algorithmica*, Vol. 10, pp. 365-382.
- [AHU58] Arrow, K. J., Hurwicz, L., and Uzawa, H., (Eds.), 1958. *Studies in Linear and Nonlinear Programming*, Stanford Univ. Press, Stanford, CA.
- [AHU61] Arrow, K. J., Hurwicz, L., and Uzawa, H., 1961. "Constraint Qualifications in Maximization Problems," *Naval Research Logistics Quarterly*, Vol. 8, pp. 175-191.
- [AaL97] Aarts, E., and Lenstra, J. K., 1997. *Local Search in Combinatorial Optimization*, Wiley, N. Y.
- [Aba67] Abadie, J., 1967. "On the Kuhn-Tucker Theorem," in *Nonlinear Programming*, Abadie, J., (Ed.), North Holland, Amsterdam.
- [AIK90] Al-Khayyal, F., and Kyparisis, J., 1990. "Finite Convergence of Algorithms for Nonlinear Programs and Variational Inequalities," *J. Opt. Theory and Appl.*, Vol. 70, pp. 319-332.
- [Ali92] Alizadeh, F., 1992. "Optimization over the Positive-Definite Cone: Interior Point Methods and Combinatorial Applications," in *Pardalos, P., (Ed.), Advances in Optimization and Parallel Computing*, North Holland, Amsterdam.
- [Ali95] Alizadeh, F., 1995. "Interior-Point Methods in Semidefinite Programming with Applications in Combinatorial Applications," *SIAM J. on Optimization*, Vol. 5, pp. 13-51.
- [AnV94] Anstreicher, K. M., and Vial, J.-P., 1994. "On the Convergence of an Infeasible Primal-Dual Interior-Point Method for Convex Programming," *Optimization Methods and Software*, Vol. 3, pp. 273-283.
- [Arm66] Armijo, L., 1966. "Minimization of Functions Having Continuous Partial Derivatives," *Pacific J. Math.*, Vol. 16, pp. 1-3.