

Nonlinear Optimization

James V. Burke
University of Washington

Contents

Chapter 1. Introduction	5
Chapter 2. Review of Matrices and Block Structures	7
1. Rows and Columns	7
2. Matrix Multiplication	9
3. Block Matrix Multiplication	11
4. Gauss-Jordan Elimination Matrices and Reduction to Reduced Echelon Form	13
5. Some Special Square Matrices	15
6. The LU Factorization	16
7. Solving Equations with the LU Factorization	18
8. The Four Fundamental Subspaces and Echelon Form	19
Chapter 3. The Linear Least Squares Problem	21
1. Applications	21
2. Optimality in the Linear Least Squares Problem	26
3. Orthogonal Projection onto a Subspace	27
4. Minimal Norm Solutions to $Ax = b$	29
5. Gram-Schmidt Orthogonalization, the QR Factorization, and Solving the Normal Equations	30
Chapter 4. Optimization of Quadratic Functions	37
1. Eigenvalue Decomposition of Symmetric Matrices	37
2. Optimality Properties of Quadratic Functions	40
3. Minimization of a Quadratic Function on an Affine Set	42
4. The Principal Minor Test for Positive Definiteness	44
5. The Cholesky Factorizations	45
6. Linear Least Squares Revisited	48
7. The Conjugate Gradient Algorithm	48
Chapter 5. Elements of Multivariable Calculus	53
1. Norms and Continuity	53
2. Differentiation	55
3. The Delta Method for Computing Derivatives	58
4. Differential Calculus	59
5. The Mean Value Theorem	59
Chapter 6. Optimality Conditions for Unconstrained Problems	63
1. Existence of Optimal Solutions	63
2. First-Order Optimality Conditions	64
3. Second-Order Optimality Conditions	65
4. Convexity	66
Chapter 7. Optimality Conditions for Constrained Optimization	73
1. First-Order Conditions	73
2. Regularity and Constraint Qualifications	76
3. Second-Order Conditions	78
4. Optimality Conditions in the Presence of Convexity	79
5. Convex Optimization, Saddle Point Theory, and Lagrangian Duality	83

Chapter 8. Line Search Methods	89
1. The Basic Backtracking Algorithm	89
2. The Wolfe Conditions	94
Chapter 9. Search Directions for Unconstrained Optimization	99
1. Rate of Convergence	99
2. Newton's Method for Solving Equations	99
3. Newton's Method for Minimization	102
4. Matrix Secant Methods	103
Index	109

Introduction

In mathematical optimization we seek to either minimize or maximize a function over a set of alternatives. The function is called the *objective function*, and we allow it to be transfinite in the sense that at each point its value is either a real number or it is one of the infinite values $\pm\infty$. The set of alternatives is called the *constraint region*. Since every maximization problem can be restated as a minimization problem by simply replacing the objective f_0 by its negative $-f_0$ (and visa versa), we choose to focus only on minimization problems. We denote such problems using the notation

$$(1) \quad \begin{array}{l} \text{minimize } f_0(x) \\ \quad \quad \quad x \in X \\ \text{subject to } x \in \Omega, \end{array}$$

where $f_0 : X \rightarrow \mathbb{R} \cup \{\pm\infty\}$ is the objective function, X is the space over which the optimization occurs, and $\Omega \subset X$ is the constraint region. This is a very general description of an optimization problem and as one might imagine there is a taxonomy of optimization problems depending on the underlying structural features that the problem possesses, e.g., properties of the space X , is it the integers, the real numbers, the complex numbers, matrices, or an infinite dimensional space of functions, properties of the function f_0 , is it discrete, continuous, or differentiable, the geometry of the set Ω , how Ω is represented, properties of the underlying applications and how they fit into a broader context, methods of solution or approximate solution, For our purposes, we assume that Ω is a subset of \mathbb{R}^n and that $f_0 : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\pm\infty\}$. This severely restricts the kind of optimization problems that we study, however, it is sufficiently broad to include a wide variety of applied problems of great practical importance and interest. For example, this framework includes *linear programming* (LP).

Linear Programming

In the case of LP, the objective function is linear, that is, there exists $c \in \mathbb{R}^n$ such that

$$f_0(x) = c^T x = \sum_{j=1}^n c_j x_j,$$

and the constraint region is representable as the set of solution to a finite system of linear equation and inequalities,

$$(2) \quad \Omega = \left\{ x \in \mathbb{R}^n \mid \sum_{j=1}^n a_{ij} x_j \leq b_i, \quad i = 1, \dots, s, \quad \sum_{j=1}^n a_{ij} x_j = b_i, \quad i = s + 1, \dots, m \right\},$$

where $A := [a_{ij}] \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$.

However, in this course we are primarily concerned with nonlinear problems, that is, problems that cannot be encoded using finitely many linear function alone. A natural generalization of the LP framework to the nonlinear setting is to simply replace each of the linear functions with a nonlinear function. This leads to the general *nonlinear programming* (NLP) problem which is the problem of central concern in these notes.

Nonlinear Programming

In nonlinear programming we are given nonlinear functions $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$, $i = 1, 2, \dots, m$, where f_0 is the objective function in (1) and the functions f_i , $i = 1, 2, \dots, m$ are called the constraint functions which are used to define the constrain region in (1) by setting

$$(3) \quad \Omega = \{x \in \mathbb{R}^n \mid f_i(x) \leq 0, \quad i = 1, \dots, s, \quad f_i(x) = 0, \quad i = s + 1, \dots, m\} .$$

If $\Omega = \mathbb{R}^n$, then we say that the problem (1) is an *unconstrained optimization problem*; otherwise, it called a constrained problem. We begin or study with unconstrained problems. They are simpler to handle since we are only

concerned with minimizing the objective function and we need not concern ourselves with the constraint region. However, since we allow the objective to take infinite values, we shall see that every explicitly constrained problem can be restated as an ostensibly unconstrained problem.

In the following section, we begin our study of unconstrained optimization which is arguably the most widely studied and used class of unconstrained nonlinear optimization problems. This is the class of *linear least squares* problems. The theory and techniques we develop for this class of problems provides a template for how we address and exploit structure in a wide variety of other problem classes.

Linear Least Squares

A *linear least squares problem* is one of the form

$$(4) \quad \underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2} \|Ax - b\|_2^2,$$

where

$$A \in \mathbb{R}^{m \times n}, \quad b \in \mathbb{R}^m, \quad \text{and} \quad \|y\|_2^2 := y_1^2 + y_2^2 + \cdots + y_m^2.$$

Problems of this type arise in a diverse range of application, some of which are discussed in later chapters. Whole books have been written about this problem, and various instances of this problem remain a very active area of research. This problem formulation is usually credited to Legendre and Gauss who made careful studies of the method around 1800. But others had applied the basic approach in a ad hoc manner in the previous 50 years to observational data and, in particular, to studying the motion of the planets.

The second class most important class of unconstrained nonlinear optimization problems is the minimization of *quadratic functions*. As we will see, the linear least squares problem is a member of this class of problems. It is an important for a wide variety of reasons, not the least of which is the relationship to the second-order Taylor approximations for functions mapping \mathbb{R}^n into \mathbb{R} .

Quadratic Functions

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is said to be *quadratic* if there exists $\alpha \in \mathbb{R}$, $g \in \mathbb{R}^n$ and $H \in \mathbb{R}^{n \times n}$ such that

$$f(x) = \alpha + g^T x + \frac{1}{2} x^T H x.$$

The first thing to notice about such functions is that we may as well assume that the matrix H is symmetric since

$$x^T H x = \frac{1}{2} (x^T H x + x^T H x) = \frac{1}{2} ((x^T H x)^T + x^T H x) = \frac{1}{2} (x^T H^T x + x^T H x) = x^T (\frac{1}{2} (H^T + H)) x,$$

that is, we may as well replace the matrix H by its symmetric part $\frac{1}{2} (H^T + H)$.

Having quadratic functions in hand, one arrives at an important nonlinear generalization of linear programming where we simply replace the LP linear objective with a quadratic function.

Quadratic Programming

In *quadratic programming* we minimize a quadratic objective function subject convex polyhedral constraints of the form (2).

The linear least squares problem and the optimization of quadratic functions are the themes for our initial forays into optimization. The theory and methods we develop for these problems as well as certain variations on these problems form the basis for our extensions to other problem classes. For this reason, we study these problems with great care. Notice that although these problems are nonlinear, their component pieces come from linear algebra, that is matrices and vectors. Obviously, these components play a key role in understanding the structure and behavior of these problems. For this reason, our first task is to review and develop the essential elements from linear algebra that provide the basis for our investigation into these problems.

Review of Matrices and Block Structures

Numerical linear algebra lies at the heart of modern scientific computing and computational science. Today it is not uncommon to perform numerical computations with matrices having millions of components. The key to understanding how to implement such algorithms is to exploit underlying structure within the matrices. In these notes we touch on a few ideas and tools for dissecting matrix structure. Specifically we are concerned with the *block structure* matrices.

1. Rows and Columns

Let $A \in \mathbb{R}^{m \times n}$ so that A has m rows and n columns. Denote the element of A in the i th row and j th column as A_{ij} . Denote the m rows of A by $A_{1.}, A_{2.}, A_{3.}, \dots, A_{m.}$ and the n columns of A by $A_{.1}, A_{.2}, A_{.3}, \dots, A_{.n}$. For example, if

$$A = \begin{bmatrix} 3 & 2 & -1 & 5 & 7 & 3 \\ -2 & 27 & 32 & -100 & 0 & 0 \\ -89 & 0 & 47 & 22 & -21 & 33 \end{bmatrix},$$

then $A_{2,4} = -100$, and

$$A_{1.} = [3 \ 2 \ -1 \ 5 \ 7 \ 3], \quad A_{2.} = [-2 \ 27 \ 32 \ -100 \ 0 \ 0], \quad A_{3.} = [-89 \ 0 \ 47 \ 22 \ -21 \ 33]$$

and

$$A_{.1} = \begin{bmatrix} 3 \\ -2 \\ -89 \end{bmatrix}, \quad A_{.2} = \begin{bmatrix} 2 \\ 27 \\ 0 \end{bmatrix}, \quad A_{.3} = \begin{bmatrix} -1 \\ 32 \\ 47 \end{bmatrix}, \quad A_{.4} = \begin{bmatrix} 5 \\ -100 \\ 22 \end{bmatrix}, \quad A_{.5} = \begin{bmatrix} 7 \\ 0 \\ -21 \end{bmatrix}, \quad A_{.6} = \begin{bmatrix} 3 \\ 0 \\ 33 \end{bmatrix}.$$

EXERCISE 1.1. If

$$C = \begin{bmatrix} 3 & -4 & 1 & 1 & 0 & 0 \\ 2 & 2 & 0 & 0 & 1 & 0 \\ -1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 2 & 1 & 4 \\ 0 & 0 & 0 & 1 & 0 & 3 \end{bmatrix},$$

what are $C_{4,4}$, $C_{.4}$ and C_4 ? For example, $C_{2.} = [2 \ 2 \ 0 \ 0 \ 1 \ 0]$ and $C_{.2} = \begin{bmatrix} -4 \\ 2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$.

The block structuring of a matrix into its rows and columns is of fundamental importance and is extremely useful in understanding the properties of a matrix. In particular, for $A \in \mathbb{R}^{m \times n}$ it allows us to write

$$A = \begin{bmatrix} A_{1.} \\ A_{2.} \\ A_{3.} \\ \vdots \\ A_{m.} \end{bmatrix} \quad \text{and} \quad A = [A_{.1} \ A_{.2} \ A_{.3} \ \dots \ A_{.n}].$$

These are called the row and column block representations of A , respectively

1.1. Matrix vector Multiplication. Let $A \in \mathbb{R}^{m \times n}$ and $x \in \mathbb{R}^n$. In terms of its coordinates (or components),

we can also write $x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$ with each $x_j \in \mathbb{R}$. The term x_j is called the j th component of x . For example if

$$x = \begin{bmatrix} 5 \\ -100 \\ 22 \end{bmatrix},$$

then $n = 3$, $x_1 = 5$, $x_2 = -100$, $x_3 = 22$. We define the matrix-vector product Ax by

$$Ax = \begin{bmatrix} A_{1.} \bullet x \\ A_{2.} \bullet x \\ A_{3.} \bullet x \\ \vdots \\ A_{m.} \bullet x \end{bmatrix},$$

where for each $i = 1, 2, \dots, m$, $A_{i.} \bullet x$ is the dot product of the i th row of A with x and is given by

$$A_{i.} \bullet x = \sum_{j=1}^n A_{ij}x_j.$$

For example, if

$$A = \begin{bmatrix} 3 & 2 & -1 & 5 & 7 & 3 \\ -2 & 27 & 32 & -100 & 0 & 0 \\ -89 & 0 & 47 & 22 & -21 & 33 \end{bmatrix} \text{ and } x = \begin{bmatrix} 1 \\ -1 \\ 0 \\ 0 \\ 2 \\ 3 \end{bmatrix},$$

then

$$Ax = \begin{bmatrix} 24 \\ -29 \\ -32 \end{bmatrix}.$$

EXERCISE 1.2. If

$$C = \begin{bmatrix} 3 & -4 & 1 & 1 & 0 & 0 \\ 2 & 2 & 0 & 0 & 1 & 0 \\ -1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 2 & 1 & 4 \\ 0 & 0 & 0 & 1 & 0 & 3 \end{bmatrix} \text{ and } x = \begin{bmatrix} 1 \\ -1 \\ 0 \\ 0 \\ 2 \\ 3 \end{bmatrix},$$

what is Cx ?

Note that if $A \in \mathbb{R}^{m \times n}$ and $x \in \mathbb{R}^n$, then Ax is always well defined with $Ax \in \mathbb{R}^m$. In terms of components, the i th component of Ax is given by the dot product of the i th row of A (i.e. $A_{i.}$) and x (i.e. $A_{i.} \bullet x$).

The view of the matrix-vector product described above is the *row-space* perspective, where the term *row-space* will be given a more rigorous definition at a later time. But there is a very different way of viewing the matrix-vector product based on a *column-space* perspective. This view uses the notion of the linear combination of a collection of vectors.

Given k vectors $v^1, v^2, \dots, v^k \in \mathbb{R}^n$ and k scalars $\alpha_1, \alpha_2, \dots, \alpha_k \in \mathbb{R}$, we can form the vector

$$\alpha_1 v^1 + \alpha_2 v^2 + \dots + \alpha_k v^k \in \mathbb{R}^n.$$

Any vector of this kind is said to be a *linear combination* of the vectors v^1, v^2, \dots, v^k where the $\alpha_1, \alpha_2, \dots, \alpha_k$ are called the coefficients in the linear combination. The set of all such vectors formed as linear combinations of v^1, v^2, \dots, v^k is said to be the *linear span* of v^1, v^2, \dots, v^k and is denoted

$$\text{span}(v^1, v^2, \dots, v^k) := \{ \alpha_1 v^1 + \alpha_2 v^2 + \dots + \alpha_k v^k \mid \alpha_1, \alpha_2, \dots, \alpha_k \in \mathbb{R} \}.$$

Returning to the matrix-vector product, one has that

$$Ax = \begin{bmatrix} A_{11}x_1 + A_{12}x_2 + A_{13}x_3 + \cdots + A_{1n}x_n \\ A_{21}x_1 + A_{22}x_2 + A_{23}x_3 + \cdots + A_{2n}x_n \\ \vdots \\ A_{m1}x_1 + A_{m2}x_2 + A_{m3}x_3 + \cdots + A_{mn}x_n \end{bmatrix} = x_1A_{\cdot 1} + x_2A_{\cdot 2} + x_3A_{\cdot 3} + \cdots + x_nA_{\cdot n},$$

which is a linear combination of the columns of A . That is, we can view the matrix-vector product Ax as taking a linear combination of the columns of A where the coefficients in the linear combination are the coordinates of the vector x .

We now have two fundamentally different ways of viewing the matrix-vector product Ax .

Row-Space view of Ax :

$$Ax = \begin{bmatrix} A_{1\cdot} \bullet x \\ A_{2\cdot} \bullet x \\ A_{3\cdot} \bullet x \\ \vdots \\ A_{m\cdot} \bullet x \end{bmatrix}$$

Column-Space view of Ax :

$$Ax = x_1A_{\cdot 1} + x_2A_{\cdot 2} + x_3A_{\cdot 3} + \cdots + x_nA_{\cdot n}.$$

2. Matrix Multiplication

We now build on our notion of a matrix-vector product to define a notion of a matrix-matrix product which we call *matrix multiplication*. Given two matrices $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{n \times k}$ note that each of the columns of B resides in \mathbb{R}^n , i.e. $B_{\cdot j} \in \mathbb{R}^n$ $j = 1, 2, \dots, k$. Therefore, each of the matrix-vector products $AB_{\cdot j}$ is well defined for $j = 1, 2, \dots, k$. This allows us to define a matrix-matrix product that exploits the block column structure of B by setting

$$(5) \quad AB := [AB_{\cdot 1} \quad AB_{\cdot 2} \quad AB_{\cdot 3} \quad \cdots \quad AB_{\cdot k}].$$

Note that the j th column of AB is $(AB)_{\cdot j} = AB_{\cdot j} \in \mathbb{R}^m$ and that $AB \in \mathbb{R}^{m \times k}$, i.e.

$$\text{if } H \in \mathbb{R}^{m \times n} \text{ and } L \in \mathbb{R}^{n \times k}, \text{ then } HL \in \mathbb{R}^{m \times k}.$$

Also note that

$$\text{if } T \in \mathbb{R}^{s \times t} \text{ and } M \in \mathbb{R}^{r \times \ell}, \text{ then the matrix product } TM \text{ is only defined when } t = r.$$

For example, if

$$A = \begin{bmatrix} 3 & 2 & -1 & 5 & 7 & 3 \\ -2 & 27 & 32 & -100 & 0 & 0 \\ -89 & 0 & 47 & 22 & -21 & 33 \end{bmatrix} \text{ and } B = \begin{bmatrix} 2 & 0 \\ -2 & 2 \\ 0 & 3 \\ 0 & 0 \\ 1 & 1 \\ 2 & -1 \end{bmatrix},$$

then

$$AB = \begin{bmatrix} A & \begin{bmatrix} 2 \\ -2 \\ 0 \\ 0 \\ 1 \\ 2 \end{bmatrix} \\ A & \begin{bmatrix} 0 \\ -2 \\ 3 \\ 0 \\ 1 \\ -1 \end{bmatrix} \end{bmatrix} = \begin{bmatrix} 15 & 5 \\ -58 & 150 \\ -133 & 87 \end{bmatrix}.$$

EXERCISE 2.1. if

$$C = \begin{bmatrix} 3 & -4 & 1 & 1 \\ 2 & 2 & 0 & 0 \\ -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2 \\ 0 & 1 & 0 & 1 \end{bmatrix} \text{ and } D = \begin{bmatrix} -1 & 0 & 2 & 4 & 3 \\ 0 & -2 & -1 & 4 & 5 \\ 5 & 2 & -4 & 1 & 1 \\ 3 & 0 & 1 & 0 & 0 \end{bmatrix},$$

is CD well defined and if so what is it?

The formula (5) can be used to give further insight into the individual components of the matrix product AB . By the definition of the matrix-vector product we have for each $j = 1, 2, \dots, k$

$$AB_{\cdot j} = \begin{bmatrix} A_{1\cdot} \bullet B_{\cdot j} \\ A_{2\cdot} \bullet B_{\cdot j} \\ \vdots \\ A_{m\cdot} \bullet B_{\cdot j} \end{bmatrix}.$$

Consequently,

$$(AB)_{ij} = A_{i\cdot} \bullet B_{\cdot j} \quad \forall i = 1, 2, \dots, m, j = 1, 2, \dots, k.$$

That is, the element of AB in the i th row and j th column, $(AB)_{ij}$, is the dot product of the i th row of A with the j th column of B .

2.1. Elementary Matrices. We define the *elementary unit coordinate matrices* in $\mathbb{R}^{m \times n}$ in much the same way as we define the elementary unit coordinate vectors. Given $i \in \{1, 2, \dots, m\}$ and $j \in \{1, 2, \dots, n\}$, the elementary unit coordinate matrix $E_{ij} \in \mathbb{R}^{m \times n}$ is the matrix whose ij entry is 1 with all other entries taking the value zero. This is a slight abuse of notation since the notation E_{ij} is supposed to represent the ij th entry in the matrix E . To avoid confusion, we reserve the use of the letter E when speaking of matrices to the elementary matrices.

EXERCISE 2.2. (*Multiplication of square elementary matrices*) Let $i, k \in \{1, 2, \dots, m\}$ and $j, \ell \in \{1, 2, \dots, m\}$. Show the following for elementary matrices in $\mathbb{R}^{m \times m}$ first for $m = 3$ and then in general.

$$(1) E_{ij}E_{k\ell} = \begin{cases} E_{i\ell} & , \text{ if } j = k, \\ 0 & , \text{ otherwise.} \end{cases}$$

(2) For any $\alpha \in \mathbb{R}$, if $i \neq j$, then $(I_{m \times m} - \alpha E_{ij})(I_{m \times m} + \alpha E_{ij}) = I_{m \times m}$ so that

$$(I_{m \times m} + \alpha E_{ij})^{-1} = (I_{m \times m} - \alpha E_{ij}).$$

(3) For any $\alpha \in \mathbb{R}$ with $\alpha \neq 0$, $(I + (\alpha^{-1} - 1)E_{ii})(I + (\alpha - 1)E_{ii}) = I$ so that

$$(I + (\alpha - 1)E_{ii})^{-1} = (I + (\alpha^{-1} - 1)E_{ii}).$$

EXERCISE 2.3. (*Elementary permutation matrices*) Let $i, \ell \in \{1, 2, \dots, m\}$ and consider the matrix $P_{ij} \in \mathbb{R}^{m \times m}$ obtained from the identity matrix by interchanging its i and ℓ th rows. We call such a matrix an *elementary permutation matrix*. Again we are abusing notation, but again we reserve the letter P for permutation matrices (and, later, for projection matrices). Show the following are true first for $m = 3$ and then in general.

$$(1) P_{i\ell}P_{i\ell} = I_{m \times m} \text{ so that } P_{i\ell}^{-1} = P_{i\ell}.$$

$$(2) P_{i\ell}^T = P_{i\ell}.$$

$$(3) P_{i\ell} = I - E_{ii} - E_{\ell\ell} + E_{i\ell} + E_{\ell i}.$$

EXERCISE 2.4. (*Three elementary row operations as matrix multiplication*) In this exercise we show that the three elementary row operations can be performed by left multiplication by an invertible matrix. Let $A \in \mathbb{R}^{m \times n}$, $\alpha \in \mathbb{R}$ and let $i, \ell \in \{1, 2, \dots, m\}$ and $j \in \{1, 2, \dots, n\}$. Show that the following results hold first for $m = n = 3$ and then in general.

(1) (*row interchanges*) Given $A \in \mathbb{R}^{m \times n}$, the matrix $P_{ij}A$ is the same as the matrix A except with the i and j th rows interchanged.

(2) (*row multiplication*) Given $\alpha \in \mathbb{R}$ with $\alpha \neq 0$, show that the matrix $(I + (\alpha - 1)E_{ii})A$ is the same as the matrix A except with the i th row replaced by α times the i th row of A .

(3) Show that matrix $E_{ij}A$ is the matrix that contains the j th row of A in its i th row with all other entries equal to zero.

(4) (*replace a row by itself plus a multiple of another row*) Given $\alpha \in \mathbb{R}$ and $i \neq j$, show that the matrix $(I + \alpha E_{ij})A$ is the same as the matrix A except with the i th row replaced by itself plus α times the j th row of A .

2.2. Associativity of matrix multiplication. Note that the definition of matrix multiplication tells us that this operation is associative. That is, if $A \in \mathbb{R}^{m \times n}$, $B \in \mathbb{R}^{n \times k}$, and $C \in \mathbb{R}^{k \times s}$, then $AB \in \mathbb{R}^{m \times k}$ so that $(AB)C$ is well defined and $BC \in \mathbb{R}^{n \times s}$ so that $A(BC)$ is well defined, and, moreover,

$$(6) \quad (AB)C = [(AB)C_{.1} \quad (AB)C_{.2} \quad \cdots \quad (AB)C_{.s}]$$

where for each $\ell = 1, 2, \dots, s$

$$\begin{aligned} (AB)C_{. \ell} &= [AB_{.1} \quad AB_{.2} \quad AB_{.3} \quad \cdots \quad AB_{.k}] C_{. \ell} \\ &= C_{1 \ell} AB_{.1} + C_{2 \ell} AB_{.2} + \cdots + C_{k \ell} AB_{.k} \\ &= A [C_{1 \ell} B_{.1} + C_{2 \ell} B_{.2} + \cdots + C_{k \ell} B_{.k}] \\ &= A(BC_{. \ell}). \end{aligned}$$

Therefore, we may write (6) as

$$\begin{aligned} (AB)C &= [(AB)C_{.1} \quad (AB)C_{.2} \quad \cdots \quad (AB)C_{.s}] \\ &= [A(BC_{.1}) \quad A(BC_{.2}) \quad \cdots \quad A(BC_{.s})] \\ &= A [BC_{.1} \quad BC_{.2} \quad \cdots \quad BC_{.s}] \\ &= A(BC). \end{aligned}$$

Due to this associativity property, we may dispense with the parentheses and simply write ABC for this triple matrix product. Obviously longer products are possible.

EXERCISE 2.5. Consider the following matrices:

$$\begin{aligned} A &= \begin{bmatrix} 2 & 3 & 1 \\ 1 & 0 & -3 \end{bmatrix} & B &= \begin{bmatrix} 4 & -1 \\ 0 & -7 \end{bmatrix} & C &= \begin{bmatrix} -2 & 3 & 2 \\ 1 & 1 & -3 \\ 2 & 1 & 0 \end{bmatrix} \\ D &= \begin{bmatrix} 2 & 3 \\ 1 & 0 \\ 8 & -5 \end{bmatrix} & F &= \begin{bmatrix} 2 & 1 & 1 & 2 \\ 1 & 0 & -4 & 0 \\ 3 & 0 & -2 & 0 \\ 5 & 1 & 1 & 1 \end{bmatrix} & G &= \begin{bmatrix} 2 & 3 & 1 & -2 \\ 1 & 0 & -3 & 0 \end{bmatrix}. \end{aligned}$$

Using these matrices, which pairs can be multiplied together and in what order? Which triples can be multiplied together and in what order (e.g. the triple product BAC is well defined)? Which quadruples can be multiplied together and in what order? Perform all of these multiplications.

3. Block Matrix Multiplication

To illustrate the general idea of block structures consider the following matrix.

$$A = \begin{bmatrix} 3 & -4 & 1 & 1 & 0 & 0 \\ 0 & 2 & 2 & 0 & 1 & 0 \\ 1 & 0 & -1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 2 & 1 & 4 \\ 0 & 0 & 0 & 1 & 0 & 3 \end{bmatrix}.$$

Visual inspection tells us that this matrix has structure. But what is it, and how can it be represented? We re-write the the matrix given above *blocking* out some key structures:

$$A = \left[\begin{array}{ccc|ccc} 3 & -4 & 1 & 1 & 0 & 0 \\ 0 & 2 & 2 & 0 & 1 & 0 \\ 1 & 0 & -1 & 0 & 0 & 1 \\ \hline 0 & 0 & 0 & 2 & 1 & 4 \\ 0 & 0 & 0 & 1 & 0 & 3 \end{array} \right] = \left[\begin{array}{c|c} B & I_{3 \times 3} \\ \hline 0_{2 \times 3} & C \end{array} \right],$$

where

$$B = \begin{bmatrix} 3 & -4 & 1 \\ 0 & 2 & 2 \\ 1 & 0 & -1 \end{bmatrix}, \quad C = \begin{bmatrix} 2 & 1 & 4 \\ 1 & 0 & 3 \end{bmatrix},$$

$I_{3 \times 3}$ is the 3×3 identity matrix, and $0_{2 \times 3}$ is the 2×3 zero matrix. Having established this structure for the matrix A , it can now be exploited in various ways. As a simple example, we consider how it can be used in matrix multiplication.

Consider the matrix

$$M = \begin{bmatrix} 1 & 2 \\ 0 & 4 \\ -1 & -1 \\ 2 & -1 \\ 4 & 3 \\ -2 & 0 \end{bmatrix}.$$

The matrix product AM is well defined since A is 5×6 and M is 6×2 . We show how to compute this matrix product using the structure of A . To do this we must first *block decompose* M conformally with the block decomposition of A . Another way to say this is that we must give M a block structure that allows us to do block matrix multiplication with the blocks of A . The correct block structure for M is

$$M = \begin{bmatrix} X \\ Y \end{bmatrix},$$

where

$$X = \begin{bmatrix} 1 & 2 \\ 0 & 4 \\ -1 & -1 \end{bmatrix}, \quad \text{and} \quad Y = \begin{bmatrix} 2 & -1 \\ 4 & 3 \\ -2 & 0 \end{bmatrix},$$

since then X can multiply $\begin{bmatrix} B \\ 0_{2 \times 3} \end{bmatrix}$ and Y can multiply $\begin{bmatrix} I_{3 \times 3} \\ C \end{bmatrix}$. This gives

$$\begin{aligned} AM &= \begin{bmatrix} B & I_{3 \times 3} \\ 0_{2 \times 3} & C \end{bmatrix} \begin{bmatrix} X \\ Y \end{bmatrix} = \begin{bmatrix} BX + Y \\ CY \end{bmatrix} \\ &= \begin{bmatrix} \begin{bmatrix} 2 & -11 \\ 2 & 12 \\ -1 & -2 \end{bmatrix} + \begin{bmatrix} -2 & 6 \\ 4 & 3 \\ -2 & 0 \end{bmatrix} \\ \begin{bmatrix} 0 & 1 \\ -4 & -1 \end{bmatrix} \end{bmatrix} \\ &= \begin{bmatrix} 4 & -12 \\ 2 & 9 \\ 0 & 3 \\ 0 & 1 \\ -4 & -1 \end{bmatrix}. \end{aligned}$$

Block structured matrices and their matrix product is a very powerful tool in matrix analysis. Consider the matrices $M \in \mathbb{R}^{n \times m}$ and $T \in \mathbb{R}^{m \times k}$ given by

$$M = \begin{bmatrix} A_{n_1 \times m_1} & B_{n_1 \times m_2} \\ C_{n_2 \times m_1} & D_{n_2 \times m_2} \end{bmatrix}$$

and

$$T = \begin{bmatrix} E_{m_1 \times k_1} & F_{m_1 \times k_2} & G_{m_1 \times k_3} \\ H_{m_2 \times k_1} & J_{m_2 \times k_2} & K_{m_2 \times k_3} \end{bmatrix},$$

where $n = n_1 + n_2$, $m = m_1 + m_2$, and $k = k_1 + k_2 + k_3$. The block structures for the matrices M and T are said to be *conformal* with respect to matrix multiplication since

$$MT = \begin{bmatrix} AE + BH & AF + BJ & AG + BK \\ CE + DH & CF + DJ & CG + DK \end{bmatrix}.$$

Similarly, one can conformally block structure matrices with respect to matrix addition (how is this done?).

EXERCISE 3.1. Consider the matrix

$$H = \begin{bmatrix} -2 & 3 & 2 & 0 & 0 & 0 & 0 \\ 1 & 1 & -3 & 0 & 0 & 0 & 0 \\ 2 & 1 & -1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 4 & -1 & 0 & 0 \\ 0 & 0 & 0 & 2 & -7 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 2 & 3 \\ 0 & 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 8 & -5 \end{bmatrix}.$$

Does H have a natural block structure that might be useful in performing a matrix-matrix multiply, and if so describe it by giving the blocks? Describe a conformal block decomposition of the matrix

$$M = \begin{bmatrix} 1 & 2 \\ 3 & -4 \\ -5 & 6 \\ 1 & -2 \\ -3 & 4 \\ 1 & 1 \\ 1 & 1 \end{bmatrix}$$

that would be useful in performing the matrix product HM . Compute the matrix product HM using this conformal decomposition.

EXERCISE 3.2. Let $T \in \mathbb{R}^{m \times n}$ with $T \neq 0$ and let I be the $m \times m$ identity matrix. Consider the block structured matrix $A = \begin{bmatrix} I & T \end{bmatrix}$.

- (i) If $A \in \mathbb{R}^{k \times s}$, what are k and s ?
- (ii) Construct a non-zero $s \times n$ matrix B such that $AB = 0$.

The examples given above illustrate how block matrix multiplication works and why it might be useful. One of the most powerful uses of block structures is in understanding and implementing standard *matrix factorizations* or reductions.

4. Gauss-Jordan Elimination Matrices and Reduction to Reduced Echelon Form

In this section, we show that Gaussian-Jordan elimination can be represented as a consequence of left multiplication by a specially designed matrix called a *Gaussian-Jordan elimination matrix*.

Consider the vector $v \in \mathbb{R}^m$ block decomposed as

$$v = \begin{bmatrix} a \\ \alpha \\ b \end{bmatrix}$$

where $a \in \mathbb{R}^s$, $\alpha \in \mathbb{R}$, and $b \in \mathbb{R}^t$ with $m = s + 1 + t$. In this vector we refer to the α entry as the *pivot* and assume that $\alpha \neq 0$. We wish to determine a matrix G such that

$$Gv = e_{s+1}$$

where for $j = 1, \dots, n$, e_j is the unit coordinate vector having a one in the j th position and zeros elsewhere. We claim that the matrix

$$G = \begin{bmatrix} I_{s \times s} & -\alpha^{-1}a & 0 \\ 0 & \alpha^{-1} & 0 \\ 0 & -\alpha^{-1}b & I_{t \times t} \end{bmatrix}$$

does the trick. Indeed,

$$(7) \quad Gv = \begin{bmatrix} I_{s \times s} & -\alpha^{-1}a & 0 \\ 0 & \alpha^{-1} & 0 \\ 0 & -\alpha^{-1}b & I_{t \times t} \end{bmatrix} \begin{pmatrix} a \\ \alpha \\ b \end{pmatrix} = \begin{bmatrix} a - a \\ \alpha^{-1}\alpha \\ -b + b \end{bmatrix} = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} = e_{s+1}.$$

The matrix G is called a *Gaussian-Jordan Elimination Matrix*, or GJEM for short. Note that G is invertible since

$$G^{-1} = \begin{bmatrix} I & a & 0 \\ 0 & \alpha & 0 \\ 0 & b & I \end{bmatrix},$$

Moreover, for any vector of the form $w = \begin{pmatrix} x \\ 0 \\ y \end{pmatrix}$ where $x \in \mathbb{R}^s$ $y \in \mathbb{R}^t$, we have

$$Gw = w.$$

The GJEM matrices perform precisely the operations required in order to execute Gauss-Jordan elimination. That is, each elimination step can be realized as left multiplication of the augmented matrix by the appropriate GJEM.

For example, consider the linear system

$$\begin{aligned} 2x_1 + x_2 + 3x_3 &= 5 \\ 2x_1 + 2x_2 + 4x_3 &= 8 \\ 4x_1 + 2x_2 + 7x_3 &= 11 \\ 5x_1 + 3x_2 + 4x_3 &= 10 \end{aligned}$$

and its associated augmented matrix

$$A = \begin{bmatrix} 2 & 1 & 3 & 5 \\ 2 & 2 & 4 & 8 \\ 4 & 2 & 7 & 11 \\ 5 & 3 & 4 & 10 \end{bmatrix}.$$

The first step of Gauss-Jordan elimination is to transform the first column of this augmented matrix into the first unit coordinate vector. The procedure described in (7) can be employed for this purpose. In this case the pivot is the (1, 1) entry of the augmented matrix and so

$$s = 0, \text{ } a \text{ is void, } \alpha = 2, \text{ } t = 3, \text{ and } b = \begin{bmatrix} 2 \\ 4 \\ 5 \end{bmatrix},$$

which gives

$$G_1 = \begin{bmatrix} 1/2 & 0 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ -2 & 0 & 1 & 0 \\ -5/2 & 0 & 0 & 1 \end{bmatrix}.$$

Multiplying these two matrices gives

$$G_1 A = \begin{bmatrix} 1/2 & 0 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ -2 & 0 & 1 & 0 \\ -5/2 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 2 & 1 & 3 & 5 \\ 2 & 2 & 4 & 8 \\ 4 & 2 & 7 & 11 \\ 5 & 3 & 4 & 10 \end{bmatrix} = \begin{bmatrix} 1 & 1/2 & 3/2 & 5/2 \\ 0 & 1 & 1 & 3 \\ 0 & 0 & 1 & 1 \\ 0 & 1/2 & -7/2 & -5/2 \end{bmatrix}.$$

We now repeat this process to transform the second column of this matrix into the second unit coordinate vector. In this case the (2, 2) position becomes the pivot so that

$$s = 1, \text{ } a = 1/2, \text{ } \alpha = 1, \text{ } t = 2, \text{ and } b = \begin{bmatrix} 0 \\ 1/2 \end{bmatrix}$$

yielding

$$G_2 = \begin{bmatrix} 1 & -1/2 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & -1/2 & 0 & 1 \end{bmatrix}.$$

Again, multiplying these two matrices gives

$$G_2 G_1 A = \begin{bmatrix} 1 & -1/2 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & -1/2 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1/2 & 3/2 & 5/2 \\ 0 & 1 & 1 & 3 \\ 0 & 0 & 1 & 1 \\ 0 & 1/2 & -7/2 & -5/2 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 3 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & -4 & -4 \end{bmatrix}.$$

Repeating the process on the third column transforms it into the third unit coordinate vector. In this case the pivot is the (3, 3) entry so that

$$s = 2, a = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \alpha = 1, t = 1, \text{ and } b = -4$$

yielding

$$G_3 = \begin{bmatrix} 1 & 0 & -1 & 0 \\ 0 & 1 & -1 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 4 & 1 \end{bmatrix}.$$

Multiplying these matrices gives

$$G_3 G_2 G_1 A = \begin{bmatrix} 1 & 0 & -1 & 0 \\ 0 & 1 & -1 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 4 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 3 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & -4 & -4 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 2 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix},$$

which is in reduced echelon form. Therefore the system is consistent and the unique solution is

$$x = \begin{bmatrix} 0 \\ 2 \\ 1 \end{bmatrix}.$$

Observe that

$$G_3 G_2 G_1 = \begin{bmatrix} 3 & -1/2 & -1 & 0 \\ 1 & 1 & -1 & 0 \\ -2 & 0 & 1 & 0 \\ -10 & -1/2 & 4 & 1 \end{bmatrix}$$

and that

$$\begin{aligned} (G_3 G_2 G_1)^{-1} &= G_1^{-1} G_2^{-1} G_3^{-1} \\ &= \begin{bmatrix} 2 & 0 & 0 & 0 \\ 2 & 1 & 0 & 0 \\ 4 & 0 & 1 & 0 \\ 5 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1/2 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1/2 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & -4 & 1 \end{bmatrix} \\ &= \begin{bmatrix} 2 & 1 & 3 & 0 \\ 2 & 2 & 4 & 0 \\ 4 & 2 & 7 & 0 \\ 5 & 3 & 4 & 1 \end{bmatrix}. \end{aligned}$$

In particular, reduced Gauss-Jordan form can always be achieved by multiplying the augmented matrix on the left by an invertible matrix which can be written as a product of Gauss-Jordan elimination matrices.

EXERCISE 4.1. What are the Gauss-Jordan elimination matrices that transform the vector $\begin{bmatrix} 2 \\ 3 \\ -2 \\ 5 \end{bmatrix}$ in to e_j for

$j = 1, 2, 3, 4$, and what are the inverses of these matrices?

5. Some Special Square Matrices

We say that a matrix A is square if there is a positive integer n such that $A \in \mathbb{R}^{n \times n}$. For example, the Gauss-Jordan elimination matrices are a special kind of square matrix. Below we give a list of some square matrices with special properties that are very useful to our future work.

Diagonal Matrices: The diagonal of a matrix $A = [A_{ij}]$ is the vector $(A_{11}, A_{22}, \dots, A_{nn})^T \in \mathbb{R}^n$. A matrix in $\mathbb{R}^{n \times n}$ is said to be diagonal if the only non-zero entries of the matrix are the diagonal entries. Given a vector $v \in \mathbb{R}^n$, we write $\text{diag}(v)$ to denote the diagonal matrix whose diagonal is the vector v .

The Identity Matrix: The identity matrix is the diagonal matrix whose diagonal entries are all ones. We denote the identity matrix in \mathbb{R}^k by I_k . If the dimension of the identity is clear, we simply write I . Note that for any matrix $A \in \mathbb{R}^{m \times n}$ we have $I_m A = A = A I_n$.

Inverse Matrices: The inverse of a matrix $X \in \mathbb{R}^{n \times n}$ is any matrix $Y \in \mathbb{R}^{n \times n}$ such that $XY = I$ in which case we write $X^{-1} := Y$. It is easily shown that if Y is an inverse of X , then Y is unique and $YX = I$.

Permutation Matrices: A matrix $P \in \mathbb{R}^{n \times n}$ is said to be a permutation matrix if P is obtained from the identity matrix by either permuting the columns of the identity matrix or permuting its rows. It is easily seen that $P^{-1} = P^T$.

Unitary Matrices: A matrix $U \in \mathbb{R}^{n \times n}$ is said to be a unitary matrix if $U^T U = I$, that is $U^T = U^{-1}$. Note that every permutation matrix is unitary. But the converse is not true since for any vector u with $\|u\|_2 = 1$ the matrix $I - 2uu^T$ is unitary.

Symmetric Matrices: A matrix $M \in \mathbb{R}^{n \times n}$ is said to be symmetric if $M^T = M$.

Skew Symmetric Matrices: A matrix $M \in \mathbb{R}^{n \times n}$ is said to be skew symmetric if $M^T = -M$.

6. The LU Factorization

In this section we revisit the reduction to echelon form, but we incorporate permutation matrices into the pivoting process. Recall that a matrix $P \in \mathbb{R}^{m \times m}$ is a *permutation matrix* if it can be obtained from the identity matrix by permuting either its rows or columns. It is straightforward to show that $P^T P = I$ so that the inverse of a permutation matrix is its transpose. Multiplication of a matrix on the left permutes the rows of the matrix while multiplication on the right permutes the columns. We now apply permutation matrices in the Gaussian elimination process in order to avoid zero pivots.

Let $A \in \mathbb{R}^{m \times n}$ and assume that $A \neq 0$. Set $\tilde{A}_0 := A$. If the $(1, 1)$ entry of \tilde{A}_0 is zero, then apply permutation matrices P_{l0} and P_{r0} to the left and right of \tilde{A}_0 , respectively, to bring *any* non-zero element of \tilde{A}_0 into the $(1, 1)$ position (e.g., the one with largest magnitude) and set $A_0 := P_{l0} \tilde{A}_0 P_{r0}$. Write A_0 in block form as

$$A_0 = \begin{bmatrix} \alpha_1 & v_1^T \\ u_1 & \tilde{A}_1 \end{bmatrix} \in \mathbb{R}^{m \times n},$$

with $0 \neq \alpha_1 \in \mathbb{R}$, $u_1 \in \mathbb{R}^{n-1}$, $v_1 \in \mathbb{R}^{m-1}$, and $\tilde{A}_1 \in \mathbb{R}^{(m-1) \times (n-1)}$. Then using α_1 to zero out u_1 amounts to left multiplication of the matrix A_0 by the Gaussian elimination matrix

$$\begin{bmatrix} 1 & 0 \\ -\frac{u_1}{\alpha_1} & I \end{bmatrix}$$

to get

$$(8) \quad \begin{bmatrix} 1 & 0 \\ -\frac{u_1}{\alpha_1} & I \end{bmatrix} \begin{bmatrix} \alpha_1 & v_1^T \\ u_1 & \tilde{A}_1 \end{bmatrix} = \begin{bmatrix} \alpha_1 & v_1^T \\ 0 & \tilde{A}_1 \end{bmatrix} \in \mathbb{R}^{m \times n},$$

where

$$\tilde{A}_1 = \tilde{A}_1 - u_1 v_1^T / \alpha_1.$$

Define

$$\tilde{L}_1 = \begin{bmatrix} 1 & 0 \\ \frac{u_1}{\alpha_1} & I \end{bmatrix} \in \mathbb{R}^{m \times m} \quad \text{and} \quad \tilde{U}_1 = \begin{bmatrix} \alpha_1 & v_1^T \\ 0 & \tilde{A}_1 \end{bmatrix} \in \mathbb{R}^{m \times n}.$$

and observe that

$$\tilde{L}_1^{-1} = \begin{bmatrix} 1 & 0 \\ -\frac{u_1}{\alpha_1} & I \end{bmatrix}$$

Hence (8) becomes

$$(9) \quad \tilde{L}_1^{-1} P_{l0} \tilde{A}_0 P_{r0} = \tilde{U}_1, \quad \text{or equivalently,} \quad A = P_{l0}^T \tilde{L}_1 \tilde{U}_1 P_{r0}^T.$$

Note that \tilde{L}_1 is *unit* lower triangular (ones on the main diagonal) and \tilde{U}_1 is block upper-triangular with one nonsingular 1×1 block and one $(m-1) \times (n-1)$ block on the block diagonal.

Next consider the matrix \tilde{A}_1 in \tilde{U}_1 . If the $(1, 1)$ entry of \tilde{A}_1 is zero, then apply permutation matrices $\tilde{P}_{l1} \in \mathbb{R}^{(m-1) \times (m-1)}$ and $\tilde{P}_{r1} \in \mathbb{R}^{(n-1) \times (n-1)}$ to the left and right of $\tilde{A}_1 \in \mathbb{R}^{(m-1) \times (n-1)}$, respectively, to bring *any* non-zero element of \tilde{A}_1 into the $(1, 1)$ position (e.g., the one with largest magnitude) and set $A_1 := \tilde{P}_{l1} \tilde{A}_1 \tilde{P}_{r1}$. If the element of \tilde{A}_1 is zero, then stop. Define

$$P_{l1} := \begin{bmatrix} 1 & 0 \\ 0 & \tilde{P}_{l1} \end{bmatrix} \quad \text{and} \quad P_{r1} := \begin{bmatrix} 1 & 0 \\ 0 & \tilde{P}_{r1} \end{bmatrix}$$

so that P_{l1} and P_{r1} are also permutation matrices and

$$(10) \quad P_{l1}\tilde{U}_1P_{r1} = \begin{bmatrix} 1 & 0 \\ 0 & \tilde{P}_{l1} \end{bmatrix} \begin{bmatrix} \alpha_1 & v_1^T \\ 0 & \tilde{A}_1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & \tilde{P}_{r1} \end{bmatrix} = \begin{bmatrix} \alpha_1 & v_1^T P_{r1} \\ 0 & \tilde{P}_{l1}\tilde{A}_1P_{r1} \end{bmatrix} = \begin{bmatrix} \alpha_1 & \tilde{v}_1^T \\ 0 & \tilde{A}_1 \end{bmatrix},$$

where $\tilde{v}_1 := P_{r1}^T v_1$. Define

$$U_1 := \begin{bmatrix} \alpha_1 & \tilde{v}_1^T \\ 0 & \tilde{A}_1 \end{bmatrix}, \quad \text{where } A_1 = \begin{bmatrix} \alpha_2 & v_2^T \\ u_2 & \tilde{A}_2 \end{bmatrix} \in \mathbb{R}^{(m-1) \times (n-1)},$$

with $0 \neq \alpha_2 \in \mathbb{R}$, $u_2 \in \mathbb{R}^{n-2}$, $v_2 \in \mathbb{R}^{m-2}$, and $\tilde{A}_2 \in \mathbb{R}^{(m-2) \times (n-2)}$. In addition, define

$$L_1 := \begin{bmatrix} 1 & 0 \\ \tilde{P}_{l1} \frac{u_1}{\alpha_1} & I \end{bmatrix},$$

so that

$$\begin{aligned} P_{l1}^T L_1 &= \begin{bmatrix} 1 & 0 \\ 0 & \tilde{P}_{l1}^T \end{bmatrix} \begin{bmatrix} 1 & 0 \\ \tilde{P}_{l1} \frac{u_1}{\alpha_1} & I \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 \\ \frac{u_1}{\alpha_1} & \tilde{P}_{l1}^T \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 \\ \frac{u_1}{\alpha_1} & I \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & \tilde{P}_{l1}^T \end{bmatrix} \\ &= \tilde{L}_1 P_{l1}^T, \end{aligned}$$

and consequently

$$L_1^{-1} P_{l1} = P_{l1} \tilde{L}_1^{-1}.$$

Plugging this into (9) and using (10), we obtain

$$L_1^{-1} P_{l1} P_{l0} \tilde{A}_0 P_{r0} P_{r1} = P_{l1} \tilde{L}_1^{-1} P_{l0} \tilde{A}_0 P_{r0} P_{r1} = P_{l1} \tilde{U}_1 P_{r1} = U_1,$$

or equivalently,

$$P_{l1} P_{l0} A P_{r0} P_{r1} = L_1 U_1.$$

We can now repeat this process on the matrix A_1 since the (1,1) entry of this matrix is non-zero. The process can run for no more than the number of rows of A which is m . However, it may terminate after $k < m$ steps if the matrix \hat{A}_k is the zero matrix. In either event, we obtain the following result.

THEOREM 6.1. *[The LU Factorization] Let $A \in \mathbb{R}^{m \times n}$. If $k = \text{rank}(A)$, then there exist permutation matrices $P_l \in \mathbb{R}^{m \times m}$ and $P_r \in \mathbb{R}^{n \times n}$ such that*

$$P_l A P_r = LU,$$

where $L \in \mathbb{R}^{m \times m}$ is a lower triangular matrix having ones on its diagonal and

$$U = \begin{bmatrix} U_1 & U_2 \\ 0 & 0 \end{bmatrix}$$

with $U_1 \in \mathbb{R}^{k \times k}$ a nonsingular upper triangular matrix.

Note that a column permutation is only required if the first column of \hat{A}_k is zero for some k before termination. In particular, this implies that the $\text{rank}(A) < m$. Therefore, if $\text{rank}(A) = m$, column permutations are not required, and $P_r = I$. If one implements the LU factorization so that a column permutation is *only* employed in the case when the first column of \hat{A}_k is zero for some k , then we say the LU factorization is obtained through partial pivoting.

EXAMPLE 6.1. *We now use the procedure outlined above to compute the LU factorization of the matrix*

$$A = \begin{bmatrix} 1 & 1 & 2 \\ 2 & 4 & 2 \\ -1 & 1 & 3 \end{bmatrix}.$$

$$\begin{aligned} L_1^{-1}A &= \begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 2 \\ 2 & 4 & 2 \\ -1 & 1 & 3 \end{bmatrix} \\ &= \begin{bmatrix} 1 & 1 & 2 \\ 0 & 2 & -3 \\ 0 & 2 & 5 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} L_2^{-1}L_1^{-1}A &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 2 \\ 0 & 2 & -3 \\ 0 & 2 & 5 \end{bmatrix} \\ &= \begin{bmatrix} 1 & 1 & 2 \\ 0 & 2 & -3 \\ 0 & 0 & 8 \end{bmatrix} \end{aligned}$$

We now have

$$U = \begin{bmatrix} 1 & 1 & 2 \\ 0 & 2 & -3 \\ 0 & 0 & 8 \end{bmatrix},$$

and

$$L = L_1L_2 = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ -1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ -1 & 1 & 1 \end{bmatrix}.$$

7. Solving Equations with the LU Factorization

Consider the equation $Ax = b$. In this section we show how to solve this equation using the LU factorization. Recall from Theorem 6.1 that the algorithm of the previous section produces a factorization of A of the form $P_l \in \mathbb{R}^{m \times m}$ and $P_r \in \mathbb{R}^{n \times n}$ such that

$$A = P_l^T L U P_r^T,$$

where $P_l \in \mathbb{R}^{m \times m}$ and $P_r \in \mathbb{R}^{n \times n}$ are permutation matrices, $L \in \mathbb{R}^{m \times m}$ is a lower triangular matrix having ones on its diagonal, and

$$U = \begin{bmatrix} U_1 & U_2 \\ 0 & 0 \end{bmatrix}$$

with $U_1 \in \mathbb{R}^{k \times k}$ a nonsingular upper triangular matrix. Hence we may write the equation $Ax = b$ as

$$P_l^T L U P_r^T x = b.$$

Multiplying through by P_l and replacing $U P_r^T x$ by w gives the equation

$$Lw = \hat{b}, \quad \text{where } \hat{b} := P_l b.$$

This equation is easily solved by forward substitution since L is a nonsingular lower triangular matrix. Denote the solution by \bar{w} . To obtain a solution x we must still solve $U P_r^T x = \bar{w}$. Set $y = P_r x$. Then this equation becomes

$$\bar{w} = U y = \begin{bmatrix} U_1 & U_2 \\ 0 & 0 \end{bmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix},$$

where we have decomposed y to conform to the decomposition of U . Doing the same for \bar{w} gives

$$\begin{pmatrix} \bar{w}_1 \\ \bar{w}_2 \end{pmatrix} = \begin{bmatrix} U_1 & U_2 \\ 0 & 0 \end{bmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix},$$

or equivalently,

$$\begin{aligned} \bar{w}_1 &= U_1 y_1 + U_2 y_2 \\ \bar{w}_2 &= 0. \end{aligned}$$

Hence, if $\bar{w}_2 \neq 0$, the system is inconsistent, i.e., no solution exists. On the other hand, if $\bar{w}_2 = 0$, we can take $y_2 = 0$ and solve the equation

$$(11) \quad \bar{w}_1 = U_1 y_1$$

for \bar{y}_1 , then

$$\bar{x} = P_r^T \begin{pmatrix} \bar{y}_1 \\ 0 \end{pmatrix}$$

is a solution to $Ax = b$. The equation (11) is also easy to solve since U_1 is an upper triangular nonsingular matrix so that (11) can be solved by back substitution.

8. The Four Fundamental Subspaces and Echelon Form

Recall that a subset W of \mathbb{R}^n is a subspace if and only if it satisfies the following three conditions:

- (1) The origin is an element of W .
- (2) The set W is closed with respect to addition, i.e. if $u \in W$ and $v \in W$, then $u + v \in W$.
- (3) The set W is closed with respect to scalar multiplication, i.e. if $\alpha \in \mathbb{R}$ and $u \in W$, then $\alpha u \in W$.

EXERCISE 8.1. Given $v^1, v^2, \dots, v^k \in \mathbb{R}^n$, show that the linear span of these vectors,

$$\text{span}(v^1, v^2, \dots, v^k) := \{\alpha_1 v^1 + \alpha_2 v^2 + \dots + \alpha_k v^k \mid \alpha_1, \alpha_2, \dots, \alpha_k \in \mathbb{R}\}$$

is a subspace.

EXERCISE 8.2. Show that for any set S in \mathbb{R}^n , the set

$$S^\perp = \{v : w^T v = 0 \text{ for all } w \in S\}$$

is a subspace. If S is itself a subspace, then S^\perp is called the subspace orthogonal (or perpendicular) to the subspace S .

EXERCISE 8.3. If S is any subset of \mathbb{R}^n (not necessarily a subspace), show that $(S^\perp)^\perp = \text{span}(S)$.

EXERCISE 8.4. If $S \subset \mathbb{R}^n$ is a subspace, show that $S = (S^\perp)^\perp$.

A set of vectors $v^1, v^2, \dots, v^k \in \mathbb{R}^n$ are said to be *linearly independent* if $0 = a_1 v^1 + \dots + a_k v^k$ if and only if $0 = a_1 = a_2 = \dots = a_k$. A *basis* for a subspace is any maximal linearly independent set. An elementary fact from linear algebra is that the subspace equals the linear span of any basis for the subspace and that every basis of a subspace has the same number of vectors in it. We call this number the *dimension* for the subspace. If S is a subspace, we denote the dimension of S by $\dim S$.

EXERCISE 8.5. If $S \subset \mathbb{R}^n$ is a subspace, then any basis of S can contain only finitely many vectors.

EXERCISE 8.6. Show that every subspace can be represented as the linear span of a basis for that subspace.

EXERCISE 8.7. Show that every basis for a subspace contains the same number of vectors.

EXERCISE 8.8. If $S \subset \mathbb{R}^n$ is a subspace, show that

$$(12) \quad \mathbb{R}^n = S + S^\perp$$

and that

$$(13) \quad n = \dim S + \dim S^\perp.$$

Let $A \in \mathbb{R}^{m \times n}$. We associate with A its four fundamental subspaces:

$$\begin{aligned} \text{Ran}(A) &:= \{Ax \mid x \in \mathbb{R}^n\} & \text{Null}(A) &:= \{x \mid Ax = 0\} \\ \text{Ran}(A^T) &:= \{A^T y \mid y \in \mathbb{R}^m\} & \text{Null}(A^T) &:= \{y \mid A^T y = 0\}. \end{aligned}$$

where

$$(14) \quad \begin{aligned} \text{rank}(A) &:= \dim \text{Ran}(A) & \text{nullity}(A) &:= \dim \text{Null}(A) \\ \text{rank}(A^T) &:= \dim \text{Ran}(A^T) & \text{nullity}(A^T) &:= \dim \text{Null}(A^T) \end{aligned}$$

EXERCISE 8.9. Show that the four fundamental subspaces associated with a matrix are indeed subspaces.

Observe that

$$\begin{aligned}
 \text{Null}(A) &:= \{x \mid Ax = 0\} \\
 &= \{x \mid A_i \bullet x = 0, i = 1, 2, \dots, m\} \\
 &= \{A_{1\cdot}, A_{2\cdot}, \dots, A_{m\cdot}\}^\perp \\
 &= \text{span}(A_{1\cdot}, A_{2\cdot}, \dots, A_{m\cdot})^\perp \\
 &= \text{Ran}(A^T)^\perp.
 \end{aligned}$$

Since for any subspace $S \subset \mathbb{R}^n$, we have $(S^\perp)^\perp = S$, we obtain

$$(15) \quad \text{Null}(A)^\perp = \text{Ran}(A^T) \text{ and } \text{Null}(A^T) = \text{Ran}(A)^\perp.$$

The equivalences in (15) are called the *Fundamental Theorem of the Alternative*.

One of the big consequences of echelon form is that

$$(16) \quad n = \text{rank}(A) + \text{nullity}(A).$$

By combining (16), (13) and (15), we obtain the equivalence

$$\text{rank}(A^T) = \dim \text{Ran}(A^T) = \dim \text{Null}(A)^\perp = n - \text{nullity}(A) = \text{rank}(A).$$

That is, the row rank of a matrix equals the column rank of a matrix, i.e., the dimensions of the row and column spaces of a matrix are the same!

The Linear Least Squares Problem

In this chapter we study the linear least squares problem introduced in (4). Since this is such a huge and important topic, we will only be able to briefly touch on a few aspects of this problem. But our introduction should give the reader some idea of the scope of this topic and an indication of a few current areas of research.

1. Applications

1.1. Polynomial Fitting. In many data fitting applications, one assumes a functional relationship between a set of “inputs” and a set of “outputs”. For example, a patient is injected with a drug and the research wishes to understand the clearance of the drug as a function of time. One way to do this is to draw blood samples over time and to measure the concentration of the drug in the drawn serum. The goal is to then provide a functional description of the concentration at any point in time.

Suppose the observed data is $y_i \in \mathbb{R}$ for each time point t_i , $i = 1, 2, \dots, N$, respectively. The underlying assumption is that there is some function of time $f : \mathbb{R} \rightarrow \mathbb{R}$ such that $y_i = f(t_i)$, $i = 1, 2, \dots, N$. The goal is to provide and estimate of the function f . One way to do this is to try to approximate f by a polynomial of a fixed degree, say n :

$$p(t) = x_0 + x_1 t + x_2 t^2 + \dots + x_n t^n.$$

We now wish to determine the values of the coefficients that “best” fit the data.

If were possible to exactly fit the data, then there would exist a value for the coefficient, say $\bar{x} = (\bar{x}_0, \bar{x}_1, \bar{x}_2, \dots, \bar{x}_n)$ such that

$$y_i = \bar{x}_0 + \bar{x}_1 t_i + \bar{x}_2 t_i^2 + \dots + \bar{x}_n t_i^n, \quad i = 1, 2, \dots, N.$$

But if N is larger than n , then it is unlikely that such an \bar{x} exists; while if N is less than n , then there are probably many choices for \bar{x} for which we can achieve a perfect fit. We discuss these two scenarios and their consequences in more depth at a future date, but, for the moment, we assume that N is larger than n . That is, we wish to approximate f with a low degree polynomial.

When $n \ll N$, we cannot expect to fit the data perfectly and so there will be errors. In this case, we must come up with a notion of what it means to “best” fit the data. In the context of least squares, “best” means that we wish to minimize the sum of the squares of the errors in the fit:

$$(17) \quad \underset{x \in \mathbb{R}^{n+1}}{\text{minimize}} \frac{1}{2} \sum_{i=1}^N (x_0 + x_1 t_i + x_2 t_i^2 + \dots + x_n t_i^n - y_i)^2.$$

The leading one half in the objective is used to simplify certain computations that occur in the analysis to come. This minimization problem has the form

$$\underset{x \in \mathbb{R}^{n+1}}{\text{minimize}} \frac{1}{2} \|Vx - y\|_2^2,$$

where

$$y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}, \quad x = \begin{pmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \quad \text{and} \quad V = \begin{bmatrix} 1 & t_1 & t_1^2 & \dots & t_1^n \\ 1 & t_2 & t_2^2 & \dots & t_2^n \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & t_N & t_N^2 & \dots & t_N^n \end{bmatrix},$$

since

$$Vx = \begin{pmatrix} x_0 + x_1 t_1 + x_2 t_1^2 + \cdots + x_n t_1^n \\ x_0 + x_1 t_2 + x_2 t_2^2 + \cdots + x_n t_2^n \\ \vdots \\ x_0 + x_1 t_N + x_2 t_N^2 + \cdots + x_n t_N^n \end{pmatrix}.$$

That is, the polynomial fitting problem (17) is an example of a linear least squares problem (4). The matrix V is called the *Vandermonde matrix* associated with this problem.

This is a neat way to approximate functions. However, polynomials are a very poor way to approximate the clearance data discussed in our motivation to this approach. The concentration of a drug in serum typically rises quickly after injection to a maximum concentration and falls off gradually decaying exponentially. There is only one place where such a function is zero, and this occurs at time zero. On the other hand, a polynomial of degree n has n zeros (counting multiplicity). Therefore, it would seem that exponential functions would provide a better basis for estimating clearance. This motivates our next application.

1.2. Function Approximation by Bases Functions. In this application we expand on the basic ideas behind polynomial fitting to allow other kinds of approximations, such as approximation by sums of exponential functions. In general, suppose we are given data points $(z_i, y_i) \in \mathbb{R}^2$, $i = 1, 2, \dots, N$ where it is assumed that the observation y_i is a function of an unknown function $f: \mathbb{R} \rightarrow \mathbb{R}$ evaluated at the point z_i for each $i = 1, 2, \dots, N$. Based on other aspects of the underlying setting from which this data arises may lead us to believe that f comes from a certain space \mathcal{F} of functions, such as the space of continuous or differentiable functions on an interval. This space of functions may itself be a vector space in the sense that the zero function is in the space ($0 \in \mathcal{F}$), two functions in the space can be added pointwise to obtain another function in the space (\mathcal{F} is closed with respect to addition), and any real multiple of a function is in the space (\mathcal{F} is closed with respect to scalar multiplication). In this case, we may select from X a finite subset of functions, say $\phi_1, \phi_2, \dots, \phi_k$, and try to approximate f as a linear combination of these functions:

$$f(x) \sim x_1 \phi_1(z) + x_2 \phi_2(z) + \cdots + x_n \phi_k(z).$$

This is exactly what we did in the polynomial fitting application discussed above. There $\phi_i(z) = z^i$ but we started the indexing at $i = 0$. Therefore, this idea is essentially the same as the polynomial fitting case. But the functions z^i have additional properties. First, they are linearly independent in the sense that the only linear combination that yields the zero function is the one where all of the coefficients are zero. In addition, any continuous function on an interval can be approximated “arbitrarily well” by a polynomial assuming that we allow the polynomials to be of arbitrarily high degree (think Taylor approximations). In this sense, polynomials form a basis for the continuous function on an interval. By analogy, we would like our functions ϕ_i to be linearly independent and to come from basis of functions. There are many possible choices of bases, but a discussion of these would take us too far afield from this course.

Let now suppose that the functions $\phi_1, \phi_2, \dots, \phi_k$ are linearly independent and arise from a set of basis function that reflect a deeper intuition about the behavior of the function f , e.g. it is well approximated as a sum of exponentials (or trig functions). Then the task to find those coefficient x_1, x_2, \dots, x_n that best fits the data in the least squares sense:

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \frac{1}{2} \sum_{i=1}^N (x_1 \phi_1(z_i) + x_2 \phi_2(z_i) + \cdots + x_n \phi_k(z_i) - y_i)^2.$$

This can be recast as the linear least squares problem

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \frac{1}{2} \|Ax - y\|_2^2,$$

where

$$y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}, \quad x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \quad \text{and} \quad A = \begin{bmatrix} \phi_1(z_1) & \phi_2(z_1) & \cdots & \phi_n(z_1) \\ \phi_1(z_2) & \phi_2(z_2) & \cdots & \phi_n(z_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_1(z_N) & \phi_2(z_N) & \cdots & \phi_n(z_N) \end{bmatrix}.$$

Many possible further generalizations of this basic idea are possible. For example, the data may be multi-dimensional: $(z_i, y_i) \in \mathbb{R}^s \times \mathbb{R}^t$. In addition, constraints may be added, e.g., the function must be monotone (either increasing or decreasing), it must be unimodal (one “bump”), etc. But the essential features are that we estimate

using linear combinations and errors are measured using sums of squares. In many cases, the sum of squares error metric is not a good choice. But it can be motivated by assuming that the error are distributed using the Gaussian, or normal, distribution.

1.3. Linear Regression and Maximum Likelihood. Suppose we are considering a new drug therapy for reducing inflammation in a targeted population, and we have a relatively precise way of measuring inflammation for each member of this population. We are trying to determine the dosing to achieve a target level of inflammation. Of course, the dose needs to be adjusted for each individual due to the great amount of variability from one individual to the next. One way to model this is to assume that the resultant level of inflammation is on average a linear function of the dose and other individual specific covariates such as sex, age, weight, body surface area, gender, race, blood iron levels, disease state, etc. We then sample a collection of N individuals from the target population, register their dose z_{i0} and the values of their individual specific covariates $z_{i1}, z_{i2}, \dots, z_{in}$, $i = 1, 2, \dots, N$. After dosing we observe that the resultant inflammation for the i th subject to be y_i , $i = 1, 2, \dots, N$. By saying that the “resultant level of inflammation is on average a linear function of the dose and other individual specific covariates”, we mean that there exist coefficients $x_0, x_1, x_2, \dots, x_n$ such that

$$y_i = x_0 z_{i0} + x_1 z_{i1} + x_2 z_{i2} + \dots + x_n z_{in} + v_i,$$

where v_i is an instance of a random variable representing the individual's deviation from the linear model. Assume that the random variables v_i are independently identically distributed $N(0, \sigma^2)$ (norm with zero mean and variance σ^2). The probability density function for the normal distribution $N(0, \sigma^2)$ is

$$\frac{1}{\sigma\sqrt{2\pi}} \text{EXP}[-v^2/(2\sigma^2)].$$

Given values for the coefficients x_i , the likelihood function for the sample y_i , $i = 1, 2, \dots, N$ is the joint probability density function evaluated at this observation. The independence assumption tells us that this joint pdf is given by

$$L(x; y) = \left(\frac{1}{\sigma\sqrt{2\pi}} \right)^n \text{EXP} \left[-\frac{1}{2\sigma^2} \sum_{i=1}^N (x_0 z_{i0} + x_1 z_{i1} + x_2 z_{i2} + \dots + x_n z_{in} - y_i)^2 \right].$$

We now wish to choose those values of the coefficients x_0, x_2, \dots, x_n that make the observation y_1, y_2, \dots, y_n most probable. One way to try to do this is to maximize the likelihood function $L(x; y)$ over all possible values of x . This is called *maximum likelihood estimation*:

$$(18) \quad \underset{x \in \mathbb{R}^{n+1}}{\text{maximize}} L(x; y).$$

Since the natural logarithm is nondecreasing on the range of the likelihood function, the problem (18) is equivalent to the problem

$$\underset{x \in \mathbb{R}^{n+1}}{\text{maximize}} \ln(L(x; y)),$$

which in turn is equivalent to the minimization problem

$$(19) \quad \underset{x \in \mathbb{R}^{n+1}}{\text{minimize}} -\ln(L(x; y)).$$

Finally, observe that

$$-\ln(L(x; y)) = K + \frac{1}{2\sigma^2} \sum_{i=1}^N (x_0 z_{i0} + x_1 z_{i1} + x_2 z_{i2} + \dots + x_n z_{in} - y_i)^2,$$

where $K = n \ln(\sigma\sqrt{2\pi})$ is constant. Hence the problem (19) is equivalent to the linear least squares problem

$$\underset{x \in \mathbb{R}^{n+1}}{\text{minimize}} \frac{1}{2} \|Ax - y\|_2^2,$$

where

$$y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}, \quad x = \begin{pmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \quad \text{and} \quad A = \begin{bmatrix} z_{10} & z_{11} & z_{12} & \dots & z_{1n} \\ z_{20} & z_{21} & z_{22} & \dots & z_{2n} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ z_{N0} & z_{N1} & z_{N2} & \dots & z_{Nn} \end{bmatrix}.$$

This is the first step in trying to select an optimal dose for each individual across a target population. What is missing from this analysis is some estimation of the variability in inflammation response due to changes in the covariates. Understanding this sensitivity to variations in the covariates is an essential part of any regression analysis. However, a discussion of this step lies beyond the scope of this brief introduction to linear regression.

1.4. System Identification in Signal Processing. We consider a standard problem in signal processing concerning the behavior of a stable, causal, linear, continuous-time, time-invariant system with input signal $u(t)$ and output signal $y(t)$. Assume that these signals can be described by the convolution integral

$$(20) \quad y(t) = (g * u)(t) := \int_0^{+\infty} g(\tau)u(t - \tau)d\tau .$$

In applications, the goal is to obtain an estimate of g by observing outputs y from a variety of known input signals u . For example, returning to our drug dosing example, the function u may represent the input of a drug into the body through a drug pump any y represent the concentration of the drug in the body at any time t . The relationship between the two is clearly causal (and can be shown to be stable). The transfer function g represents what the body is doing to the drug. In the way, the model (20) is a common model used in pharmaco-kinetics.

The problem of estimating g in (20) is an infinite dimensional problem. Below we describe a way to approximate g using the the FIR, or *finite impulse response* filter. In this model we discretize time by choosing a fixed number N of time points t_i to observe y from a known input u , and a finite time horizon $n < N$ over which to approximate the integral in (20). To simplify matters we index time on the integers, that is, we equate t_i with the integer i . After selecting the data points and the time horizon, we obtain the FIR model

$$(21) \quad y(t) = \sum_{k=1}^n g(k)u(t - k),$$

where we try to find the “best” values for $g(k)$, $k = 0, 1, 2, \dots, n$ to fit the system

$$y(t) = \sum_{k=0}^n g(k)u(t - k), \quad t = 1, 2, \dots, N.$$

Notice that this requires knowledge of the values $u(t)$ for $t = 1 - n, 2 - n, \dots, N$. One often assumes a observational error in this model that is $N(0, \sigma^2)$ for a given value of σ^2 . In this case, the FIR model (21) becomes

$$(22) \quad y(t) = \sum_{k=1}^n g(k)u(t - k) + v(t),$$

where $v(t)$, $t = 1, \dots, N$ are iid $N(0, \sigma^2)$. In this case, the corresponding maximum likelihood estimation problem becomes the linear least squares problem

$$\underset{g \in \mathbb{R}^{n+1}}{\text{minimize}} \frac{1}{2} \|Hg - y\|_2^2,$$

where

$$y = \begin{pmatrix} y(1) \\ y(2) \\ \vdots \\ y(N) \end{pmatrix}, \quad g = \begin{pmatrix} g(0) \\ g(1) \\ g(2) \\ \vdots \\ g(n) \end{pmatrix} \quad \text{and} \quad H = \begin{bmatrix} u(1) & u(0) & u(-1) & u(-2) & \dots & u(1-n) \\ u(2) & u(1) & u(0) & u(-1) & \dots & u(2-n) \\ u(3) & u(2) & u(1) & u(0) & \dots & u(3-n) \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ u(N) & u(N-1) & u(N-2) & u(N-3) & \dots & u(N-n) \end{bmatrix}.$$

Notice that the matrix H has constant “diagonals”. Such matrices are called *Toeplitz matrices*.

1.5. Kalman Smoothing. Kalman smoothing is a fundamental topic in signal processing and control literature, with numerous applications in navigation, tracking, healthcare, finance, and weather. Contributions to theory and algorithms related to Kalman smoothing, and to dynamic system inference in general, have come from statistics, engineering, numerical analysis, and optimization. Here, the term ‘Kalman smoother’ includes any method of inference on any dynamical system fitting the graphical representation of Figure 1.

The combined mathematical, statistical, and probabilistic model corresponding to Figure 1 is specified as follows:

$$(23) \quad \begin{aligned} \mathbf{x}_1 &= g_1(x_0) + \mathbf{w}_1, \\ \mathbf{x}_k &= g_k(\mathbf{x}_{k-1}) + \mathbf{w}_k \quad k = 2, \dots, N, \\ \mathbf{z}_k &= h_k(\mathbf{x}_k) + \mathbf{v}_k \quad k = 1, \dots, N, \end{aligned}$$

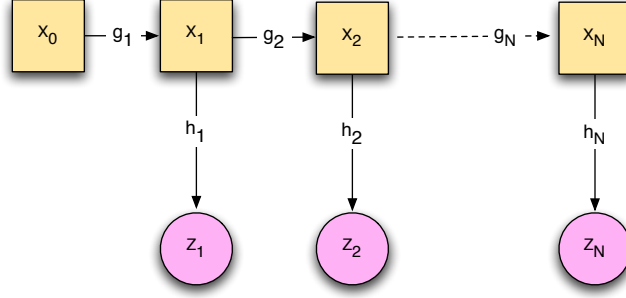


FIGURE 1. Dynamic systems amenable to Kalman smoothing methods.

where \mathbf{w}_k , \mathbf{v}_k are mutually independent random variables with known positive definite covariance matrices Q_k and R_k , respectively. Here, \mathbf{w}_k often, but not always, arises from a probabilistic model (discretization of an underlying stochastic differential equation) and \mathbf{v}_k comes from a statistical model for observations. We have $\mathbf{x}_k, \mathbf{w}_k \in \mathbb{R}^n$, and $\mathbf{z}_k, \mathbf{v}_k \in \mathbb{R}^{m(k)}$, so dimensions can vary between time points. Here the sequence $\{\mathbf{x}_k\}$ is called the state-space sequence and $\{\mathbf{z}_k\}$ is the observation sequence. The functions g_k and h_k as well as the matrices Q_k and R_k are known and given. In addition, the observation sequence $\{\mathbf{z}_k\}$ is also known. The goal is to estimate the unobserved state sequence $\{\mathbf{x}_k\}$. For example, in our drug dosing example, the amount of the drug remaining in the body at time t is the unknown state sequence while the observation sequence is the observed concentration of the drug in each of our blood draws.

The classic case is obtained by making the following assumptions:

- (1) x_0 is known, and g_k , h_k are known *linear* functions, which we denote by

$$(24) \quad g_k(x_{k-1}) = G_k x_{k-1} \quad h_k(x_k) = H_k x_k$$

where $G_k \in \mathbb{R}^{n \times n}$ and $H_k \in \mathbb{R}^{m(k) \times n}$,

- (2) \mathbf{w}_k , \mathbf{v}_k are mutually independent *Gaussian* random variables.

In the classical setting, the connection to the linear least squares problem is obtained by formulating the maximum *a posteriori* (MAP) problem under linear and Gaussian assumptions. As in the linear regression and signal processing applications, this yields the following linear least squares problem:

$$(25) \quad \min_{\{x_k\}} f(\{x_k\}) := \sum_{k=1}^N \frac{1}{2} (z_k - H_k x_k)^T R_k^{-1} (z_k - H_k x_k) + \frac{1}{2} (x_k - G_k x_{k-1})^T Q_k^{-1} (x_k - G_k x_{k-1}).$$

To simplify this expression, we introduce data structures that capture the entire state sequence, measurement sequence, covariance matrices, and initial conditions. Given a sequence of column vectors $\{u_k\}$ and matrices $\{T_k\}$ we use the notation

$$\text{vec}(\{u_k\}) = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_N \end{bmatrix}, \quad \text{diag}(\{T_k\}) = \begin{bmatrix} T_1 & 0 & \cdots & 0 \\ 0 & T_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & T_N \end{bmatrix}.$$

We now make the following definitions:

$$(26) \quad \begin{aligned} R &= \text{diag}(\{R_k\}) & x &= \text{vec}(\{x_k\}) \\ Q &= \text{diag}(\{Q_k\}) & w &= \text{vec}(\{g_0, 0, \dots, 0\}) \\ H &= \text{diag}(\{H_k\}) & z &= \text{vec}(\{z_1, z_2, \dots, z_N\}) \end{aligned} \quad G = \begin{bmatrix} I & 0 & & \\ -G_2 & I & \ddots & \\ & \ddots & \ddots & 0 \\ & & & -G_N & I \end{bmatrix},$$

where $g_0 := g_1(x_0) = G_1 x_0$. With definitions in (26), problem (25) can be written

$$(27) \quad \min_x f(x) = \frac{1}{2} \|Hx - z\|_{R^{-1}}^2 + \frac{1}{2} \|Gx - w\|_{Q^{-1}}^2,$$

where $\|a\|_M^2 = a^\top M a$.

Since the number of time steps N can be quite large, it is essential that the underlying tri-diagonal structure is exploited in any solution procedure. This is especially true when the state-space dimension n is also large which occurs when making PET scan movies of brain metabolics or reconstructing weather patterns on a global scale.

2. Optimality in the Linear Least Squares Problem

We now turn to a discussion of optimality in the least squares problem (4) which we restate here for ease of reference:

$$(28) \quad \underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2} \|Ax - b\|_2^2,$$

where

$$A \in \mathbb{R}^{m \times n}, \quad b \in \mathbb{R}^m, \quad \text{and} \quad \|y\|_2^2 := y_1^2 + y_2^2 + \cdots + y_m^2.$$

In particular, we will address the question of when a solution to this problem exists and how they can be identified or characterized.

Suppose that \bar{x} is a solution to (28), i.e.,

$$(29) \quad \|A\bar{x} - b\|_2 \leq \|Ax - b\|_2 \quad \forall x \in \mathbb{R}^n.$$

Using this inequality, we derive necessary and sufficient conditions for the optimality of \bar{x} . A useful identity for our derivation is

$$(30) \quad \|u + v\|_2^2 = (u + v)^\top (u + v) = u^\top u + 2u^\top v + v^\top v = \|u\|_2^2 + 2u^\top v + \|v\|_2^2.$$

Let x be any other vector in \mathbb{R}^n . Then, using (30) with $u = A(\bar{x} - x)$ and $v = Ax - b$ we obtain

$$(31) \quad \begin{aligned} \|A\bar{x} - b\|_2^2 &= \|A(\bar{x} - x) + (Ax - b)\|_2^2 \\ &= \|A(\bar{x} - x)\|_2^2 + 2(A(\bar{x} - x))^\top (Ax - b) + \|Ax - b\|_2^2 \\ &\geq \|A(\bar{x} - x)\|_2^2 + 2(A(\bar{x} - x))^\top (Ax - b) + \|A\bar{x} - b\|_2^2 \quad (\text{by (29)}). \end{aligned}$$

Therefore, by canceling $\|A\bar{x} - b\|_2^2$ from both sides, we know that, for all $x \in \mathbb{R}^n$,

$$0 \geq \|A(\bar{x} - x)\|_2^2 + 2(A(\bar{x} - x))^\top (Ax - b).$$

By setting $x = \bar{x} + tw$ for $t \in \mathbb{R}$ and $w \in \mathbb{R}^n$, we find that

$$0 \geq t^2 \|Aw\|_2^2 - 2t(Aw)^\top (A\bar{x} - b + tAw) \quad \forall t \in \mathbb{R} \quad \text{and} \quad w \in \mathbb{R}^n.$$

Dividing by $t \neq 0$ and letting t tend to zero, we find that

$$0 \geq -2w^\top A^\top (A\bar{x} - b) \quad \forall w \in \mathbb{R}^n,$$

This immediately implies $A^\top (A\bar{x} - b) = 0$ (why?), or equivalently,

$$(32) \quad A^\top A\bar{x} = A^\top b.$$

The system of equations (32) is called the *normal equations* associated with the linear least squares problem (28). This derivation leads to the following theorem.

THEOREM 2.1. [*Linear Least Squares and the Normal Equations*]
The vector \bar{x} solves the problem (28), i.e.,

$$\|A\bar{x} - b\|_2 \leq \|Ax - b\|_2 \quad \forall x \in \mathbb{R}^n,$$

if and only if $A^\top A\bar{x} = A^\top b$.

PROOF. We have just shown that if \bar{x} is a solution to (28), then the normal equations are satisfied, so we need only establish the reverse implication. Assume that (32) is satisfied. Then, for all $x \in \mathbb{R}^n$,

$$\begin{aligned} \|Ax - b\|_2^2 &= \|(Ax - A\bar{x}) + (A\bar{x} - b)\|_2^2 \\ &= \|A(x - \bar{x})\|_2^2 + 2(A(x - \bar{x}))^\top (A\bar{x} - b) + \|A\bar{x} - b\|_2^2 \quad (\text{by (30)}) \\ &\geq 2(x - \bar{x})^\top A^\top (A\bar{x} - b) + \|A\bar{x} - b\|_2^2 \quad (\text{since } \|A(x - \bar{x})\|_2^2 \geq 0) \\ &= \|A\bar{x} - b\|_2^2 \quad (\text{since } A^\top (A\bar{x} - b) = 0), \end{aligned}$$

or equivalently, \bar{x} solves (28). □

This theorem provides a nice characterization of solutions to (28), but it does not tell us if a solution exists. For this we use the following elementary result from linear algebra.

LEMMA 2.1. *For every matrix $A \in \mathbb{R}^{m \times n}$ we have*

$$\text{Null}(A^T A) = \text{Null}(A) \quad \text{and} \quad \text{Ran}(A^T A) = \text{Ran}(A^T) .$$

PROOF. Note that if $x \in \text{Null}(A)$, then $Ax = 0$ and so $A^T Ax = 0$, that is, $x \in \text{Null}(A^T A)$. Therefore, $\text{Null}(A) \subset \text{Null}(A^T A)$. Conversely, if $x \in \text{Null}(A^T A)$, then

$$A^T Ax = 0 \implies x^T A^T Ax = 0 \implies (Ax)^T (Ax) = 0 \implies \|Ax\|_2^2 = 0 \implies Ax = 0,$$

or equivalently, $x \in \text{Null}(A)$. Therefore, $\text{Null}(A^T A) \subset \text{Null}(A)$, and so $\text{Null}(A^T A) = \text{Null}(A)$.

Since $\text{Null}(A^T A) = \text{Null}(A)$, the Fundamental Theorem of the Alternative tells us that

$$\text{Ran}(A^T A) = \text{Ran}((A^T A)^T) = \text{Null}(A^T A)^\perp = \text{Null}(A)^\perp = \text{Ran}(A^T),$$

which proves the lemma. □

This lemma immediately gives us the following existence result.

THEOREM 2.2. *[Existence and Uniqueness for the Linear Least Squares Problem]*

Consider the linear least squares problem (28).

- (1) *A solution to the normal equations (32) always exists.*
- (2) *A solution to the linear least squares problem (28) always exists.*
- (3) *The linear least squares problem (28) has a unique solution if and only if $\text{Null}(A) = \{0\}$ in which case $(A^T A)^{-1}$ exists and the unique solution is given by $\bar{x} = (A^T A)^{-1} A^T b$.*
- (4) *If $\text{Ran}(A) = \mathbb{R}^m$, then $(AA^T)^{-1}$ exists and $\bar{x} = A^T (AA^T)^{-1} b$ solves (28), indeed, $A\bar{x} = b$.*

PROOF. (1) Lemma 2.1 tells us that $\text{Ran}(A^T A) = \text{Ran}(A^T)$; hence, a solution to $A^T Ax = A^T b$ must exist.

(2) This follows from Part (1) and Theorem 2.1.

(3) By Theorem 2.1, \bar{x} solves the linear least squares problem if and only if \bar{x} solves the normal equations. Hence, the linear least squares problem has a unique solution if and only if the normal equations have a unique solution. Since $A^T A \in \mathbb{R}^{n \times n}$ is a square matrix, this is equivalent to saying that $A^T A$ is invertible, or equivalently, $\text{Null}(A^T A) = \{0\}$. However, by Lemma 2.1, $\text{Null}(A) = \text{Null}(A^T A)$. Therefore, the linear least squares problem has a unique solution if and only if $\text{Null}(A) = \{0\}$ in which case $A^T A$ is invertible and the unique solution is given by $\bar{x} = (A^T A)^{-1} A^T b$.

(4) By the hypotheses, Lemma 2.1, and the Fundamental Theorem of the Alternative, $\{0\} = (\mathbb{R}^m)^\perp = (\text{Ran}(A))^\perp = \text{Null}(A^T) = \text{Null}(AA^T)$; hence, $AA^T \in \mathbb{R}^{m \times m}$ is invertible. Consequently, $\bar{x} = A^T (AA^T)^{-1} b$ is well-defined and satisfies $A\bar{x} = b$ □

The results given above establish the existence and uniqueness of solutions, provide necessary and sufficient conditions for optimality, and, in some cases, give a formula for the solution to the linear least squares problem. However, these results do not indicate how a solution can be computed. Here the dimension of the problem, or the problem size, plays a key role. In addition, the level of accuracy in the solution as well as the greatest accuracy possible are also issues of concern. Linear least squares problems range in size from just a few variables and equations to millions. Some are so large that all of the computing resources at our disposal today are insufficient to solve them, and in many cases the matrix A is not even available to us although, with effort, we can obtain Ax for a given vector x . Therefore, great care and inventiveness is required in the numerical solution of these problems. Research into how to solve this class of problems is still a very hot research topic today.

In our study of numerical solution techniques we present two classical methods. But before doing so, we study other aspects of the problem in order to gain further insight into its geometric structure.

3. Orthogonal Projection onto a Subspace

In this section we view the linear least squares problem from the perspective of a least distance problem to a subspace, or equivalently, as a projection problem for a subspace. Suppose $S \subset \mathbb{R}^m$ is a given subspace and $b \notin S$. The least distance problem for S and b is to find that element of S that is as close to b as possible. That is we wish to solve the problem

$$(33) \quad \min_{z \in S} \frac{1}{2} \|z - b\|_2^2 ,$$

or equivalently, we wish to find the point $\bar{z} \in S$ such that

$$\|\bar{z} - b\|_2 \leq \|z - b\|_2 \quad \forall z \in S.$$

If we now take the subspace to be the range of A , $S = \text{Ran}(A)$, then the problem (33) is closely related to the problem (28) since

$$(34) \quad \bar{z} \in \mathbb{R}^m \text{ solves (33) if and only if there is an } \bar{x} \in \mathbb{R}^n \text{ with } \bar{z} = A\bar{x} \text{ such that } \bar{x} \text{ solves (28).} \quad (\text{why?})$$

Below we discuss this connection and its relationship to the notion of an *orthogonal projection* onto a subspace.

A matrix $P \in \mathbb{R}^{m \times m}$ is said to be a *projection* if and only if $P^2 = P$. In this case we say that P is a projection onto the subspace $S = \text{Ran}(P)$, the range of P . Note that if $x \in \text{Ran}(P)$, then there is a $w \in \mathbb{R}^m$ such that $x = Pw$, therefore, $Px = P(Pw) = P^2w = Pw = x$. That is, P leaves all elements of $\text{Ran}(P)$ fixed. Also, note that, if P is a projection, then

$$(I - P)^2 = I - P - P + P^2 = I - P,$$

and so $(I - P)$ is also a projection. Since for all $w \in \mathbb{R}^m$,

$$w = Pw + (I - P)w,$$

we have

$$\mathbb{R}^m = \text{Ran}(P) + \text{Ran}(I - P).$$

In this case we say that the subspaces $\text{Ran}(P)$ and $\text{Ran}(I - P)$ are *complementary subspaces* since their sum is the whole space and their intersection is the origin, i.e., $\text{Ran}(P) \cap \text{Ran}(I - P) = \{0\}$ (why?).

Conversely, given any two subspaces S_1 and S_2 that are complementary, that is, $S_1 \cap S_2 = \{0\}$ and $S_1 + S_2 = \mathbb{R}^m$, there is a projection P such that $S_1 = \text{Ran}(P)$ and $S_2 = \text{Ran}(I - P)$. We do not show how to construct these projections here, but simply note that they can be constructed with the aid of bases for S_1 and S_2 .

The relationship between projections and complementary subspaces allows us to define a notion of *orthogonal projection*. Recall that for every subspace $S \subset \mathbb{R}^m$, we have defined

$$S^\perp := \{x \mid x^T y = 0 \quad \forall y \in S\}$$

as the subspace orthogonal to S . Clearly, S and S^\perp are complementary:

$$S \cap S^\perp = \{0\} \quad \text{and} \quad S + S^\perp = \mathbb{R}^m. \quad (\text{why?})$$

Therefore, there is a projection P such that $\text{Ran}(P) = S$ and $\text{Ran}(I - P) = S^\perp$, or equivalently,

$$(35) \quad ((I - P)y)^T (Pw) = 0 \quad \forall y, w \in \mathbb{R}^m.$$

The orthogonal projection plays a very special role among all possible projections onto a subspace. For this reason, we denote the orthogonal projection onto the subspace S by P_S .

We now use the condition (35) to derive a simple test of whether a linear transformation is an orthogonal projection. For brevity, we write $P := P_S$ and set $M = (I - P)^T P$. Then, by (35),

$$0 = e_i^T M e_j = M_{ij} \quad \forall i, j = 1, \dots, n,$$

i.e., M is the zero matrix. But then, since $0 = (I - P)^T P = P - P^T P$,

$$P = P^T P = (P^T P)^T = P^T.$$

Conversely, if $P = P^T$ and $P^2 = P$, then $(I - P)^T P = 0$. Therefore, a matrix P is an orthogonal projection if and only if $P^2 = P$ and $P = P^T$.

An orthogonal projection for a given subspace S can be constructed from any orthonormal basis for that subspace. Indeed, if the columns of the matrix Q form an orthonormal basis for S , then the matrix $P = QQ^T$ satisfies

$$P^2 = QQ^T QQ^T \stackrel{\text{why?}}{=} Q I_k Q^T = QQ^T = P \quad \text{and} \quad P^T = (QQ^T)^T = QQ^T = P,$$

where $k = \dim(S)$, and so P is the orthogonal projection onto S since, by construction, $\text{Ran}(QQ^T) = \text{Ran}(Q) = S$. We catalogue these observations in the following lemma.

LEMMA 3.1. [*Orthogonal Projections*]

- (1) The projection $P \in \mathbb{R}^{n \times n}$ is orthogonal if and only if $P = P^T$.
- (2) If the columns of the matrix $Q \in \mathbb{R}^{n \times k}$ form an orthonormal basis for the subspace $S \subset \mathbb{R}^n$, then $P := QQ^T$ is the orthogonal projection onto S .

Let us now apply these projection ideas to the problem (33). Let $P := P_S$ be the orthogonal projection onto the subspace S , and let $\bar{z} = Pb$. Then, for every $z \in S$,

$$\begin{aligned} \|z - b\|_2^2 &= \|Pz - Pb - (I - P)b\|_2^2 && \text{(since } z \in S\text{)} \\ &= \|P(z - b) + (I - P)b\|_2^2 \\ &= \|P(z - b)\|_2^2 + 2(z - b)^T P^T (I - P)b + \|(I - P)b\|_2^2 \\ &= \|P(z - b)\|_2^2 + \|(I - P)b\|_2^2 && \text{(since } P = P^T \text{ and } P = P^2\text{)} \\ &\geq \|(P - I)b\|_2^2 && \text{(since } \|P(z - b)\|_2^2 \geq 0\text{)} \\ &= \|\bar{z} - b\|_2^2. \end{aligned}$$

Consequently, $\|\bar{z} - b\|_2 \leq \|z - b\|_2$ for all $z \in S$, that is, $\bar{z} = Pb$ solves (33).

THEOREM 3.1. [*Subspace Projection Theorem*]

Let $S \subset \mathbb{R}^m$ be a subspace and let $b \in \mathbb{R}^m \setminus S$. Then the unique solution to the least distance problem

$$\underset{z \in S}{\text{minimize}} \|z - b\|_2$$

is $\bar{z} := P_S b$, where P_S is the orthogonal projector onto S .

PROOF. Everything but the uniqueness of the solution has been established in the discussion preceding the theorem. For this we make use of the identity

$$\|(1 - t)u + tv\|_2^2 = (1 - t)\|u\|_2^2 + t\|v\|_2^2 - t(1 - t)\|u - v\|_2^2 \quad \forall 0 \leq t \leq 1. \quad \text{(Verify!)}$$

Let $z^1, z^2 \in \mathbb{R}^m$ be two points that solve the minimum distance problem. Then, $\|z^1 - b\|_2 = \|z^2 - b\|_2 =: \eta > 0$, and so by the identity given above,

$$\begin{aligned} \|\tfrac{1}{2}(z^1 + z^2) - b\|_2^2 &= \|\tfrac{1}{2}(z^1 - b) + \tfrac{1}{2}(z^2 - b)\|_2^2 \\ &= \tfrac{1}{2}\|z^1 - b\|_2^2 + \tfrac{1}{2}\|z^2 - b\|_2^2 - \tfrac{1}{4}\|z^1 - z^2\|_2^2 \\ &= \eta^2 - \tfrac{1}{4}\|z^1 - z^2\|_2^2. \end{aligned}$$

Since $\eta = \inf \{\|z - b\|_2 \mid z \in S\}$, we must have $z^1 = z^2$. □

Let us now reconsider the linear least-squares problem (28) as it relates to our new found knowledge about orthogonal projections and their relationship to least distance problems for subspaces. Consider the case where $m \gg n$ and $\text{Null}(A) = \{0\}$. In this case, Theorem 2.2 tells us that $\bar{x} = (A^T A)^{-1} A^T b$ solves (28), and $\bar{z} = P_S b$ solves (35) where P_S is the orthogonal projector onto $S = \text{Ran}(A)$. Hence, by (34),

$$P_S b = \bar{z} = A\bar{x} = A(A^T A)^{-1} A^T b.$$

Since this is true for all possible choices of the vector b , we have

$$(36) \quad P_S = P_{\text{Ran}(A)} = A(A^T A)^{-1} A^T !$$

That is, the matrix $A(A^T A)^{-1} A^T$ is the orthogonal projector onto the range of A . One can also check this directly by showing that the matrix $M = A(A^T A)^{-1} A^T$ satisfies $M^2 = M$, $M^T = M$, and $\text{Ran}(M) = \text{Ran}(A)$.

PROPOSITION 3.1. Let $A \in \mathbb{R}^{m \times n}$ with $m \leq n$ and $\text{Null}(A) = \{0\}$. Then

$$P_{\text{Ran}(A)} = A(A^T A)^{-1} A^T.$$

4. Minimal Norm Solutions to $Ax = b$

Again let $A \in \mathbb{R}^{m \times n}$, but now we suppose that $m \ll n$. In this case A is short and fat so the matrix A most likely has rank m , or equivalently,

$$(37) \quad \text{Ran}(A) = \mathbb{R}^m.$$

But regardless of the range of A and the choice of the vector $b \in \mathbb{R}^m$, the set of solutions to $Ax = b$ will be infinite since the nullity of A is $n - m$. Indeed, if x^0 is any particular solution to $Ax = b$, then the set of solutions is given

by $x^0 + \text{Null}(A) := \{x^0 + z \mid z \in \text{Null}(A)\}$. In this setting, one might prefer the solution to the system having least norm. This solution is found by solving the problem

$$(38) \quad \min_{z \in \text{Null}(A)} \frac{1}{2} \|z + x^0\|_2^2 .$$

This problem is of the form (33). Consequently, the solution is given by $\bar{z} = -P_S x^0$ where P_S is now the orthogonal projection onto $S := \text{Null}(A)$.

In this context, note that $(I - P_{\text{Null}(A)})$ is the orthogonal projector onto $\text{Null}(A)^\perp = \text{Ran}(A^T)$. Recall that the formula (36) shows that if $M \in \mathbb{R}^{k \times s}$ is such that $\text{Null}(M) = \{0\}$, then the orthogonal projector onto $\text{Ran}(M)$ is given by

$$(39) \quad P_{\text{Ran}(M)} = M(M^T M)^{-1} M^T .$$

In our case, $M = A^T$ and $M^T M = AA^T$. Our working assumption (37) implies that

$$\text{Null}(M) = \text{Null}(A^T) = \text{Ran}(A)^\perp = (\mathbb{R}^m)^\perp = \{0\}$$

and consequently, by (39), the orthogonal projector onto $\text{Ran}(A^T)$ is given by

$$P_{\text{Ran}(A^T)} = A^T (AA^T)^{-1} A .$$

Therefore, the orthogonal projector onto $\text{Null}(A) = \text{Ran}(A^T)^\perp$ is

$$P_{\text{Null}(A)} = I - P_{\text{Ran}(A^T)} = I - P_{\text{Ran}(A^T)} = I - A^T (AA^T)^{-1} A .$$

Putting this all together, we find that the solution to (38) is

$$\bar{z} = P_{\text{Null}(A)}(-x^0) = (A^T (AA^T)^{-1} A - I)x^0 ,$$

and the solution to $Ax = b$ of least norm is

$$\bar{x} = x^0 + \bar{z} = A^T (AA^T)^{-1} Ax^0 ,$$

where x^0 is any particular solution to $Ax = b$, i.e., $Ax^0 = b$. Plugging \bar{x} into $Ax = b$ gives

$$A\bar{x} = AA^T (AA^T)^{-1} Ax^0 = Ax^0 = b .$$

THEOREM 4.1. [*Least Norm Solution to Linear Systems*] Let $A \in \mathbb{R}^{m \times n}$ be such that $m \leq n$ and $\text{Ran}(A) = \mathbb{R}^m$.

(1) The matrix AA^T is invertible.

(2) The orthogonal projection onto $\text{Null}(A)$ is given by

$$P_{\text{Null}(A)} = I - A^T (AA^T)^{-1} A .$$

(3) For every $b \in \mathbb{R}^m$, the system $Ax = b$ is consistent, and the least norm solution to this system is uniquely given by

$$\bar{x} = A^T (AA^T)^{-1} b .$$

5. Gram-Schmidt Orthogonalization, the QR Factorization, and Solving the Normal Equations

5.1. Gram-Schmidt Orthogonalization. We learned in the previous sections the important role orthogonal projections play in the linear least squares problem. In addition, we found that if the matrix U contains an orthonormal basis for the subspace S , then the matrix UU^T is the orthogonal projection onto S . Hence, one way to obtain these projections is to compute an orthogonal basis for a subspace. This is precisely what the *Gram-Schmidt orthogonalization* process does.

Let us recall the Gram-Schmidt orthogonalization process for a sequence of linearly independent vectors $a_1, \dots, a_k \in \mathbb{R}^m$ (note that this implies that $n \leq m$ (why?)). In this process we define vectors q_1, \dots, q_n inductively, as follows: set

$$p_1 = a_1, \quad q_1 = p_1 / \|p_1\| ,$$

$$p_j = a_j - \sum_{i=1}^{j-1} \langle a_j, q_i \rangle q_i \quad \text{and} \quad q_j = p_j / \|p_j\| \quad \text{for} \quad 2 \leq j \leq n .$$

For $1 \leq j \leq n$, $q_j \in \text{Span}\{a_1, \dots, a_j\}$, so $p_j \neq 0$ by the linear independence of a_1, \dots, a_j . An elementary induction argument shows that the q_j 's form an orthonormal basis for $\text{span}(a_1, \dots, a_n)$.

If we now define

$$r_{jj} = \|p_j\| \neq 0 \quad \text{and} \quad r_{ij} = \langle a_j, q_i \rangle \quad \text{for} \quad 1 \leq i < j \leq n,$$

then

$$\begin{aligned} a_1 &= r_{11} q_1, \\ a_2 &= r_{12} q_1 + r_{22} q_2, \\ a_3 &= r_{13} q_1 + r_{23} q_2 + r_{33} q_3, \\ &\vdots \\ a_n &= \sum_{i=1}^n r_{in} q_i. \end{aligned}$$

Set

$$A = [a_1 \ a_2 \ \dots \ a_n] \in \mathbb{R}^{m \times n}, \quad R = [r_{ij}] \in \mathbb{R}^{n \times n}, \quad \text{and} \quad Q = [q_1 \ q_2 \ \dots \ q_n] \in \mathbb{R}^{m \times n},$$

where $r_{ij} = 0$, $i > j$. Then

$$A = QR,$$

where Q is unitary and R is an upper triangular $n \times n$ matrix. In addition, R is invertible since the diagonal entries r_{jj} are non-zero. This is called the *QR factorization* of the matrix A .

REMARK 5.1. *If the a_j 's for $j = 1, \dots, n$ are linearly dependent, then, for some value(s) of j ,*

$$a_j \in \text{Span}\{a_1, \dots, a_{j-1}\}, \quad \text{and so} \quad p_j = 0.$$

The process can be modified by setting $r_{jj} = 0$, not defining a new q_j for this iteration, but continuing to define $r_{ij} = \langle a_j, q_i \rangle$ for $1 \leq i < j$, and proceeding. We still end up with orthonormal vectors $\{q_1, q_2, \dots, q_k\}$, but now $k < n$. In general, after n iterations, there will be $1 \leq k \leq n$ vectors $\{q_1, \dots, q_k\}$ that form an orthonormal basis for $\text{Span}\{a_1, \dots, a_n\}$, where $n - k$ is the number of diagonal entries r_{jj} that take the value zero. Again we obtain $A = QR$, but now Q may not be square and the matrix R may have zero diagonal entries in which case it is not invertible.

REMARK 5.2. *The classical Gram-Schmidt algorithm as described above does not behave well computationally. This is due to the accumulation of round-off error. The computed q_j 's are not orthogonal: $\langle q_j, q_k \rangle$ is small for $j \neq k$ with j near k , but not so small for $j \ll k$ or $j \gg k$.*

An alternate version, "Modified Gram-Schmidt," is equivalent in exact arithmetic, but behaves better numerically. In the following "pseudo-codes," p denotes a temporary storage vector used to accumulate the sums defining the p_j 's.

<u>Classic Gram-Schmidt</u>	<u>Modified Gram-Schmidt</u>
For $j = 1, \dots, n$ do	For $j = 1, \dots, n$ do
$p := a_j$	$p := a_j$
For $i = 1, \dots, j - 1$ do	For $i = 1, \dots, j - 1$ do
$r_{ij} = \langle a_j, q_i \rangle$	$r_{ij} = \langle p, q_i \rangle$
$p := p - r_{ij} q_i$	$p := p - r_{ij} q_i$
$r_{jj} := \ p\ $	$r_{jj} := \ p\ $
$q_j := p/r_{jj}$	$q_j := p/r_{jj}$

The only difference is in the computation of r_{ij} : in Modified Gram-Schmidt, we orthogonalize the accumulated partial sum for p_j against each q_i successively.

THEOREM 5.1. *[The Full QR Factorization] Suppose $A \in \mathbb{R}^{m \times n}$ with $m \geq n$. Then there exists a permutation matrix $P \in \mathbb{R}^{n \times n}$, a unitary matrix $Q \in \mathbb{R}^{m \times m}$, and an upper triangular matrix $R \in \mathbb{R}^{m \times n}$ such that $AP = QR$.*

Let $Q_1 \in \mathbb{R}^{m \times n}$ denote the first n columns of Q , Q_2 the remaining $(m - n)$ columns of Q , and $R_1 \in \mathbb{R}^{n \times n}$ the first n rows of R , then

$$(40) \quad AP = QR = [Q_1 \ Q_2] \begin{bmatrix} R_1 \\ 0 \end{bmatrix} = Q_1 R_1.$$

Moreover, we have the following:

- (a) We may choose R to have nonnegative diagonal entries.
- (b) If A is of full rank, then we can choose R with positive diagonal entries, in which case we obtain the condensed factorization $A = Q_1 R_1$, where $R_1 \in \mathbb{R}^{n \times n}$ invertible and the columns of Q_1 forming an orthonormal basis for the range of A .
- (c) If $\text{rank}(A) = k < n$, then

$$R_1 = \begin{bmatrix} R_{11} & R_{12} \\ 0 & 0 \end{bmatrix},$$

where R_{11} is a $k \times k$ invertible upper triangular matrix and $R_{12} \in \mathbb{R}^{k \times (n-k)}$. In particular, this implies that $AP = Q_{11}[R_{11} \ R_{12}]$, where Q_{11} are the first k columns of Q . In this case, the columns of Q_{11} form an orthonormal basis for the range of A .

REMARK 5.3. We call the factorization $AP = Q_{11}[R_{11} \ R_{12}]$ in Part (c) above the condensed QR Factorization. Note that if P is a permutation matrix, then so is P^T with $P^{-1} = P^T$ (i.e. permutation matrices are unitary). The role of the permutation matrix is to make the first $k = \text{rank}(A)$ columns of AP linearly independent.

To distinguish the condensed QR Factorization from the factorization in (40) with Q an $m \times m$ unitary matrix, we will refer the factorization where Q is unitary as the full QR factorization.

PROOF. If necessary, permute the columns of A so that the first $k = \text{rank}(A)$ columns of A are linearly independent and let P denote the permutation matrix that accomplishes this task so the the first k columns of AP are linearly independent. Apply the Gram-Schmidt orthogonalization process to obtain the matrix

$$Q_1 = [q_1, \dots, q_k] \in \mathbb{R}^{m \times k} \quad \text{and the upper triangular matrix} \quad \tilde{R}_{11} = [r_{ij}] \in \mathbb{R}^{k \times k}$$

so that $Q_1 R_1$ gives the first k columns of A . The write the remaining columns of A as linear combinations of the columns of Q_1 to obtain the coefficient matrix $R_{12} \in \mathbb{R}^{k \times (n-k)}$ yielding $AP = Q_1 [R_{11} \ R_{12}]$. Finally, extend $\{q_1, \dots, q_k\}$ to an orthonormal basis $\{q_1, \dots, q_m\}$ of \mathbb{R}^m , and set

$$Q = [q_1, \dots, q_m] \quad \text{and} \quad R = \begin{bmatrix} R_1 \\ 0 \end{bmatrix} \in \mathbb{R}^{m \times n}, \quad \text{so } AP = QR.$$

As $r_{jj} > 0$ in the Gram-Schmidt process, we have (b). □

REMARK 5.4. There are more efficient and better computationally behaved ways of calculating the Q and R factors. The idea is to create zeros below the diagonal (successively in columns 1, 2, ...) as in Gaussian Elimination, except instead of doing this by successive left multiplication by Gaussian elimination matrices, we left multiply by unitary matrices. Below, we show how this can be done with Householder transformations. But another popular approach is to use Givens rotations.

In practice, every $A \in \mathbb{R}^{m \times n}$ has a QR-factorization, even when $m < n$. This follows immediately from Part (c) Theorem 5.1.

COROLLARY 5.1.1. [The General Condensed QR Factorization] Let $A \in \mathbb{R}^{m \times n}$ have $\text{rank } k \leq \min\{m, n\}$. Then there exist

$$\begin{aligned} Q &\in \mathbb{R}^{m \times k} \quad \text{with orthonormal columns,} \\ R &\in \mathbb{R}^{k \times n} \quad \text{full rank upper triangular, and} \\ P &\in \mathbb{R}^{n \times n} \quad \text{a permutation matrix} \end{aligned}$$

such that

$$AP = QR.$$

In particular, the columns of the matrix Q form a basis for the range of A . Moreover, the matrix R can be written in the form

$$R = [R_1 \ R_2],$$

where $R_1 \in \mathbb{R}^{k \times k}$ is nonsingular.

REMARK 5.5. *The permutation P in the corollary above can be taken to be any permutation that re-orders the columns of A so that the first k columns of A are linearly independent, where k is the rank of A (similarly for \tilde{P} in permuting the columns of A^T).*

COROLLARY 5.1.2. *[Orthogonal Projections onto the Four Fundamental Subspaces] Let $A \in \mathbb{R}^{m \times n}$ have rank $k \leq \min\{m, n\}$. Let A and A^T have generalized QR factorizations*

$$AP = Q[R_1 \ R_2] \quad \text{and} \quad A^T \tilde{P} = \tilde{Q}[\tilde{R}_1 \ \tilde{R}_2].$$

Since row rank equals column rank, $P \in \mathbb{R}^{n \times n}$ is a permutation matrix, $\tilde{P} \in \mathbb{R}^{m \times m}$ is a permutation matrix, $Q \in \mathbb{R}^{m \times k}$ and $\tilde{Q} \in \mathbb{R}^{n \times k}$ have orthonormal columns, $R_1, \tilde{R}_1 \in \mathbb{R}^{k \times k}$ are both upper triangular nonsingular matrices, $R_2 \in \mathbb{R}^{k \times (n-k)}$, and $\tilde{R}_2 \in \mathbb{R}^{k \times (m-k)}$. Moreover,

$$\begin{aligned} QQ^T & \text{ is the orthogonal projection onto } \text{Ran}(A), \\ I - QQ^T & \text{ is the orthogonal projection onto } \text{Null}(A^T), \\ \tilde{Q}\tilde{Q}^T & \text{ is the orthogonal projection onto } \text{Ran}(A^T), \text{ and} \\ I - \tilde{Q}\tilde{Q}^T & \text{ is the orthogonal projection onto } \text{Null}(A)^\perp. \end{aligned}$$

PROOF. The result follows immediately from Corollary 5.1.1 and the Fundamental Theorem of the Alternative. \square

EXERCISE 5.1. *Verify the representations of the orthogonal projections onto $\text{Ran}(A)$ and $\text{Null}(A)$ given in Corollary 5.1.2 correspond to those given in Proposition 3.1 and Theorem 4.1.*

5.2. Solving the Normal Equations with the QR Factorization. Let's now reconsider the linear least squares problem (28) and how the QR factorization can be used in its solution. Specifically, we examine how it can be used to solve the normal equations $A^T A x = A^T b$. Let A and b be as in (28), and let

$$AP = Q[R_1 \ R_2]$$

be the general condensed QR factorization of A , where $P \in \mathbb{R}^{n \times n}$ is a permutation matrix, $Q \in \mathbb{R}^{m \times k}$ has orthonormal columns, $R_1 \in \mathbb{R}^{k \times k}$ is nonsingular and upper triangular, and $R_2 \in \mathbb{R}^{k \times (n-k)}$ with $k = \text{rank}(A) \leq \min\{n, m\}$. Replacing A by $A = Q[R_1 \ R_2]P^T$ in the normal equations gives the following equivalent system:

$$P^T \begin{bmatrix} R_1^T \\ R_2^T \end{bmatrix} Q^T Q [R_1 \ R_2] P x = P^T \begin{bmatrix} R_1^T \\ R_2^T \end{bmatrix} [R_1 \ R_2] P x = A^T b = P^T \begin{bmatrix} R_1^T \\ R_2^T \end{bmatrix} Q^T b,$$

since $Q^T Q = I_k$ the $k \times k$ identity matrix. By multiplying on the left by P , replacing b by $\hat{b} := Q^T b \in \mathbb{R}^k$ and x by

$$(41) \quad z := [R_1 \ R_2] P x,$$

we obtain

$$\begin{bmatrix} R_1^T \\ R_2^T \end{bmatrix} z = \begin{bmatrix} R_1^T \\ R_2^T \end{bmatrix} \hat{b}.$$

Let us see if we can reconstruct a solution to the normal equations by choosing the most obvious solution to the this system, namely, $\bar{z} := \hat{b}$. If this is to yield a solution to the normal equations, then, by (41), we need to solve the system

$$[R_1 \ R_2] P x = \hat{b}.$$

Set

$$\begin{bmatrix} w_1 \\ w_2 \end{bmatrix} := P x,$$

where $w_1 \in \mathbb{R}^k$ and $w_2 \in \mathbb{R}^{(n-k)}$, and consider the system

$$R_1 w_1 = \hat{b} \in \mathbb{R}^k.$$

Since $R_1 \in \mathbb{R}^{k \times k}$ is invertible, this system has a unique solution $\bar{w}_1 := R_1^{-1} \hat{b}$. Indeed, this system is very easy to solve using *back substitution* since R_1 is upper triangular. Next set $\bar{w}_2 = 0 \in \mathbb{R}^{(n-k)}$ and

$$\bar{x} := P^T \bar{w} = P^T \begin{bmatrix} R_1^{-1} \hat{b} \\ 0 \end{bmatrix}.$$

Then

$$\begin{aligned}
A^T A \bar{x} &= A^T A P^T \begin{bmatrix} R_1^{-1} \hat{b} \\ 0 \end{bmatrix} \\
&= A^T Q [R_1 \ R_2] P P^T \begin{bmatrix} R_1^{-1} \hat{b} \\ 0 \end{bmatrix} \\
&= A^T Q R_1 R_1^{-1} \hat{b} && \text{(since } P P^T = I) \\
&= A^T Q \hat{b} \\
&= A^T Q Q^T b \\
&= P^T \begin{bmatrix} R_1^T \\ R_2^T \end{bmatrix} Q^T Q Q^T b && \text{(since } A^T = P^T \begin{bmatrix} R_1^T \\ R_2^T \end{bmatrix} Q^T) \\
&= P^T \begin{bmatrix} R_1^T \\ R_2^T \end{bmatrix} Q^T b && \text{(since } Q^T Q = I) \\
&= A^T b,
\end{aligned}$$

that is, \bar{x} solves the normal equations!

Let us now consider the computational cost of obtaining the solution to the linear least squares problem in this way. The key steps in this computation are as follows:

$$\begin{array}{lll}
AP = Q[R_1 \ R_2] & \text{the general condensed QR factorization} & o(m^2 n) \\
\hat{b} = Q^T b & \text{a matrix-vector product} & o(km) \\
\bar{w}_1 = R_1^{-1} \hat{b} & \text{a back solve} & o(k^2) \\
\bar{x} = P^T \begin{bmatrix} R_1^{-1} \hat{b} \\ 0 \end{bmatrix} & \text{a matrix-vector product} & o(kn).
\end{array}$$

Therefore, the majority of the numerical effort is in the computation of the QR factorization.

5.3. Computing the Full QR Factorization using Householder Reflections. In subsection 5.1 we showed how to compute the QR factorization using the Gram-Schmidt orthogonalization procedure. We also indicated that due to numerical round-off error this procedure has difficulty in preserving the orthogonality of the columns of the matrix Q . To address this problem we presented the mathematically equivalent *modified* Gram-Schmidt process which has improved performance. We now present a very different method for obtaining the full QR factorization. The approach we describe is very much like Gauss-Jordan Elimination to obtain reduced echelon form. However, now we successively multiply A on the left by unitary matrices, rather than Gauss-Jordan elimination matrices, which eventually put A into upper triangular form. The matrices we multiply by are the *Householder reflection matrices*.

Given $w \in \mathbb{R}^n$ we can associate the matrix

$$U = I - 2 \frac{w w^T}{w^T w}$$

which reflects \mathbb{R}^n across the hyperplane $\text{Span}\{w\}^\perp$. The matrix U is called the Householder reflection across this hyperplane.

Given a pair of vectors x and y with

$$\|x\|_2 = \|y\|_2, \quad \text{and} \quad x \neq y,$$

the Householder reflection

$$U = I - 2 \frac{(x - y)(x - y)^T}{(x - y)^T (x - y)}$$

is such that $y = Ux$, since

$$\begin{aligned}
Ux &= x - 2(x - y) \frac{\|x\|^2 - y^T x}{\|x\|^2 - 2y^T x + \|y\|^2} \\
&= x - 2(x - y) \frac{\|x\|^2 - y^T x}{2(\|x\|^2 - y^T x)} \quad (\text{since } \|x\| = \|y\|) \\
&= y.
\end{aligned}$$

We now show how Householder reflections can be used to obtain the QR factorization. We begin by describing the basic *deflation* step in the QR-factorization of the matrix $A_0 \in \mathbb{R}^{m \times n}$ which we block decompose as

$$A_0 = \begin{bmatrix} \alpha_0 & a_0^T \\ b_0 & \tilde{A}_0 \end{bmatrix}, \text{ with } \tilde{A}_0 \in \mathbb{R}^{(m-1) \times (n-1)},$$

and set

$$\nu_0 = \left\| \begin{pmatrix} \alpha_0 \\ b_0 \end{pmatrix} \right\|_2.$$

If $\nu_0 = 0$, then multiply A_0 on the left by a permutation matrix P_0 to bring a non-zero (largest magnitude) column in A_0 into the first column and the zero column to the last column. Then block decompose $A_0 P_0$ as above with

$$A_0 P_0 = \begin{bmatrix} \alpha_0 & a_0^T \\ b_0 & \tilde{A}_0 \end{bmatrix}, \text{ with } \tilde{A}_0 \in \mathbb{R}^{(m-1) \times (n-1)},$$

and set

$$\nu_0 = \left\| \begin{pmatrix} \alpha_0 \\ b_0 \end{pmatrix} \right\|_2 \neq 0.$$

Let H_0 be the Householder transformation that maps

$$\begin{pmatrix} \alpha_0 \\ b_0 \end{pmatrix} \mapsto \nu_0 e_1 \quad :$$

$$H_0 = I - 2 \frac{ww^T}{w^T w} \quad \text{where} \quad w = \begin{pmatrix} \alpha_0 \\ b_0 \end{pmatrix} - \nu_0 e_1 = \begin{pmatrix} \alpha_0 - \nu_0 \\ b_0 \end{pmatrix}.$$

Then,

$$H_0 A = \begin{bmatrix} \nu_0 & a_1^T \\ 0 & A_1 \end{bmatrix}.$$

Now repeat with A_1 .

If the above method is implemented by always permuting the column of greatest magnitude into the current pivot column, then

$$AP = QR$$

gives a QR-factorization with the diagonal entries of R nonnegative and listed in the order of descending magnitude. Since Q is unitary, this is the full QR factorization in (40).

Optimization of Quadratic Functions

In this chapter we study the problem

$$(42) \quad \underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2}x^T Hx + g^T x + \beta,$$

where $H \in \mathbb{R}^{n \times n}$ is symmetric, $g \in \mathbb{R}^n$, and $\beta \in \mathbb{R}$. It has already been observed that we may as well assume that H is symmetric since

$$x^T Hx = \frac{1}{2}x^T Hx + \frac{1}{2}(x^T Hx)^T = x^T \left[\frac{1}{2}(H + H^T) \right] x,$$

where $\frac{1}{2}(H + H^T)$ is called the *symmetric part* of H . Therefore, in this chapter we assume that H is symmetric. In addition, we have also noted that an objective function can always be shifted by a constant value without changing the solution set to the optimization problem. Therefore, we assume that $\beta = 0$ for most of our discussion. However, just as in the case of integration theory where it is often helpful to choose a particular constant of integration, in many applications there is a “natural” choice for β that helps one interpret the problem as well as its solution.

The class of problems (42) is important for many reasons. Perhaps the most common instance of this problem is the linear least squares problem:

$$(43) \quad \underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2} \|Ax - b\|_2^2,$$

where $A \in \mathbb{R}^{m \times n}$, and $b \in \mathbb{R}^m$. By expanding the objective function in (43), we see that

$$(44) \quad \frac{1}{2} \|Ax - b\|_2^2 = \frac{1}{2}x^T(A^T A)x - (A^T b)^T x + \frac{1}{2} \|b\|_2^2 = \frac{1}{2}x^T Hx + g^T x + \beta,$$

where $H = A^T A$, $g = -A^T b$, and $\beta = \frac{1}{2} \|b\|_2^2$. This connection to the linear least squares problem will be explored in detail later in this chapter. For the moment, we continue to exam the general problem (42). As in the case of the linear least squares problem, we begin by discussing characterizations of the solutions as well as their existence and uniqueness. In this discussion we try to follow the approach taken for the the linear least squares problem. However, in the case of (43), the matrix $H := A^T A$ and the vector $g = -A^T b$ possess special features that allowed us to establish very strong results on optimality conditions as well as on the existence and uniqueness of solutions. In the case of a general symmetric matrix H and vector g it is possible to obtain similar results, but there are some twists. Symmetric matrices have many special properties that can be exploited to help us achieve our goal. Therefore, we begin by recalling a few of these properties, specifically those related to eigenvalue decomposition.

1. Eigenvalue Decomposition of Symmetric Matrices

Given a matrix $A \in \mathbb{R}^{n \times n}$, we say that the scalar λ is an eigenvalue of A if there is a non-zero vector x such that $Ax = \lambda x$, or equivalently, $\text{Null}(\lambda I - A) \neq \{0\}$. Observe that $\text{Null}(\lambda I - A) \neq \{0\}$ if and only if $(\lambda I - A)$ is singular, that is, $\det(\lambda I - A) = 0$. Consequently, λ is an eigenvalue of A if and only if $\det(\lambda I - A) = 0$. If we now think of λ as a variable, this says that we can find all eigenvalues of A by finding all roots of the equation $\det(\lambda I - A) = 0$. The function $p(\lambda) := \det(\lambda I - A)$ is easily seen to be a polynomial of degree n in λ which we call the *characteristic polynomial* of A . By the Fundamental Theorem of Algebra, we know that $p(\lambda)$ has n roots over the complex numbers if we count the multiplicities of these roots. Hence, when we discuss eigenvalues and eigenvectors we are forced in the setting of complex numbers. For this reason we may as well assume that $A \in \mathbb{C}^{n \times n}$.

Working on \mathbb{C}^n requires us to re-examine our notion of the Euclidean norm and its associated dot product. Recall that for a complex number $\zeta := x + iy$, with $x, y \in \mathbb{R}$ and $i := \sqrt{-1}$, the magnitude of ζ is given by $|\zeta| = \sqrt{\zeta \bar{\zeta}}$, where $\bar{\zeta} := x - iy$ is the complex conjugate of ζ . If we now define the Euclidean norm of a vector $z \in \mathbb{C}^n$ to be the square root of the sum of the squares of magnitude of its components, then

$$\|z\|_2 = \sqrt{\sum_{k=1}^n |z_k|^2} = \sqrt{\sum_{k=1}^n \bar{z}_k z_k} = \sqrt{\bar{z}^T z} = \sqrt{z^* z},$$

where we define

$$z^* z = (\bar{z})^T z,$$

that is, z^* takes z to its *conjugate transpose*. When $z \in \mathbb{R}^n$, we have $z^* = z^T$, and we recover the usual formulas. With the $*$ operation, we can extend our notion of dot product (or, inner product) by writing

$$\langle z, y \rangle := z^* y \in \mathbb{C}.$$

When z and y are real vectors we recover usual notion of dot product for such vectors. Finally, for matrices $A \in \mathbb{C}^{n \times n}$, we define

$$A^* := \bar{A}^T,$$

that is, we conjugate every element of A and then take the transpose. This notation is very helpful in a number of ways. For example, we have

$$\langle Ay, x \rangle = (Ay)^* x = y^* A^* x \quad \text{and} \quad \|Ax\|_2^2 = x^* A^* Ax.$$

We call A^* the *adjoint* of A .

Recall that a matrix $H \in \mathbb{R}^{n \times n}$ is said to be symmetric if $H^T = H$. By extension, we say that an matrix $Q \in \mathbb{C}^{n \times n}$ is *self-adjoint* if $Q^* = Q$. Thus, in particular, every real symmetric matrix is self adjoint. We have the following remarkable fact about self-adjoint matrices.

LEMMA 1.1. *If $Q \in \mathbb{C}^{n \times n}$ is self-adjoint, then Q has only real eigenvalues. In particular, if H is a real symmetric matrix, then H has only real eigenvalues and for each such eigenvalue there is a real eigenvector. Moreover, if (λ_1, v^1) and (λ_2, v^2) are two eigenvalue-eigenvectors pairs for H with $\lambda_1 \neq \lambda_2$, then $(v^1)^T v^2 = 0$.*

PROOF. Let $\lambda \in \mathbb{C}$ be an eigenvalue of Q . Then there is a non-zero eigenvector $x \in \mathbb{C}^n$ such that $Qx = \lambda x$. Therefore,

$$\lambda \|x\|_2^2 = \lambda x^* x = x^* Qx = x^* Q^* x = (x^* Qx)^* = (\lambda \|x\|_2^2)^* = \bar{\lambda} \|x\|_2^2,$$

so that $\lambda = \bar{\lambda}$ which can only occur if λ is a real number.

If H is real symmetric matrix, then it is self adjoint so all of its eigenvalues are real. If λ is one such eigenvalue with associated eigenvector $z = x + iy$ with $x, y \in \mathbb{R}^n$, then

$$Hx + iHy = Hz = \lambda z = \lambda x + i\lambda y.$$

Consequently, $Hx = \lambda x$ and $Hy = \lambda y$ since both Hx and Hy are real vectors. Since $z \neq 0$, either x or y or both are non-zero, in any case we have a real eigenvector for H corresponding to λ .

Next let (λ_1, v^1) and (λ_2, v^2) be eigenvalue-eigenvectors pairs for H with $\lambda_1 \neq \lambda_2$. Then

$$\lambda_1 (v^1)^T v^2 = (Hv^1)^T v^2 = (v^1)^T H v^2 = \lambda_2 (v^1)^T v^2,$$

since $\lambda_1 \neq \lambda_2$, we must have $(v^1)^T v^2 = 0$. □

Next, suppose λ_1 is an eigenvalue for the real symmetric matrix $H \in \mathbb{R}^{n \times n}$ and let the columns of the matrix $U_1 \in \mathbb{R}^{n \times k}$ form an orthonormal basis for the subspace $\text{Null}(\lambda_1 I - H)$, where $k = \dim(\text{Null}(\lambda_1 I - H)) \geq 1$. Let the columns of $U_2 \in \mathbb{R}^{n \times (n-k)}$ form an orthonormal basis for the subspace $\text{Null}(\lambda_1 I - H)^\perp$ and set $\tilde{U} = [U_1 \ U_2] \in \mathbb{R}^{n \times n}$. Then $\tilde{U}^T \tilde{U} = I$, that is, \tilde{U} is a *unitary matrix*. In particular, $\tilde{U}^{-1} = \tilde{U}^T$ and so $\tilde{U} \tilde{U}^T = I$ as well. We have the following relationships between U_1 , U_2 , and H :

$$HU_1 = \lambda_1 U_1, \quad U_1^T H U_1 = \lambda_1 U_1^T U_1 = \lambda_1 I_k \quad \text{and} \quad (U_1^T H U_2)^T = U_2^T H U_1 = \lambda_1 U_2^T U_1 = 0_{(n-k) \times k}.$$

Consequently,

$$(45) \quad \tilde{U}^T H \tilde{U} = \begin{bmatrix} U_1^T \\ U_2^T \end{bmatrix} H [U_1 \ U_2] = \begin{bmatrix} U_1^T \\ U_2^T \end{bmatrix} [H U_1 \ H U_2] = \begin{bmatrix} U_1^T H U_1 & U_1^T H U_2 \\ U_2^T H U_1 & U_2^T H U_2 \end{bmatrix} = \begin{bmatrix} \lambda_1 I_k & 0 \\ 0 & U_2^T H U_2 \end{bmatrix},$$

and so

$$(46) \quad H = \tilde{U} \tilde{U}^T H \tilde{U} \tilde{U}^T = [U_1 \ U_2] \begin{bmatrix} \lambda_1 I_k & 0 \\ 0 & U_2^T H U_2 \end{bmatrix} \begin{bmatrix} U_1^T \\ U_2^T \end{bmatrix}.$$

These observations provide the foundation for the following eigenvalue theorem for real symmetric matrices.

THEOREM 1.1. [*Eigenvalue Decomposition for Symmetric Matrices*] *Let $H \in \mathbb{R}^{n \times n}$ be a real symmetric matrix. Then there is a unitary matrix U such that*

$$H = U \Lambda U^T,$$

where $\Lambda := \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$ with $\lambda_1, \lambda_2, \dots, \lambda_n$ being the eigenvalues of H repeated according to multiplicity.

PROOF. We proceed by induction on the dimension. The result is trivially true for $n = 1$. Assume that the result is true for all dimensions $k < n$ with $n > 1$ and show it is true for all $n \times n$ symmetric matrices. Let $H \in \mathbb{R}^{n \times n}$ be symmetric and let λ_1 be any eigenvalue of H with $k = \dim(\text{Null}(\lambda_1 I - H)) \geq 1$. Let $U_1 \in \mathbb{R}^{n \times k}$ and $U_2 \in \mathbb{R}^{n \times (n-k)}$ be as in (45) and (46) above. If $k = n$, the result follows from (45) so we can assume that $k < n$.

Since (45) is a similarity transformation of H , $\tilde{U}^T H \tilde{U}$ has the same characteristic polynomial as H :

$$\det(\lambda I_n - H) = (\lambda - \lambda_1)^k q(\lambda), \quad \text{where} \quad q(\lambda) = \det(\lambda I_{n-k} - U_2^T H U_2).$$

Therefore, the eigenvalues of $U_2^T H U_2$ are necessarily those of H that are not equal to λ_1 and each has the same multiplicity as they have for H .

Apply the induction hypothesis to the $(n-k) \times (n-k)$ matrix $U_2^T H U_2$ to obtain a real unitary matrix $V \in \mathbb{R}^{(n-k) \times (n-k)}$ such that

$$U_2^T H U_2 = V \tilde{\Lambda} V^T,$$

where $\tilde{\Lambda} = \text{diag}(\mu_1, \mu_2, \dots, \mu_{(n-k)})$ with $\mu_1, \mu_2, \dots, \mu_{(n-k)}$ being the eigenvalues of H that are not equal to λ_1 with each having the same multiplicity as they have for H . Then, by (46)

$$H = [U_1 \quad U_2] \begin{bmatrix} \lambda_1 I_k & 0 \\ 0 & U_2^T H U_2 \end{bmatrix} \begin{bmatrix} U_1^T \\ U_2^T \end{bmatrix} = [U_1 \quad U_2] \begin{bmatrix} I_k & 0 \\ 0 & V \end{bmatrix} \begin{bmatrix} \lambda_1 I_k & 0 \\ 0 & \tilde{\Lambda} \end{bmatrix} \begin{bmatrix} I_k & 0 \\ 0 & V^T \end{bmatrix} \begin{bmatrix} U_1^T \\ U_2^T \end{bmatrix}.$$

The result is obtained by setting

$$U = [U_1 \quad U_2] \begin{bmatrix} I_k & 0 \\ 0 & V \end{bmatrix} = [U_1 \quad U_2 V]$$

and observing that $U^T U = I$. □

One important consequence of this result is the following theorem

THEOREM 1.2. *[The Rayleigh-Ritz Theorem] Let the symmetric matrix $H \in \mathbb{R}^{n \times n}$ have smallest eigenvalue $\lambda_{\min}(H)$ and largest eigenvalue $\lambda_{\max}(H)$. Then, for all $u \in \mathbb{R}^n$,*

$$\lambda_{\min}(H) \|u\|_2^2 \leq u^T H u \leq \lambda_{\max}(H) \|u\|_2^2,$$

with equality holding on the left for every eigenvector u for $\lambda_{\min}(H)$ and equality holding on the right for every eigenvector u for $\lambda_{\max}(H)$.

PROOF. Let $H = U \Lambda U^T$ be the eigenvalue decomposition of H in Theorem 1.1 with $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$. Then the columns of U form an orthonormal basis for \mathbb{R}^n . Therefore, given any $u \in \mathbb{R}^n \setminus \{0\}$, there is a $z \in \mathbb{R}^n \setminus \{0\}$ such that $u = Uz$. Hence

$$u^T H u = (Uz)^T U \Lambda U^T (Uz) = z^T \Lambda z = \sum_{j=1}^n \lambda_j z_j^2.$$

Clearly,

$$\lambda_{\min}(H) \|z\|_2^2 = \sum_{j=1}^n \lambda_{\min}(H) z_j^2 \leq \sum_{j=1}^n \lambda_j z_j^2 \leq \sum_{j=1}^n \lambda_{\max}(H) z_j^2 = \lambda_{\max}(H) \|z\|_2^2.$$

The result now follows since $\|z\|_2^2 = z^T z = z^T U^T U z = u^T u = \|u\|_2^2$. □

The following definition describes some important concepts associated with symmetric matrices that are important for optimization.

DEFINITION 1.1. *Let $H \in \mathbb{R}^{n \times n}$.*

- (1) H is said to be positive definite if $x^T H x > 0$ for all $x \in \mathbb{R}^n \setminus \{0\}$.
- (2) H is said to be positive semi-definite if $x^T H x \geq 0$ for all $x \in \mathbb{R}^n$.
- (3) H is said to be negative definite if $x^T H x < 0$ for all $x \in \mathbb{R}^n \setminus \{0\}$.
- (4) H is said to be negative semi-definite if $x^T H x \leq 0$ for all $x \in \mathbb{R}^n$.
- (5) H is said to be indefinite if H is none of the above.

We denote the set of real $n \times n$ symmetric matrices by \mathcal{S}^n , the set of positive semi-definite real $n \times n$ symmetric matrices by \mathcal{S}_+^n , and the set of positive definite real $n \times n$ symmetric matrices by \mathcal{S}_{++}^n . It is easily seen that \mathcal{S}^n is a vector space.

Theorem 1.2 provides necessary and sufficient conditions under which a symmetric matrix H is positive/negative definite/semi-definite. For example, since $\lambda_{\min}(H) \|u\|_2^2 \leq u^T H u$ with equality when u is an eigenvector associated

with $\lambda_{\min}(H)$, we have that H is positive definite if and only if $\lambda_{\min}(H) > 0$. Similar results can be obtained for the other cases.

An additional property of positive semi-definite matrices is that they possess *square roots*. If $H \in \mathbb{R}^{n \times n}$ is symmetric and positive semi-definite, then Theorem 1.1 tells us that $H = U\Lambda U^T$, where U is unitary and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ with $\lambda_i \geq 0$, $i = 1, \dots, n$. If we define $\Lambda^{1/2} := \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_n})$ and $H^{1/2} = U\Lambda^{1/2}U^T$, then $H = U\Lambda U^T = U\Lambda^{1/2}U^T U\Lambda^{1/2}U^T = H^{1/2}H^{1/2}$, so $H^{1/2}$ provides a natural notion of the square root of a matrix. However, $H^{1/2}$ is not uniquely defined since we can always re-order the diagonal elements and their corresponding columns to produce the same effect. In addition, $H^{1/2}$ is always symmetric while in some instances choosing a non-symmetric square root may be beneficial. For example, if we consider the linear least squares problem (43), then $H = A^T A$. Should A be considered a square root of H ? In order to cover the full range of possible considerations, we make the following definition for the square root of a symmetric matrix.

DEFINITION 1.2. *[Square Roots of Positive Semi-Definite Matrices] Let $H \in \mathbb{R}^{n \times n}$ be a symmetric positive semi-definite matrix. We say that the matrix $L \in \mathbb{R}^{n \times n}$ is a square root of H if $H = LL^T$.*

2. Optimality Properties of Quadratic Functions

Recall that for the linear least squares problem, we were able to establish a necessary and sufficient condition for optimality, namely the normal equations, by working backward from a known solution. We now try to apply this same approach to quadratic functions, in particular, we try to extend the derivation in (31) to the objective function in (47). Suppose \bar{x} is a local solution to the quadratic optimization problem

$$(47) \quad \underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2}x^T Hx + g^T x,$$

where $H \in \mathbb{R}^{n \times n}$ is symmetric and $g \in \mathbb{R}^n$, i.e., there is an $\epsilon > 0$ such that

$$(48) \quad \frac{1}{2}\bar{x}^T H\bar{x} + g^T \bar{x} \leq \frac{1}{2}x^T Hx + g^T x \quad \forall x \in \bar{x} + \epsilon\mathbb{B}_2,$$

where $\bar{x} + \epsilon\mathbb{B}_2 := \{\bar{x} + \epsilon u \mid u \in \mathbb{B}_2\}$ and $\mathbb{B}_2 := \{u \mid \|u\|_2 \leq 1\}$ (hence, $\bar{x} + \epsilon\mathbb{B}_2 = \{x \mid \|\bar{x} - x\|_2 \leq \epsilon\}$). Note that, for all $x \in \mathbb{R}^n$,

$$(49) \quad \begin{aligned} \bar{x}^T H\bar{x} &= (x + (\bar{x} - x))^T H(x + (\bar{x} - x)) \\ &= x^T Hx + 2x^T H(\bar{x} - x) + (\bar{x} - x)^T H(\bar{x} - x) \\ &= x^T Hx + 2(\bar{x} + (x - \bar{x}))^T H(\bar{x} - x) + (\bar{x} - x)^T H(\bar{x} - x) \\ &= x^T Hx + 2\bar{x}^T H(\bar{x} - x) + 2(x - \bar{x})^T H(\bar{x} - x) + (\bar{x} - x)^T H(\bar{x} - x) \\ &= x^T Hx + 2\bar{x}^T H(\bar{x} - x) - (\bar{x} - x)^T H(\bar{x} - x). \end{aligned}$$

Therefore, for all $x \in \bar{x} + \epsilon\mathbb{B}_2$,

$$\begin{aligned} \frac{1}{2}\bar{x}^T H\bar{x} + g^T \bar{x} &= \left(\frac{1}{2}x^T Hx + g^T x\right) + (H\bar{x} + g)^T(\bar{x} - x) - \frac{1}{2}(\bar{x} - x)^T H(\bar{x} - x) \\ &\geq \left(\frac{1}{2}\bar{x}^T H\bar{x} + g^T \bar{x}\right) + (H\bar{x} + g)^T(\bar{x} - x) - \frac{1}{2}(\bar{x} - x)^T H(\bar{x} - x), \end{aligned} \quad (\text{since } \bar{x} \text{ is a local solution})$$

and so

$$(50) \quad \frac{1}{2}(\bar{x} - x)^T H(\bar{x} - x) \geq (H\bar{x} + g)^T(\bar{x} - x) \quad \forall x \in \bar{x} + \epsilon\mathbb{B}_2.$$

Let $0 \leq t \leq \epsilon$ and $v \in \mathbb{B}_2$ and define $x = \bar{x} + tv \in \bar{x} + \epsilon\mathbb{B}_2$. If we plug $x = \bar{x} + tv$ into (50), then

$$(51) \quad \frac{t^2}{2}v^T Hv \geq -t(H\bar{x} + g)^T v.$$

Dividing this expression by $t > 0$ and taking the limit as $t \downarrow 0$ tells us that

$$0 \leq (H\bar{x} + g)^T v \quad \forall v \in \mathbb{B}_2,$$

which implies that $H\bar{x} + g = 0$. Plugging this information back into (51) gives

$$\frac{t^2}{2}v^T Hv \geq 0 \quad \forall v \in \mathbb{B}_2.$$

Dividing by $t^2/2$ for $t \neq 0$ tells us that

$$v^T Hv \geq 0 \quad \forall v \in \mathbb{B}_2$$

or equivalently, that H is positive semi-definite. These observations motivate the following theorem.

THEOREM 2.1. *[Existence and Uniqueness in Quadratic Optimization] Let $H \in \mathbb{R}^{n \times n}$ and $g \in \mathbb{R}^n$ be as in (47).*

- (1) *A local solution to the problem (47) exists if and only if H is positive semi-definite and there exists a solution \bar{x} to the equation $Hx + g = 0$ in which case \bar{x} is a local solution to (47).*
- (2) *If \bar{x} is a local solution to (47), then it is a global solution to (47).*
- (3) *The problem (47) has a unique global solution if and only if H is positive definite in which case this solution is given by $\bar{x} = -H^{-1}g$.*
- (4) *If either H is not positive semi-definite or there is no solution to the equation $Hx + g = 0$ (or both), then*

$$-\infty = \inf_{x \in \mathbb{R}^n} \frac{1}{2}x^T Hx + g^T x .$$

PROOF. (1) We have already shown that if a local solution \bar{x} to (47) exists, then $H\bar{x} + g = 0$ and H is positive semi-definite. On the other hand, suppose that H is positive semi-definite and \bar{x} is a solution to $Hx + g = 0$. Then, for all $x \in \mathbb{R}^n$, we can interchange the roles of x and \bar{x} in the second line of (49) to obtain

$$x^T Hx = \bar{x}^T H\bar{x} + 2\bar{x}^T H(x - \bar{x}) + (x - \bar{x})^T H(x - \bar{x}).$$

Hence, for all $x \in \mathbb{R}^n$,

$$\frac{1}{2}x^T Hx + g^T x = \frac{1}{2}\bar{x}^T H\bar{x} + g^T \bar{x} + (H\bar{x} + g)^T(x - \bar{x}) + \frac{1}{2}(x - \bar{x})^T H(x - \bar{x}) \geq \frac{1}{2}\bar{x}^T H\bar{x} + g^T \bar{x} ,$$

since $H\bar{x} + g = 0$ and H is positive semi-definite. That is, \bar{x} is a global solution to (47) and hence a local solution.

(2) Suppose \bar{x} is a local solution to (47) so that, by Part (1), H is positive semi-definite and $H\bar{x} + g = 0$, and there is an $\epsilon > 0$ such that (48) holds. Next observe that, for all $x, y \in \mathbb{R}^n$ and $\lambda \in \mathbb{R}$, we have

$$\begin{aligned} ((1 - \lambda)x + \lambda y)^T H((1 - \lambda)x + \lambda y) - (1 - \lambda)x^T Hx - \lambda y^T Hy & \\ &= (1 - \lambda)^2 x^T Hx + 2\lambda(1 - \lambda)x^T Hy + \lambda^2 y^T Hy - (1 - \lambda)x^T Hx - \lambda y^T Hy \\ &= -\lambda(1 - \lambda)x^T Hx + 2\lambda(1 - \lambda)x^T Hy - \lambda(1 - \lambda)y^T Hy \\ &= -\lambda(1 - \lambda)(x - y)^T H(x - y), \end{aligned}$$

or equivalently,

$$(52) \quad ((1 - \lambda)x + \lambda y)^T H((1 - \lambda)x + \lambda y) \leq (1 - \lambda)x^T Hx + \lambda y^T Hy - \lambda(1 - \lambda)(x - y)^T H(x - y).$$

Since H is positive semi-definite, this implies that

$$(53) \quad ((1 - \lambda)x + \lambda y)^T H((1 - \lambda)x + \lambda y) \leq (1 - \lambda)x^T Hx + \lambda y^T Hy \quad \forall \lambda \in [0, 1].$$

If \bar{x} is not a global solution, then there is an \hat{x} such that $f(\hat{x}) < f(\bar{x})$, where $f(x) := \frac{1}{2}x^T Hx + g^T x$. By (48), we must have $\|\bar{x} - \hat{x}\|_2 > \epsilon$. Set $\lambda := \frac{\epsilon}{2\|\bar{x} - \hat{x}\|_2}$ so that $0 < \lambda < 1$, and define $x_\lambda := (1 - \lambda)\bar{x} + \lambda\hat{x} = \bar{x} + \lambda(\hat{x} - \bar{x})$ so that $x_\lambda \in \bar{x} + \epsilon\mathbb{B}_2$. But then, by (53),

$$f(x_\lambda) \leq (1 - \lambda)f(\bar{x}) + \lambda f(\hat{x}) < (1 - \lambda)f(\bar{x}) + \lambda f(\bar{x}) = f(\bar{x}),$$

which contradicts (48). Hence, no such \hat{x} can exist so that \bar{x} is a global solution to (47).

(3) If (47) has a unique global solution \bar{x} , then \bar{x} must be the unique solution to the equation $Hx + g = 0$. This can only happen if H is invertible. Hence, H is invertible and positive semi-definite which implies that H is positive definite. On the other hand, if H is positive definite, then it is positive semi-definite and there is a unique solution to the equation $Hx + g = 0$, i.e., (49) has a unique global solution.

(4) The result follows if we can show that $f(x) := \frac{1}{2}x^T Hx + g^T x$ is unbounded below when either H is not positive semi-definite or there is no solution to the equation $Hx + g = 0$ (or both). Let us first suppose that H is not positive semi-definite, or equivalently, $\lambda_{\min}(H) < 0$. Let $u \in \mathbb{R}^n$ be an eigenvector associated with the eigenvalue $\lambda_{\min}(H)$ with $\|u\|_2 = 1$. Then, for $x := tu$ with $t > 0$, we have $f(tu) = \lambda_{\min}(H)\frac{t^2}{2} + tg^T u \xrightarrow{t \uparrow \infty} -\infty$ since $\lambda_{\min}(H) < 0$, so f is unbounded below.

Next suppose that there is no solution to the equation $Hx + g = 0$. In particular, $g \neq 0$ and $g \notin \text{Ran}(H) = \text{Null}(H)^\perp$. Then the orthogonal projection of g onto $\text{Null}(H)$ cannot be zero: $\hat{g} := P_{\text{Null}(H)}(g) \neq 0$. Hence, for $x := -t\hat{g}$ with $t > 0$, we have $f(-t\hat{g}) = -t\|\hat{g}\|_2^2 \xrightarrow{t \uparrow \infty} -\infty$, so again f is unbounded below. \square

The identity (52) is a very powerful tool in the analysis of quadratic functions. It was the key tool in showing that every local solution to (47) is necessarily a global solution. It is also remarkable, that a local solution exists if and only if H is positive definite and there is a solution to the equation $Hx + g = 0$. We now show how these results can be extended to problems with linear equality constraints.

3. Minimization of a Quadratic Function on an Affine Set

In this section we consider the problem

$$(54) \quad \begin{aligned} & \text{minimize } \frac{1}{2}x^T Hx + g^T x \\ & \text{subject to } Ax = b, \end{aligned}$$

where $H \in \mathbb{R}^{n \times n}$ is symmetric, $g \in \mathbb{R}^n$, $A \in \mathbb{R}^{m \times n}$, and $b \in \mathbb{R}^m$. We assume that the system $Ax = b$ is consistent. That is, there exists $\hat{x} \in \mathbb{R}^n$ such that $A\hat{x} = b$ in which case

$$\{x \mid Ax = b\} = \hat{x} + \text{Null}(A).$$

Consequently, the problem (54) is of the form

$$(55) \quad \text{minimize}_{x \in \hat{x} + S} \frac{1}{2}x^T Hx + g^T x,$$

where S is a subspace of \mathbb{R}^n . This representation of the problem shows that the problem (54) is trivial if $\text{Null}(A) = \{0\}$ since then the unique solution \hat{x} to $Ax = b$ is the unique solution to (54). Hence, when considering the problem (54) it is always assumed that $\text{Null}(A) \neq \{0\}$, and furthermore, that $m < n$.

DEFINITION 3.1. [Affine Sets] A subset K of \mathbb{R}^n is said to be an affine set if there exists a point $\hat{x} \in \mathbb{R}^n$ and a subspace $S \subset \mathbb{R}^n$ such that $K = \hat{x} + S = \{\hat{x} + u \mid u \in S\}$.

We now develop necessary and sufficient optimality conditions for the problem (55), that is, for the minimization of a quadratic function over an affine set. For this we assume that we have a basis v^1, v^2, \dots, v^k for S so that $\dim(S) = k$. Let $V \in \mathbb{R}^{n \times k}$ be the matrix whose columns are the vectors v^1, v^2, \dots, v^k so that $S = \text{Ran}(V)$. Then $\hat{x} + S = \{\hat{x} + Vz \mid z \in \mathbb{R}^k\}$. This allows us to rewrite the problem (55) as

$$(56) \quad \text{minimize}_{z \in \mathbb{R}^k} \frac{1}{2}(\hat{x} + Vz)^T H(\hat{x} + Vz) + g^T(\hat{x} + Vz).$$

PROPOSITION 3.1. Consider the two problems (55) and (56), where the columns of the matrix V form a basis for the subspace S . The set of optimal solution to these problems are related as follows:

$$\{\bar{x} \mid \bar{x} \text{ solves (55)}\} = \{\hat{x} + V\bar{z} \mid \bar{z} \text{ solves (56)}\}.$$

By expanding the objective function in (56), we obtain

$$\frac{1}{2}(\hat{x} + Vz)^T H(\hat{x} + Vz) + g^T(\hat{x} + Vz) = \frac{1}{2}z^T V^T H V z + (V^T(H\hat{x} + g))^T z + f(\hat{x}),$$

where $f(x) := \frac{1}{2}x^T Hx + g^T x$. If we now set $\hat{H} := V^T H V$, $\hat{g} := V^T(H\hat{x} + g)$, and $\beta := f(\hat{x})$, then problem (56) has the form of (42):

$$(57) \quad \text{minimize}_{z \in \mathbb{R}^k} \frac{1}{2}z^T \hat{H}z + \hat{g}^T z,$$

where, as usual, we have dropped the constant term $\beta = f(\hat{x})$. Since we have already developed necessary and sufficient conditions for optimality in this problem, we can use them to state similar conditions for the problem (55).

THEOREM 3.1. [Optimization of Quadratics on Affine Sets]

Consider the problem (55).

- (1) A local solution to the problem (55) exists if and only if $u^T H u \geq 0$ for all $u \in S$ and there exists a vector $\bar{x} \in \hat{x} + S$ such that $H\bar{x} + g \in S^\perp$, in which case \bar{x} is a local solution to (55).
- (2) If \bar{x} is a local solution to (55), then it is a global solution.
- (3) The problem (55) has a unique global solution if and only if $u^T H u > 0$ for all $u \in S \setminus \{0\}$. Moreover, if $V \in \mathbb{R}^{n \times k}$ is any matrix such that $\text{Ran}(V) = S$ where $k = \dim(S)$, then a unique solution to (55) exists if and only if the matrix $V^T H V$ is positive definite in which case the unique solution \bar{x} is given by

$$\bar{x} = [I - V(V^T H V)^{-1}V^T H]\hat{x} - V(V^T H V)^{-1}V^T g.$$

- (4) If either there exists $\bar{u} \in S$ such that $\bar{u}^T H \bar{u} < 0$ or there does not exist $\bar{x} \in \hat{x} + S$ such that $H\bar{x} + g \in S^\perp$ (or both), then

$$-\infty = \inf_{x \in \hat{x} + S} \frac{1}{2} x^T H x + g^T x .$$

PROOF. (1) By Proposition 3.1, a solution to (55) exists if and only if a solution to (56) exists. By Theorem 2.1, a solution to (56) exists if and only if $V^T H V$ is positive semi-definite and there is a solution \bar{z} to the equation $V^T(H(\hat{x} + Vz) + g) = 0$ in which case \bar{z} solves (56), or equivalently, by Proposition 3.1, $\bar{x} = \hat{x} + V\bar{z}$ solves (55). The condition that $V^T H V$ is positive semi-definite is equivalent to the statement that $z^T V^T H V z \geq 0$ for all $z \in \mathbb{R}^k$, or equivalently, $u^T H u \geq 0$ for all $u \in S$. The condition, $V^T(H(\hat{x} + V\bar{z}) + g) = 0$ is equivalent to $H\bar{x} + g \in \text{Null}(V^T) = \text{Ran}(V)^\perp = S^\perp$.

(2) This is an immediate consequence of Proposition 3.1 and Part (2) of Theorem 2.1.

(3) By Theorem 2.1, the problem (56) has a unique solution if and only if $V^T H V$ is positive definite in which case the solution is given by $\bar{z} = (V^T H V)^{-1} V^T (H\hat{x} + g)$. Note that $V^T H V$ is positive definite if and only if $u^T H u > 0$ for all $u \in S \setminus \{0\}$ which proves that this condition is necessary and sufficient. In addition, by Proposition 3.1, $\bar{x} = \hat{x} + V\bar{z} = [I - V(V^T H V)^{-1} V^T H] \hat{x} - V(V^T H V)^{-1} V^T g$ is the unique solution to (55).

(4) This follows the same pattern of proof using Part (4) of Theorem 2.1. \square

THEOREM 3.2. [*Optimization of Quadratics Subject to Linear Equality Constraints*]
Consider the problem (54).

- (1) A local solution to the problem (54) exists if and only if $u^T H u \geq 0$ for all $u \in \text{Null}(A)$ and there exists a vector pair $(\bar{x}, \bar{y}) \in \mathbb{R}^n \times \mathbb{R}^m$ such that $H\bar{x} + A^T \bar{y} + g = 0$, in which case \bar{x} is a local solution to (55).
- (2) If \bar{x} is a local solution to (55), then it is a global solution.
- (3) The problem (55) has a unique global solution if and only if $u^T H u > 0$ for all $u \in \text{Null}(A) \setminus \{0\}$.
- (4) If $u^T H u > 0$ for all $u \in \text{Null}(A) \setminus \{0\}$ and $\text{rank}(A) = m$, the matrix

$$M := \begin{bmatrix} H & A^T \\ A & 0 \end{bmatrix} \text{ is invertible, and the vector } \begin{bmatrix} \bar{x} \\ \bar{y} \end{bmatrix} = M^{-1} \begin{bmatrix} -g \\ b \end{bmatrix}$$

has \bar{x} as the unique global solution to (55).

- (5) If either there exists $\bar{u} \in \text{Null}(A)$ such that $\bar{u}^T H \bar{u} < 0$ or there does not exist a vector pair $(\bar{x}, \bar{y}) \in \mathbb{R}^n \times \mathbb{R}^m$ such that $H\bar{x} + A^T \bar{y} + g = 0$ (or both), then

$$-\infty = \inf_{x \in \hat{x} + S} \frac{1}{2} x^T H x + g^T x .$$

REMARK 3.1. The condition that $\text{rank}(A) = m$ in Part (4) of the theorem can always be satisfied by replacing A by first row reducing A to echelon form.

PROOF. (1) Recall that $\text{Null}(A)^\perp = \text{Ran}(A^T)$. Hence, $w \in \text{Null}(A)$ if and only if there exists $y \in \mathbb{R}^m$ such that $w = A^T y$. By setting $w = H\bar{x} + g$ the result follows from Part (1) of Theorem 3.1.

(2) Again, this is an immediate consequence of Proposition 3.1 and Part (2) of Theorem 2.1.

(3) This is just Part (3) of Theorem 3.1.

(4) Suppose $M \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$, then $Hx + A^T y = 0$ and $Ax = 0$. If we multiply $Hx + A^T y$ on the left by x^T , we obtain $0 = x^T Hx + x^T A^T y = x^T Hx$ which implies that $x = 0$ since $x \in \text{Null}(A)$. But then $A^T y = 0$, so that $y = 0$ since $\text{rank}(A) = m$. Consequently, $\text{Null}(M) = \{0\}$, i.e., M is invertible. The result now follows from Part (1).

(5) By Part (1), this is just a restatement of Theorem 3.1 Part (4). \square

The vector \bar{y} appearing in this Theorem is call a *Lagrange multiplier* vector. Lagrange multiplier vectors play an essential role in constrained optimization and lie at the heart of what is called *duality theory*. This theory is more fully developed in Chapter ??.

We now study how one might check when H is positive semi-definite as well as solving the equation $Hx + g = 0$ when H is positive semi-definite.

4. The Principal Minor Test for Positive Definiteness

Let $H \in \mathcal{S}^n$. We wish to obtain a test of when H is positive definite without having to compute its eigenvalue decomposition. First note that $H_{ii} = e_i^T H e_i$, so that H can be positive definite only if $H_{ii} > 0$. This is only a “sanity check” for whether a matrix is positive definite. That is, if any diagonal element of H is not positive, then H cannot be positive definite. In this section we develop a necessary and sufficient condition for H to be positive definite that makes use of the determinant. We begin with the following lemma.

LEMMA 4.1. *Let $H \in \mathcal{S}^n$, $u \in \mathbb{R}^n$, and $\alpha \in \mathbb{R}$, and consider the block matrix*

$$\hat{H} := \begin{bmatrix} H & u \\ u^T & \alpha \end{bmatrix} \in \mathcal{S}^{(n+1)}.$$

- (1) *The matrix \hat{H} is positive semi-definite if and only if H is positive semi-definite and there exists a vector $z \in \mathbb{R}^n$ such that $u = Hz$ and $\alpha \geq z^T Hz$.*
(2) *The matrix \hat{H} is positive definite if and only if H is positive definite and $\alpha > u^T H^{-1}u$.*

PROOF. (1) Suppose H is positive semi-definite, and there exists z such that $u = Hz$ and $\alpha \geq z^T Hz$. Then for any $\hat{x} = \begin{bmatrix} x \\ x_n \end{bmatrix}$ where $x_n \in \mathbb{R}$ and $x \in \mathbb{R}^n$, we have

$$\begin{aligned} \hat{x}^T \hat{H} \hat{x} &= x^T H x + 2x^T H x_n z + x_n^2 \alpha \\ &= (x + x_n z)^T H (x + x_n z) + x_n^2 (\alpha - z^T H z) \geq 0. \end{aligned}$$

Hence, \hat{H} is positive semi-definite.

Conversely, suppose that \hat{H} is positive semi-definite. Write $u = u_1 + u_2$ where $u_1 \in \text{Ran}(H)$ and $u_2 \in \text{Ran}(H)^\perp = \text{Null}(H)$, so that there is a $z \in \mathbb{R}^n$ such that $u_1 = Hz$. Then, for all $\hat{x} = \begin{pmatrix} x \\ x_n \end{pmatrix} \in \mathbb{R}^{(n+1)}$,

$$\begin{aligned} 0 \leq \hat{x}^T \hat{H} \hat{x} &= x^T H x + 2x_n u^T x + \alpha x_n^2 \\ &= x^T H x + 2x_n (u_1 + u_2)^T x + \alpha x_n^2 \\ &= x^T H x + 2x_n z^T H x + x_n^2 z^T H z + x_n^2 (\alpha - z^T H z) + 2x_n u_2^T x \\ &= (x + x_n z)^T H (x + x_n z) + x_n^2 (\alpha - z^T H z) + 2x_n u_2^T x. \end{aligned}$$

Taking $x_n = 0$ tells us that H is positive semi-definite, and taking $\hat{x} = \begin{pmatrix} -tu_2 \\ 1 \end{pmatrix}$ for $t \in \mathbb{R}$ gives

$$\alpha - 2t \|u_2\|_2^2 \geq 0 \quad \text{for all } t \in \mathbb{R},$$

which implies that $u_2 = 0$. Finally, taking $\hat{x} = \begin{pmatrix} -z \\ 1 \end{pmatrix}$, tells us that $z^T H z \leq \alpha$ which proves the result.

- (2) The proof follows the pattern of Part (1) but now we can take $z = H^{-1}u$. □

If the matrix H is invertible, we can apply a kind of block Gaussian elimination to the matrix \hat{H} in the lemma to obtain a matrix with block upper triangular structure:

$$\begin{bmatrix} I & 0 \\ (-H^{-1}u)^T & 1 \end{bmatrix} \begin{bmatrix} H & u \\ u^T & \alpha \end{bmatrix} = \begin{bmatrix} H & u \\ 0 & (\alpha - u^T H^{-1}u) \end{bmatrix}.$$

One consequence of this relationship is that

$$\begin{aligned} \det \begin{bmatrix} H & u \\ u^T & \alpha \end{bmatrix} &= \det \begin{bmatrix} I & 0 \\ (-H^{-1}u)^T & 1 \end{bmatrix} \det \begin{bmatrix} H & u \\ u^T & \alpha \end{bmatrix} \\ (58) \quad &= \det \left(\begin{bmatrix} I & 0 \\ (-H^{-1}u)^T & 1 \end{bmatrix} \begin{bmatrix} H & u \\ u^T & \alpha \end{bmatrix} \right) \\ &= \det \begin{bmatrix} H & u \\ 0 & (\alpha - u^T H^{-1}u) \end{bmatrix} \\ &= \det(H)(\alpha - u^T H^{-1}u). \end{aligned}$$

We use this determinant identity in conjunction with the previous lemma to establish a test for whether a matrix is positive definite based on determinants. The test requires us to introduce the following elementary definition.

DEFINITION 4.1. [Principal Minors] The k th principal minor of a matrix $B \in \mathbb{R}^{n \times n}$ is the determinant of the upper left-hand corner $k \times k$ -submatrix of B for $1 \leq k \leq n$.

PROPOSITION 4.1. [The Principal Minor Test] Let $H \in \mathcal{S}^n$. Then H is positive definite if and only if each of its principal minors is positive.

PROOF. The proof proceeds by induction on the dimension n of H . The result is clearly true for $n = 1$. We now assume the result is true for $1 \leq k \leq n$ and show it is true for dimension $n + 1$. Write

$$H := \begin{bmatrix} \hat{H} & u \\ u^T & \alpha \end{bmatrix}.$$

Then Lemma 4.1 tells us that H is positive definite if and only if \hat{H} is positive definite and $\alpha > u^T \hat{H}^{-1} u$. By the induction hypothesis, \hat{H} is positive definite if and only if all of its principal minors are positive. If we now combine this with the expression (58), we get that H is positive definite if and only if all principal minors of \hat{H} are positive and, by (58), $\det(H) = \det(\hat{H})(\alpha - u^T \hat{H}^{-1} u) > 0$, or equivalently, all principal minors of H are positive. \square

This result only applies to positive definite matrices, and does not provide insight into how to solve linear equations involving H such as $Hx + g = 0$. These two issues can be addressed through the Cholesky factorization.

5. The Cholesky Factorizations

We now consider how one might solve a quadratic optimization problem. Recall that a solution only exists when H is positive semi-definite and there is a solution to the equation $Hx + g = 0$. Let us first consider solving the equation when H is positive definite. We use a procedure similar to the LU factorization but which also takes advantage of symmetry.

Suppose

$$H = \begin{bmatrix} \alpha_1 & h_1^T \\ h_1 & \tilde{H}_1 \end{bmatrix}, \quad \text{where } \tilde{H}_1 \in \mathcal{S}^n.$$

Note that $\alpha_1 = e_1^T H e_1 > 0$ since H is positive definite (if $\alpha_1 \leq 0$, then H cannot be positive definite), so there is no need to apply a permutation. Multiply H on the left by the Gaussian elimination matrix for the first column, we obtain

$$L_1^{-1} H = \begin{bmatrix} 1 & 0 \\ -\frac{h_1}{\alpha_1} & I \end{bmatrix} \begin{bmatrix} \alpha_1 & h_1^T \\ h_1 & \tilde{H}_1 \end{bmatrix} = \begin{bmatrix} \alpha_1 & h_1^T \\ 0 & \tilde{H}_1 - \alpha_1^{-1} h_1 h_1^T \end{bmatrix}.$$

By symmetry, we have

$$L_1^{-1} H L_1^{-T} = \begin{bmatrix} \alpha_1 & h_1^T \\ 0 & \tilde{H}_1 - \alpha_1^{-1} h_1 h_1^T \end{bmatrix} \begin{bmatrix} 1 & -\frac{h_1^T}{\alpha_1} \\ 0 & I \end{bmatrix} = \begin{bmatrix} \alpha_1 & 0 \\ 0 & \tilde{H}_1 - \alpha_1^{-1} h_1 h_1^T \end{bmatrix}.$$

Set $H_1 = \tilde{H}_1 - \alpha_1^{-1} h_1 h_1^T$. Observe that for every non-zero vector $v \in \mathbb{R}^{(n-1)}$,

$$v^T H_1 v = \begin{pmatrix} 0 \\ v \end{pmatrix}^T \begin{bmatrix} \alpha_1 & 0 \\ 0 & H_1 \end{bmatrix} \begin{pmatrix} 0 \\ v \end{pmatrix} = \left(L_1^{-T} \begin{pmatrix} 0 \\ v \end{pmatrix} \right)^T H \left(L_1^{-T} \begin{pmatrix} 0 \\ v \end{pmatrix} \right) > 0,$$

which shows that H_1 is positive definite. Decomposing H_1 as we did H gives

$$H_1 = \begin{bmatrix} \alpha_2 & h_2^T \\ h_2 & \tilde{H}_2 \end{bmatrix}, \quad \text{where } \tilde{H}_2 \in \mathcal{S}^{(n-1)}.$$

Again, $\alpha_2 > 0$ since H_1 is positive definite (if $\alpha_2 \leq 0$, then H cannot be positive definite). Hence, can repeat the reduction process for H_1 . But if at any stage we discover and $\alpha_i \leq 0$, then we terminate, since H cannot be positive definite.

If we can continue this process n times, we will have constructed a lower triangular matrix

$$L := L_1 L_2 \cdots L_n \quad \text{such that} \quad L^{-1} H L^{-T} = D, \quad \text{where} \quad D := \text{diag}(\alpha_1, \alpha_2, \dots, \alpha_n)$$

is a diagonal matrix with strictly positive diagonal entries. On the other hand, if at some point in the process we discover an α_i that is not positive, then H cannot be positive definite and the process terminates. That is, this computational procedure simultaneously tests whether H is positive definite as it tries to diagonalize H . We will call this process the *Cholesky diagonalization procedure*. It is used to establish the following factorization theorem.

THEOREM 5.1. *[The Cholesky Factorization] Let $H \in \mathcal{S}_+^n$ have rank k . Then there is a lower triangular matrix $L \in \mathbb{R}^{n \times k}$ such that $H = LL^T$. Moreover, if the rank of H is n , then there is a positive diagonal matrix D and a lower triangular matrix \tilde{L} with ones on its diagonal such that $H = \tilde{L}D\tilde{L}^T$.*

PROOF. Let the columns of the matrix $V_1 \in \mathbb{R}^{n \times k}$ be an orthonormal basis for $\text{Ran}(H)$ and the columns of $V_2 \in \mathbb{R}^{n \times (n-k)}$ be an orthonormal basis for $\text{Null}(H)$ and set $V = [V_1 \ V_2] \in \mathbb{R}^{n \times n}$. Then

$$\begin{aligned} V^T H V &= \begin{bmatrix} V_1^T \\ V_2^T \end{bmatrix} H [V_1 \ V_2] \\ &= \begin{bmatrix} V_1^T H V_1 & V_1^T H V_2 \\ V_2^T H V_1 & V_2^T H V_2 \end{bmatrix} \\ &= \begin{bmatrix} V_1^T H V_1 & 0 \\ 0 & 0 \end{bmatrix}. \end{aligned}$$

Since $\text{Ran}(H) = \text{Null}(H^T)^\perp = \text{Null}(H)^\perp$, $V_1 H V_1^T \in \mathbb{R}^{k \times k}$ is symmetric and positive definite. By applying the procedure described prior to the statement of the theorem, we construct a nonsingular lower triangular matrix $\tilde{L} \in \mathbb{R}^{k \times k}$ and a diagonal matrix $D = \text{diag}(\alpha_1, \alpha_2, \dots, \alpha_k)$, with $\alpha_i > 0$, $i = 1, \dots, k$, such that $V_1 H V_1^T = \tilde{L}D\tilde{L}^T$. Set $\hat{L} = \tilde{L}D^{1/2}$ so that $V_1 H V_1^T = \hat{L}\hat{L}^T$. If $k = n$, taking $V = I$ proves the theorem by setting $L = \hat{L}$. If $k < n$,

$$H = [V_1 \ V_2] \begin{bmatrix} \hat{L}\hat{L}^T & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} V_1^T \\ V_2^T \end{bmatrix} = (V_1 \hat{L})(V_1 \hat{L})^T.$$

Let $(V_1 \hat{L})^T \in \mathbb{R}^{k \times n}$ have reduced QR factorization $(V_1 \hat{L})^T = QR$ (see Theorem 5.1). Since \hat{L}^T has rank k , $Q \in \mathbb{R}^{k \times k}$ is unitary and $R = [R_1 \ R_2]$ with $R_1 \in \mathbb{R}^{k \times k}$ nonsingular and $R_2 \in \mathbb{R}^{k \times (n-k)}$. Therefore,

$$H = (V_1 \hat{L})(V_1 \hat{L})^T = R^T Q^T Q R = R^T R.$$

The theorem follows by setting $L = R^T$. \square

When H is positive definite, the factorization $H = LL^T$ is called the Cholesky factorization of H , and when $\text{rank}(H) < n$ it is called the *generalized Cholesky factorization* of H . In the positive definite case, the Cholesky diagonalization procedure computes the Cholesky factorization of H . On the other hand, when H is only positive semi-definite, the proof of the theorem provides a guide for obtaining the generalized Cholesky factorization.

5.1. Computing the Generalized Cholesky Factorization.

Step 1: Initiate the Cholesky diagonalization procedure. If the procedure successfully completes n iterations, the Cholesky factorization has been obtained. Otherwise the procedure terminates at some iteration $k+1 < n$.

If $\alpha_{k+1} < 0$, proceed no further since the matrix H is not positive semi-definite. If $\alpha_{k+1} = 0$, proceed to Step 2.

Step 2: In Step 1, the factorization

$$\hat{L}^{-1} H \hat{L}^{-T} = \begin{bmatrix} \hat{D} & 0 \\ 0 & \hat{H} \end{bmatrix},$$

where

$$\hat{L} = \begin{bmatrix} \hat{L}_1 & 0 \\ \hat{L}_2 & I_{(n-k)} \end{bmatrix}$$

with $\hat{L}_1 \in \mathbb{R}^{k \times k}$ lower triangular with ones on the diagonal, $\hat{D} = \text{diag}(\alpha_1, \alpha_2, \dots, \alpha_k) \in \mathbb{R}^{k \times k}$ with $\alpha_i > 0$ $i = 1, \dots, k$, and $\hat{H} \in \mathbb{R}^{(n-k) \times (n-k)}$ with \hat{H} symmetric has a nontrivial null space. Let the full QR factorization of \hat{H} be given by

$$\hat{H} = [U_1 \ U_2] \begin{bmatrix} R_1 & R_2 \\ 0 & 0 \end{bmatrix} = U \begin{bmatrix} R \\ 0 \end{bmatrix},$$

where

- $U = [U_1 \ U_2] \in \mathbb{R}^{k \times k}$ is unitary,
- the columns of $U_1 \in \mathbb{R}^{k \times k_1}$ form an orthonormal basis for $\text{Ran}(\hat{H})$ with $k_1 = \text{rank}(\hat{H}) < k$,
- the columns of $U_2 \in \mathbb{R}^{k \times (k-k_1)}$ form an orthonormal basis for $\text{Null}(\hat{H})$,
- $R_1 \in \mathbb{R}^{k_1 \times k_1}$ is upper triangular and nonsingular,

- $R_2 \in \mathbb{R}^{k_1 \times (k-k_1)}$, and
- $R = [R_1 \ R_2] \in \mathbb{R}^{k_1 \times k}$.

Consequently,

$$\begin{aligned} \begin{bmatrix} U_1^T \widehat{H}U_1 & 0 \\ 0 & 0 \end{bmatrix} &= \begin{bmatrix} U_1^T \widehat{H}U_1 & U_1^T \widehat{H}U_2 \\ U_2^T \widehat{H}U_1 & U_2^T \widehat{H}U_2 \end{bmatrix} \\ &= U^T \widehat{H}U \\ &= \begin{bmatrix} R \\ 0 \end{bmatrix} [U_1 \ U_2] \\ &= \begin{bmatrix} RU_1 & RU_2 \\ 0 & 0 \end{bmatrix}, \end{aligned}$$

and so $RU_2 = 0$ and $U_1^T \widehat{H}U_1 = RU_1 \in \mathbb{R}^{k_1 \times k_1}$ is a nonsingular symmetric matrix.

Note that only the reduced QR factorization of $H = U_1 R$ is required since $U_1^T \widehat{H}U_1 = RU_1$.

Step 4: Initiate the Cholesky diagonalization procedure on $U_1^T \widehat{H}U_1$. If the procedure successfully completes k_1 iterations, the Cholesky factorization

$$U_1^T \widehat{H}U_1 = \widehat{L}_3 \widehat{L}_3^T$$

has been obtained. If this does not occur, the procedure terminates at some iteration $j < k_1$ with $\alpha_j < 0$ since $U_1^T \widehat{H}U_1$ is nonsingular. In this case, terminate the process since H cannot be positive semi-definite. Otherwise proceed to Step 5.

Step 5: We now have

$$\begin{aligned} H &= \begin{bmatrix} \widehat{L}_1 & 0 \\ \widehat{L}_2 & I_{(n-k)} \end{bmatrix} \begin{bmatrix} \widehat{D} & 0 \\ 0 & \widehat{H} \end{bmatrix} \begin{bmatrix} \widehat{L}_1^T & \widehat{L}_2^T \\ 0 & I_{(n-k)} \end{bmatrix} \\ &= \begin{bmatrix} \widehat{L}_1 & 0 \\ \widehat{L}_2 & I_{(n-k)} \end{bmatrix} \begin{bmatrix} I & 0 \\ 0 & U \end{bmatrix} \begin{bmatrix} \widehat{D} & 0 \\ 0 & U^T \widehat{H}U \end{bmatrix} \begin{bmatrix} I & 0 \\ 0 & U^T \end{bmatrix} \begin{bmatrix} \widehat{L}_1^T & \widehat{L}_2^T \\ 0 & I_{(n-k)} \end{bmatrix} \\ &= \begin{bmatrix} \widehat{L}_1 & 0 & 0 \\ \widehat{L}_2 & U_1 & U_2 \end{bmatrix} \begin{bmatrix} \widehat{D} & 0 & 0 \\ 0 & \widehat{L}_3 \widehat{L}_3^T & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \widehat{L}_1^T & \widehat{L}_2^T \\ 0 & U_1^T \\ 0 & U_2^T \end{bmatrix} \\ &= \begin{bmatrix} \widetilde{L}_1 \widehat{D}^{1/2} & 0 & 0 \\ \widehat{L}_2 \widehat{D}^{1/2} & U_1 \widehat{L}_3 & 0 \end{bmatrix} \begin{bmatrix} \widehat{D}^{1/2} \widetilde{L}_1^T & \widehat{D}^{1/2} \widehat{L}_2^T \\ 0 & \widehat{L}_3^T U_1^T \\ 0 & 0 \end{bmatrix} \\ &= \begin{bmatrix} L_1 & 0 \\ L_2 & U_1 \widehat{L}_3 \end{bmatrix} \begin{bmatrix} L_1^T & L_2^T \\ 0 & \widehat{L}_3^T U_1^T \end{bmatrix}, \end{aligned}$$

where $L_1 = \widetilde{L}_1 \widehat{D}^{1/2} \in \mathbb{R}^{k \times k}$ is lower triangular, $L_2 = \widehat{L}_2 \widehat{D}^{1/2} \in \mathbb{R}^{(n-k) \times k}$, and $U_1 \widehat{L}_3 \in \mathbb{R}^{(n-k) \times k_1}$. In particular, $k + k_1 = \text{rank}(H)$ since L_1 has rank k and $U_1 \widehat{L}_3$ has rank k_1 . Let $\widehat{L}_3^T U_1^T$ have QR factorization $\widehat{L}_3^T U_1^T = VL_3^T$, where $V \in \mathbb{R}^{k_1 \times k_1}$ is unitary and $L_3 \in \mathbb{R}^{k_1 \times (n-k)}$ is lower triangular. Then

$$H = \begin{bmatrix} L_1 & 0 \\ L_2 & U_1 \widehat{L}_3 \end{bmatrix} \begin{bmatrix} L_1^T & L_2^T \\ 0 & \widehat{L}_3^T U_1^T \end{bmatrix} = \begin{bmatrix} L_1 & 0 \\ L_2 & L_3 V^T \end{bmatrix} \begin{bmatrix} L_1^T & L_2^T \\ 0 & VL_3^T \end{bmatrix} = \begin{bmatrix} L_1 & 0 \\ L_2 & L_3 \end{bmatrix} \begin{bmatrix} L_1^T & L_2^T \\ 0 & L_3^T \end{bmatrix},$$

since $V^T V = I_{k_1}$. This is the generalized Cholesky factorization of H .

5.2. Computing Solutions to the Quadratic Optimization Problem via Cholesky Factorizations.

Step 1: Apply the procedure described in the previous section for computing the generalized Cholesky factorization of H . If it is determined that H is not positive definite, then proceed no further since the problem (42) has no solution and the optimal value is $-\infty$.

Step 2: Step 1 provides us with the generalized Cholesky factorization for $H = LL^T$ with $L^T = [L_1^T \ L_2^T]$, where $L_1 \in \mathbb{R}^{k \times k}$ and $L_2 \in \mathbb{R}^{(n-k) \times k}$ with $k = \text{rank}(H)$. Write

$$g = \begin{pmatrix} g_1 \\ g_2 \end{pmatrix},$$

where $g_1 \in \mathbb{R}^k$ and $g_2 \in \mathbb{R}^{(n-k)}$. Since $\text{Ran}(H) = \text{Ran}(L)$, the system $Hx + g = 0$ is solvable if and only if $-g \in \text{Ran}(L)$. That is, there exists $w \in \mathbb{R}^k$ such that $Lw = -g$, or equivalently,

$$L_1 w = -g_1 \quad \text{and} \quad L_2 w = -g_2.$$

Since L_1 is invertible, the system $L_1 w = -g_1$ has as its unique solution $\bar{w} = L_1^{-1} g_1$. Note that \bar{w} is easy to compute by forward substitution since L_1 is lower triangular. Having \bar{w} check to see if $L_2 \bar{w} = -g_2$. If this is not the case, then proceed no further, since the system $Hx + g = 0$ has no solution and so the optimal value in (42) is $-\infty$. Otherwise, proceed to Step 3.

Step 3: Use back substitution to solve the equation $L_1^T y = \bar{w}$ for $\bar{y} := L_1^{-T} \bar{w}$ and set

$$\bar{x} = \begin{pmatrix} \bar{y} \\ 0 \end{pmatrix}.$$

Then

$$H\bar{x} = LL^T \bar{x} = \begin{bmatrix} L_1 \\ L_2 \end{bmatrix} [L_1^T \ L_2^T] \begin{pmatrix} \bar{y} \\ 0 \end{pmatrix} = \begin{bmatrix} L_1 \\ L_2 \end{bmatrix} \bar{w} = -g.$$

Hence, \bar{x} solves the equation $Hx + g = 0$ and so is an optimal solution to the quadratic optimization problem (42).

6. Linear Least Squares Revisited

We have already see that the least squares problem is a special case of the problem of minimizing a quadratic function. But what about the reverse? Part (4) of Theorem 2.1 tells us that, in general, the reverse cannot be true since the linear least squares problem always has a solution. But what about the case when the quadratic optimization problem has a solution? In this case the matrix H is necessarily positive semi-definite and a solution to the system $Hx + g = 0$ exists. By Theorem 5.1, there is a lower triangular matrix $L \in \mathbb{R}^{n \times k}$, where $k = \text{rank}(H)$, such that $H = LL^T$. Set $A := L^T$. In particular, this implies that $\text{Ran}(H) = \text{Ran}(L) = \text{Ran}(A^T)$. Since $Hx + g = 0$, we know that $-g \in \text{Ran}(H) = \text{Ran}(A^T)$, and so there is a vector $b \in \mathbb{R}^k$ such that $-g = A^T b$. Consider the linear least squares problem

$$\min_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - b\|_2^2.$$

As in (44), expand the objective in this problem to obtain

$$\begin{aligned} \frac{1}{2} \|Ax - b\|_2^2 &= \frac{1}{2} x^T (A^T A) x - (A^T b)^T x + \frac{1}{2} \|b\|_2^2 \\ &= \frac{1}{2} x^T LL^T x + g^T x + \beta \\ &= \frac{1}{2} x^T Hx + g^T x + \beta, \end{aligned}$$

where $\beta = \frac{1}{2} \|b\|_2^2$. We have just proved the following result.

PROPOSITION 6.1. *A quadratic optimization problem of the form (42) has an optimal solution if and only if it is equivalent to a linear least squares problem.*

7. The Conjugate Gradient Algorithm

The Cholesky factorization is an important and useful tool for computing solutions to the quadratic optimization problem, but it is too costly to be employed in many very large scale applications. In some applications, the matrix H is too large to be stored or it is not available as a data structure. However, in these problems it is often the case that the matrix vector product Hx can be obtained for a given vector $x \in \mathbb{R}^n$. This occurs, for example, in a signal processing applications. In this section, we develop an algorithm for solving the quadratic optimization problem (47) that only requires access to the matrix vector products Hx . Such an algorithm is called a *matrix free* method since knowledge the whole matrix H is not required. In such cases the Cholesky factorization is inefficient to compute. The focus of this section is the study of the matrix free method known as the *conjugate gradient algorithm*. Throughout this section we assume that the matrix H is positive definite.

7.1. Conjugate Direction Methods. Consider the problem (47) where it is known that H is symmetric and positive definite. In this case it is possible to define a notion of *orthogonality* or *conjugacy* with respect to H .

DEFINITION 7.1 (Conjugacy). Let $H \in \mathcal{S}_{++}^n$. We say that the vectors $x, y \in \mathbb{R}^n \setminus \{0\}$ are H -conjugate (or H -orthogonal) if $x^T H y = 0$.

PROPOSITION 7.1. [Conjugacy implies Linear Independence]

If $H \in \mathcal{S}_{++}^n$ and the set of nonzero vectors d^0, d^1, \dots, d^k are (pairwise) H -conjugate, then these vectors are linearly independent.

PROOF. If $0 = \sum_{i=0}^k \mu_i d^i$, then for $\bar{i} \in \{0, 1, \dots, k\}$

$$0 = (d^{\bar{i}})^T H \left[\sum_{i=0}^k \mu_i d^i \right] = \mu_{\bar{i}} (d^{\bar{i}})^T H d^{\bar{i}},$$

Hence $\mu_i = 0$ for each $i = 0, \dots, k$. □

Let $x^0 \in \mathbb{R}^n$ and suppose that the vectors $d^0, d^1, \dots, d^{k-1} \in \mathbb{R}^n$ are H -conjugate. Set $S = \text{Span}(d^0, d^1, \dots, d^{k-1})$. Theorem 3.1 tells us that there is a unique optimal solution \bar{x} to the problem $\min \{ \frac{1}{2} x^T H x + g^T x \mid x \in x^0 + S \}$, and that \bar{x} is uniquely identified by the condition $H\bar{x} + g \in S^\perp$, or equivalently, $0 = (d^j)^T (H\bar{x} + g)$, $j = 0, 1, \dots, k-1$. Since $\bar{x} \in x^0 + S$, there are scalars μ_0, \dots, μ_{k-1} such that

$$(59) \quad \bar{x} = x^0 + \mu_0 d^0 + \dots + \mu_{k-1} d^{k-1},$$

and so, for each $j = 0, 1, \dots, k-1$,

$$\begin{aligned} 0 &= (d^j)^T (H\bar{x} + g) \\ &= (d^j)^T (H(x^0 + \mu_0 d^0 + \dots + \mu_{k-1} d^{k-1}) + g) \\ &= (d^j)^T (Hx^0 + g) + \mu_0 (d^j)^T H d^0 + \dots + \mu_{k-1} (d^j)^T H d^{k-1} \\ &= (d^j)^T (Hx^0 + g) + \mu_j (d^j)^T H d^j. \end{aligned}$$

Therefore,

$$(60) \quad \mu_j = \frac{-(Hx^0 + g)^T (d^j)}{(d^j)^T H d^j} \quad j = 0, 1, \dots, k-1.$$

This observation motivates the following theorem.

THEOREM 7.1. [Expanding Subspace Theorem]

Consider the problem (47) with $H \in \mathcal{S}_{++}^n$, and set $f(x) = \frac{1}{2} x^T H x + g^T x$. Let $\{d^i\}_{i=0}^{n-1}$ be a sequence of nonzero H -conjugate vectors in \mathbb{R}^n . Then, for any $x^0 \in \mathbb{R}^n$ the sequence $\{x^k\}$ generated according to

$$x^{k+1} := x^k + t_k d^k,$$

with

$$t_k := \arg \min \{ f(x^k + t d^k) : t \in \mathbb{R} \},$$

has the property that $f(x) = \frac{1}{2} x^T H x + g^T x$ attains its minimum value on the affine set $x^0 + \text{Span} \{d^0, \dots, d^{k-1}\}$ at the point x^k . In particular, if $k = n$, then x^n is the unique global solution to the problem (47).

PROOF. Let us first compute the value of the t_k 's. For $j = 0, \dots, k-1$, define $\varphi_j : \mathbb{R} \rightarrow \mathbb{R}$ by

$$\begin{aligned} \varphi_j(t) &= f(x^j + t d^j) \\ &= \frac{t^2}{2} (d^j)^T H d^j + t (g^j)^T d^j + f(x^j), \end{aligned}$$

where $g^j = Hx^j + g$. Then, for $j = 0, \dots, k-1$, $\varphi_j'(t) = t (d^j)^T H d^j + (g^j)^T d^j$ and $\varphi_j''(t) = (d^j)^T H d^j > 0$. Since $\varphi_j''(t) > 0$, our one dimensional calculus tells us that φ_j attains its global minimum value at the unique solution t_j to the equation $\varphi_j'(t) = 0$, i.e.,

$$t_j = -\frac{(g^j)^T d^j}{(d^j)^T H d^j}.$$

Therefore,

$$x^k = x^0 + t_0 d^0 + t_1 d^1 + \dots + t_k d^k$$

with

$$t_j = -\frac{(g^j)^T d^j}{(d^j)^T H d^j}, \quad j = 0, 1, \dots, k.$$

In the discussion preceding the theorem it was shown that if \bar{x} is the solution to the problem

$$\min \{ f(x) \mid x \in x^0 + \text{Span}(d^0, d^1, \dots, d^k) \},$$

then \bar{x} is given by (59) and (60). Therefore, if we can now show that $\mu_j = t_j$, $j = 0, 1, \dots, k$, then $\bar{x} = x_k$ proving the result. For each $j \in \{0, 1, \dots, k\}$ we have

$$\begin{aligned} (g^j)^T d^j &= (Hx^j + g)^T d^j \\ &= (H(x^0 + t_0 d^0 + t_1 d^1 + \dots + t_{j-1} d^{j-1}) + g)^T d^j \\ &= (Hx^0 + g)^T d^j + t_0 (d^0)^T H d^j + t_1 (d^1)^T H d^j + \dots + t_{j-1} (d^{j-1})^T H d^j \\ &= (Hx^0 + g)^T d^j \\ &= (g^0)^T d^j. \end{aligned}$$

Therefore, for each $j \in \{0, 1, \dots, k\}$,

$$t_j = \frac{-(g^j)^T d^j}{(d^j)^T H d^j} = \frac{-(g^0)^T d^j}{(d^j)^T H d^j} = \mu_j,$$

which proves the result. \square

7.2. The Conjugate Gradient Algorithm. The major drawback of the Conjugate Direction Algorithm of the previous section is that it seems to require that a set of H -conjugate directions must be obtained before the algorithm can be implemented. This is in opposition to our working assumption that H is so large that it cannot be kept in storage since any set of H -conjugate directions requires the same amount of storage as H . However, it is possible to generate the directions d^j one at a time and then discard them after each iteration of the algorithm. One example of such an algorithm is the Conjugate Gradient Algorithm.

The C-G Algorithm:

Initialization: $x^0 \in \mathbb{R}^n$, $d^0 = -g^0 = -(Hx^0 + g)$.

For $k = 0, 1, 2, \dots$

$$\begin{aligned} t_k &:= -(g^k)^T d^k / (d^k)^T H d^k \\ x^{k+1} &:= x^k + t_k d^k \\ g^{k+1} &:= Hx^{k+1} + g && \text{(STOP if } g^{k+1} = 0) \\ \beta_k &:= (g^{k+1})^T H d^k / (d^k)^T H d^k \\ d^{k+1} &:= -g^{k+1} + \beta_k d^k \\ k &:= k + 1. \end{aligned}$$

THEOREM 7.2. [CONJUGATE GRADIENT THEOREM]

The C-G algorithm is a conjugate direction method. If it does not terminate at x^k (i.e. $g^k \neq 0$), then

- (1) $\text{Span}[g^0, g^1, \dots, g^k] = \text{span}[g^0, Hg^0, \dots, H^k g^0]$
- (2) $\text{Span}[d^0, d^1, \dots, d^k] = \text{span}[g^0, Hg^0, \dots, H^k g^0]$
- (3) $(d^k)^T H d^i = 0$ for $i \leq k-1$
- (4) $t_k = (g^k)^T g^k / (d^k)^T H d^k$
- (5) $\beta_k = (g^{k+1})^T g^{k+1} / (g^k)^T g^k$.

PROOF. We first prove (1)-(3) by induction. The results are clearly true for $k = 0$. Now suppose they are true for k , we show they are true for $k + 1$. First observe that

$$g^{k+1} = g^k + t_k H d^k$$

so that $g^{k+1} \in \text{Span}[g^0, \dots, H^{k+1} g^0]$ by the induction hypothesis on (1) and (2). Also $g^{k+1} \notin \text{Span}[d^0, \dots, d^k]$, otherwise, by Theorem 3.1 Part (1), $g^{k+1} = Hx^{k+1} + g = 0$ since the method is a conjugate direction method up to step k by the induction hypothesis. Hence $g^{k+1} \notin \text{Span}[g^0, \dots, H^k g^0]$ and so $\text{Span}[g^0, g^1, \dots, g^{k+1}] = \text{Span}[g^0, \dots, H^{k+1} g^0]$, which proves (1).

To prove (2) write

$$d^{k+1} = -g^{k+1} + \beta_k d^k$$

so that (2) follows from (1) and the induction hypothesis on (2).

To see (3) observe that

$$(d^{k+1})^T H d^i = -(g^{k+1})^T H d^i + \beta_k (d^k)^T H d^i.$$

For $i = k$ the right hand side is zero by the definition of β_k . For $i < k$ both terms vanish. The term $(g^{k+1})^T H d^i = 0$ by Theorem 7.1 since $H d^i \in \text{Span}[d^0, \dots, d^k]$ by (1) and (2). The term $(d^k)^T H d^i$ vanishes by the induction hypothesis on (3).

To prove (4) write

$$-(g^k)^T d^k = (g^k)^T g^k - \beta_{k-1} (g^k)^T d^{k-1}$$

where $(g^k)^T d^{k-1} = 0$ by Theorem 7.1.

To prove (5) note that $(g^{k+1})^T g^k = 0$ by Theorem 7.1 because $g^k \in \text{Span}[d^0, \dots, d^k]$. Hence

$$(g^{k+1})^T H d^k = \frac{1}{t_k} (g^{k+1})^T [g^{k+1} - g^k] = \frac{1}{t_k} (g^{k+1})^T g^{k+1}.$$

Therefore,

$$\beta_k = \frac{1}{t_k} \frac{(g^{k+1})^T g^{k+1}}{(d^k)^T H d^k} = \frac{(g^{k+1})^T g^{k+1}}{(g^k)^T g^k}.$$

□

Remarks:

- (1) The C-G method is an example of a *descent method* since the values

$$f(x^0), f(x^1), \dots, f(x^n)$$

form a decreasing sequence.

- (2) It should be observed that due to the occurrence of round-off error the C-G algorithm is best implemented as an iterative method. That is, at the end of n steps, x^n may not be the global optimal solution and the intervening directions d^k may not be H -conjugate. Consequently, the algorithm is usually iterated until $\|g^k\|_2$ is sufficiently small. Due to the observations in the previous remark, this approach is guaranteed to continue to reduce the function value if possible since the overall method is a descent method. In this sense the C-G algorithm is self correcting.

Elements of Multivariable Calculus

1. Norms and Continuity

As we have seen the 2-norm gives us a measure of the magnitude of a vector v in \mathbb{R}^n , $\|v\|_2$. As such it also gives us a measure of the distance between two vectors $u, v \in \mathbb{R}^n$, $\|u - v\|_2$. Such measures of magnitude and distance are very useful tools for measuring model misfit as is the case in linear least squares problem. They are also essential for analyzing the behavior of sequences and functions on \mathbb{R}^n as well as on the space of matrices $\mathbb{R}^{m \times n}$. For this reason, we formalize the notion of a norm to incorporate other measures of magnitude and distance.

DEFINITION 1.1. [Vector Norm] A function $\|\cdot\| : \mathbb{R}^n \rightarrow \mathbb{R}$ is a vector norm on \mathbb{R}^n if

- (1) $\|x\| \geq 0$ for all $x \in \mathbb{R}^n$ with equality if and only if $x = 0$,
- (2) $\|\alpha x\| = |\alpha| \|x\|$ for all $x \in \mathbb{R}^n$ and $\alpha \in \mathbb{R}$, and
- (3) $\|x + y\| \leq \|x\| + \|y\|$ for all $x, y \in \mathbb{R}^n$.

EXAMPLE 1.1. Perhaps the most common examples of norms are the p -norms for $1 \leq p \leq \infty$. Given $1 \leq p < \infty$, the ℓ_p -norm on \mathbb{R}^n is defined as

$$\|x\|_p := \left[\sum_{j=1}^n |x_j|^p \right]^{1/p}.$$

For $p = \infty$, we define

$$\|x\|_\infty := \max \{ |x_i| \mid i = 1, 2, \dots, n \}.$$

This choice of notation for the ∞ -norm comes from the relation

$$\lim_{p \uparrow \infty} \|x\|_p = \|x\|_\infty \quad \forall x \in \mathbb{R}^n.$$

In applications, the most important of these norms are the $p = 1, 2, \infty$ norms as well as variations on these norms.

In finite dimensions all norms are said to be *equivalent* in the sense that one can show that for any two norms $\|\cdot\|_{(a)}$ and $\|\cdot\|_{(b)}$ on \mathbb{R}^n there exist positive constants α and β such that

$$\alpha \|x\|_a \leq \|x\|_b \leq \beta \|x\|_a \quad \forall x \in \mathbb{R}^n.$$

But we caution that in practice the numerical behavior of these norms differ greatly when the dimension is large.

Since norms can be used to measure the distance between vectors, they can be used to form notions of continuity for functions mapping \mathbb{R}^n to \mathbb{R}^m that parallel those established for mappings from \mathbb{R} to \mathbb{R} .

DEFINITION 1.2. [Continuous Functions] Let $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$.

- (1) F is said to be *continuous at a point* $\bar{x} \in \mathbb{R}^n$ if for all $\epsilon > 0$ there is a $\delta > 0$ such that

$$\|F(x) - F(\bar{x})\| \leq \epsilon \quad \text{whenever} \quad \|x - \bar{x}\| \leq \delta.$$

- (2) F is said to be *continuous on a set* $S \subset \mathbb{R}^n$ if it is continuous at every point of S .
- (3) The function F is said to be *continuous relative to a set* $S \subset \mathbb{R}^n$ if

$$\|F(x) - F(\bar{x})\| \leq \epsilon \quad \text{whenever} \quad \|x - \bar{x}\| \leq \delta \quad \text{and} \quad x \in S.$$

- (4) The function F is said to be *uniformly continuous on a set* $S \subset \mathbb{R}^n$ if for all $\epsilon > 0$ there is a $\delta > 0$ such that

$$\|F(x) - F(y)\| \leq \epsilon \quad \text{whenever} \quad \|x - y\| \leq \delta \quad \text{and} \quad x, y \in S.$$

Norms allow us to define certain topological notions that are very helpful in analyzing the behavior of sequences and functions. Since we will make frequent use of these concepts, it is helpful to have certain notational conventions associated with norms. We list a few of these below:

$$\begin{aligned} \text{the closed unit ball} & \quad \mathbb{B} := \{x \mid \|x\| \leq 1\} \\ \text{the unit vectors} & \quad \mathbb{S} := \{x \mid \|x\| = 1\} \\ \epsilon\text{-ball about } \bar{x} & \quad \bar{x} + \epsilon\mathbb{B} := \{x + \epsilon u \mid u \in \mathbb{B}\} = \{x \mid \|x - \bar{x}\| \leq \epsilon\} \end{aligned}$$

The unit ball associated with the 1, 2, and ∞ norms will be denoted by \mathbb{B}_1 , \mathbb{B}_2 , and \mathbb{B}_∞ , respectively.

A few basic topological notions are listed in the following definition. The most important of these for our purposes is *compactness*.

DEFINITION 1.3. *Let S be a subset of \mathbb{R}^n , and let $\|\cdot\|$ be a norm on \mathbb{R}^n .*

- (1) *The set S is said to be an open set if for every $\bar{x} \in S$ there is an $\epsilon > 0$ such that $\bar{x} + \epsilon\mathbb{B} \subset S$.*
- (2) *The set S is said to be a closed set if S contains every point $\bar{x} \in \mathbb{R}^n$ for which there is a sequence $\{x^k\} \subset S$ with $\lim_{k \rightarrow \infty} \|x^k - \bar{x}\| = 0$.*
- (3) *The set S is said to be a bounded set if there is a $\beta > 0$ such that $S \subset \beta\mathbb{B}$.*
- (4) *The set S is said to be a compact set if it is both closed and bounded.*
- (5) *A point $\bar{x} \in \mathbb{R}^n$ is a cluster point of the set S if there is a sequence $\{x^k\} \subset S$ with $\lim_{k \rightarrow \infty} \|x^k - \bar{x}\| = 0$.*
- (6) *A point $\bar{x} \in \mathbb{R}^n$ is said to be a boundary point of the set S if for all $\epsilon > 0$, $(\bar{x} + \epsilon\mathbb{B}) \cap S \neq \emptyset$ while $(\bar{x} + \epsilon\mathbb{B}) \not\subset S$, i.e., every ϵ ball about \bar{x} contains points that are in S and points that are not in S .*

The importance of the notion of compactness in optimization is illustrated in following basic theorems from analysis that we make extensive use of, but do not prove.

THEOREM 1.1. [*Compactness implies Uniform Continuity*] *Let $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be a continuous function on an open set $S \subset \mathbb{R}^n$. Then F is uniformly continuous on every compact subset of S .*

THEOREM 1.2. [*Weierstrass Compactness Theorem*] *A set $D \subset \mathbb{R}^n$ is compact if and only if every infinite sequence in D has a cluster point in D .*

THEOREM 1.3. [*Weierstrass Extreme Value Theorem*] *Every continuous function on a compact set attains its extreme values on that set. That is, there are points in the set at which both the infimum and the supremum of the function relative to the set are attained.*

We will also have need of a norm on the space of matrices. First note that the space of matrices $\mathbb{R}^{m \times n}$ is itself a vector space since it is closed with respect to addition and real scalar multiplication with both operations being distributive and commutative and $\mathbb{R}^{m \times n}$ contains the zero matrix. In addition, we can embed $\mathbb{R}^{m \times n}$ in \mathbb{R}^{mn} by stacking one column on top of another to get a long vector of length mn . This process of stacking the columns is denoted by the *vec* operator (column *vec*): given $A \in \mathbb{R}^{m \times n}$,

$$\text{vec}(A) = \begin{pmatrix} A_{.1} \\ A_{.2} \\ \vdots \\ A_{.n} \end{pmatrix} \in \mathbb{R}^{mn}.$$

EXAMPLE 1.2.

$$\text{vec} \begin{bmatrix} 1 & 2 & -3 \\ 0 & -1 & 4 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 2 \\ -1 \\ -3 \\ 4 \end{bmatrix}$$

Using the *vec* operation, we define an inner product on $\mathbb{R}^{m \times n}$ by taking the inner product of these vectors of length mn . Given $A, B \in \mathbb{R}^{m \times n}$ we write this inner product as $\langle A, B \rangle$. It is easy to show that this inner product obeys the formula

$$\langle A, B \rangle = \text{vec}(A)^T \text{vec}(B) = \text{tr}(A^T B).$$

This is known as the *Frobenius inner product*. It generates a corresponding norm, called the *Frobenius norm*, by setting

$$\|A\|_F := \|\text{vec}(A)\|_2 = \sqrt{\langle A, A \rangle}.$$

Note that for a given $x \in \mathbb{R}^n$ and $A \in \mathbb{R}^{m \times n}$ we have

$$\|Ax\|_2^2 = \sum_{i=1}^m (A_i \cdot x)^2 \leq \sum_{i=1}^m (\|A_i\|_2 \|x\|_2)^2 = \|x\|_2^2 \sum_{i=1}^m \|A_i\|_2^2 = \|A\|_F^2 \|x\|_2^2,$$

and so

$$(61) \quad \|Ax\|_2 \leq \|A\|_F \|x\|_2.$$

This relationship between the Frobenius norm and the 2-norm is very important and is used extensively in our development. In particular, this implies that for any two matrices $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{n \times k}$ we have

$$\|AB\|_F \leq \|A\|_F \|B\|_F.$$

2. Differentiation

In this section we use our understanding of differentiability for mappings from \mathbb{R} to \mathbb{R} to build a theory of differentiation for mappings from \mathbb{R}^n to \mathbb{R}^m . Let F be a mapping from \mathbb{R}^n to \mathbb{R}^m which we denote by $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$. Let the component functions of F be denoted by $F_i : \mathbb{R}^n \rightarrow \mathbb{R}$:

$$F(x) = \begin{pmatrix} F_1(x) \\ F_2(x) \\ \vdots \\ F_m(x) \end{pmatrix}.$$

EXAMPLE 2.1.

$$F(x) = F \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 3x_1^2 + x_1x_2x_3 \\ 2 \cos(x_1) \sin(x_2x_3) \\ \ln[\exp(x_1^2 + x_2^2 + x_3^2)] \\ 1/\sqrt{1 + (x_2x_3)^2} \end{pmatrix}.$$

In this case, $n = 3$, $m = 4$, and

$$F_1(x) = 3x_1^2 + x_1x_2x_3, \quad F_2(x) = 2 \cos(x_1) \sin(x_2x_3), \quad F_3(x) = \ln[\exp(x_1^2 + x_2^2 + x_3^2)], \quad F_4(x) = 1/\sqrt{1 + (x_2x_3)^2}.$$

The first step in understanding the differentiability of mappings on \mathbb{R}^n is to study their one dimensional properties. For this, consider a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and let x and d be elements of \mathbb{R}^n . We define the *directional derivative* of f in the direction d , when it exists, to be the one sided limit

$$f'(x; d) := \lim_{t \downarrow 0} \frac{f(x + td) - f(x)}{t}.$$

EXAMPLE 2.2. Let $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ be given by $f(x_1, x_2) := x_1 |x_2|$, and let $x = (1, 0)^T$ and $d = (2, 2)$. Then,

$$f'(x; d) = \lim_{t \downarrow 0} \frac{f(x + td) - f(x)}{t} = \lim_{t \downarrow 0} \frac{(1 + 2t) |0 + 2t| - 1 |0|}{t} = \lim_{t \downarrow 0} \frac{2(1 + 2t)t}{t} = 2,$$

while, for $d = -(2, 2)^T$,

$$f'(x; d) = \lim_{t \downarrow 0} \frac{f(x + td) - f(x)}{t} = \lim_{t \downarrow 0} \frac{(1 - 2t) |0 - 2t| - 1 |0|}{t} = \lim_{t \downarrow 0} \frac{2(1 - 2t)t}{t} = 2.$$

In general, we have

$$f'((1, 0); (d_1, d_2)) = \lim_{t \downarrow 0} \frac{(1 + d_1t) |d_2t|}{t} = |d_2|.$$

For technical reasons, we allow this limit to take the values $\pm\infty$. For example, if $f(x) = x^{1/3}$, then

$$f'(0; 1) = \lim_{t \downarrow 0} t^{-2/3} = +\infty \quad \text{and} \quad f'(0; -1) = \lim_{t \downarrow 0} -t^{-2/3} = -\infty.$$

This example as well as the one given in Example 2.2 show that the directional derivative $f'(x; d)$ is not necessarily either continuous or smooth in the d argument even if it exists for all choices of d . However, the directional derivative is always *positively homogeneous* in the sense that, given $\lambda \geq 0$, we have

$$f'(x; \lambda d) = \lim_{t \downarrow 0} \frac{f(x + \lambda td) - f(x)}{t} = \lambda \lim_{t \downarrow 0} \frac{f(x + td) - f(x)}{\lambda t} = \lambda f'(x; d).$$

The directional derivative idea can be extended to functions F mapping \mathbb{R}^n into \mathbb{R}^m by defining it componentwise: if the limit

$$F'(x; d) := \lim_{t \downarrow 0} \frac{F(x + td) - F(x)}{t} = \begin{pmatrix} \lim_{t \downarrow 0} \frac{F_1(x+td) - F_1(x)}{t} \\ \lim_{t \downarrow 0} \frac{F_2(x+td) - F_2(x)}{t} \\ \vdots \\ \lim_{t \downarrow 0} \frac{F_m(x+td) - F_m(x)}{t} \end{pmatrix}$$

exists, it is called the directional derivative of F at x in the direction d .

These elementary ideas lead to the following notions of differentiability.

DEFINITION 2.1. [Differentiable Functions] Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$.

- (1) If $f'(x; d) = \lim_{t \rightarrow 0} \frac{f(x + \lambda td) - f(x)}{t}$, then we say that f is differentiable in the direction d , in which case $f'(x; -d) = -f'(x; d)$.
- (2) Let e_j $j = 1, \dots, n$ denote the unit coordinate vectors. If f is differentiable in the direction e_j , we say that the partial derivative of f with respect to the component x_j exists and write

$$\frac{\partial f(x)}{\partial x_j} := f'(x; e_j).$$

In particular, we have

$$f(x + te_j) = f(x) + t \frac{\partial f(x)}{\partial x_j} + o(t), \quad \text{where } \lim_{t \rightarrow 0} \frac{o(t)}{t} = 0.$$

Note that $\frac{\partial f(\cdot)}{\partial x_j} : \mathbb{R}^n \rightarrow \mathbb{R}$.

- (3) We say that f is (Fréchet) differentiable at $x \in \mathbb{R}^n$ if there is a vector $g \in \mathbb{R}^n$ such that

$$\lim_{y \rightarrow x} \frac{|f(y) - f(x) - g^T(y - x)|}{\|y - x\|} = 0.$$

If such a vector g exists, we write $g = \nabla f(x)$ and call $\nabla f(x)$ the gradient of f at x . In particular, the differentiability of f at x is equivalent to the following statement:

$$f(y) = f(x) + \nabla f(x)^T(y - x) + o(\|y - x\|)$$

for all y near x , where $\lim_{y \rightarrow x} \frac{o(\|y - x\|)}{\|y - x\|} = 0$.

- (4) We say that F is (Fréchet) differentiable at $x \in \mathbb{R}^n$ if there is a matrix $J \in \mathbb{R}^{m \times n}$ such that

$$\lim_{y \rightarrow x} \frac{\|F(y) - F(x) - J(y - x)\|}{\|y - x\|} = 0.$$

If such a matrix J exists, we write $J = \nabla F(x)$ and call $\nabla F(x)$ the Jacobian of F at x . In particular, the differentiability of F at x is equivalent to the following statement:

$$F(y) = F(x) + \nabla F(x)^T(y - x) + o(\|y - x\|)$$

for all y near x , where $\lim_{y \rightarrow x} \frac{o(\|y - x\|)}{\|y - x\|} = 0$.

REMARK 2.1. Note that there is an inconsistency here in the use of the ∇ notation when $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$ with $m = 1$. The inconsistency arises due to the presence of g^T in Part (3) of Definition 2.1 and the absence of a transpose in Part (4) of this definition. For this reason, we must take extra care in interpreting this notation in this case.

REMARK 2.2. [Little-o Notation] In these notes we use the notation $o(t)$ to represent any element of a function class for which $\lim_{t \rightarrow 0} \frac{o(t)}{t} = 0$. In particular, this implies that for all $\alpha \in \mathbb{R}$

$$\alpha o(t) = o(t), \quad o(t) + o(t) = o(t), \quad \text{and} \quad t^s o(t^r) = o(t^{r+s}).$$

Several observations about these notions of differentiability are in order. First, the existence of the directional derivative $f'(x; d)$ nor the differentiability of f at x in the direction d requires the continuity of the function at that point. Second, the existence of $f'(x; d)$ in all directions d does imply the continuity of the mapping $d \mapsto f'(x; d)$. Therefore, the directional derivative, although useful, is a very weak object to describe the local variational properties of a function. On the other hand, differentiability is a very powerful statement. A few consequences of differentiability are listed in the following theorem.

THEOREM 2.1. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$.*

(1) *If f is differentiable at x , then*

$$\nabla f(x) = \begin{pmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{pmatrix},$$

and $f'(x; d) = \nabla f(x)^T d$ for all $d \in \mathbb{R}^n$.

(2) *If F is differentiable at x , then*

$$(\nabla F(x))_{ij} = \frac{\partial F_i(x)}{\partial x_j} \quad i = 1, \dots, m \quad \text{and} \quad j = 1, 2, \dots, n.$$

(3) *If F is differentiable at a point x , then it is necessarily continuous at x .*

Higher order derivatives are obtained by applying these notions of differentiability to the derivatives themselves. For example, to compute the second derivative, the derivative needs to exist at all points near the point at which the second derivative needs to be computed so that the necessary limit is well defined. From the above, we know that the partial derivative $\frac{\partial F_i(x)}{\partial x_j}$, when it exists, is a mapping from \mathbb{R}^n to \mathbb{R} . Therefore, it is possible to consider the partial derivatives of these partial derivatives. For such partial derivatives we use the notation

$$(62) \quad \frac{\partial^2 F_i(x)}{\partial x_j \partial x_k} := \frac{\partial \left(\frac{\partial F_i(x)}{\partial x_k} \right)}{\partial x_j}.$$

The second derivative of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is the derivative of the mapping $\nabla f : \mathbb{R}^n \rightarrow \mathbb{R}^n$, and we write $\nabla(\nabla f(x)) =: \nabla^2 f(x)$. We call $\nabla^2 f(x)$ the *Hessian* of f at x . By (62), we have

$$\nabla^2 f(x) = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \cdots & \frac{\partial^2 f(x)}{\partial x_n \partial x_1} \\ \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \frac{\partial^2 f(x)}{\partial x_2^2} & \cdots & \frac{\partial^2 f(x)}{\partial x_n \partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_1 \partial x_n} & \frac{\partial^2 f(x)}{\partial x_2 \partial x_n} & \cdots & \frac{\partial^2 f(x)}{\partial x_n^2} \end{bmatrix}.$$

We have the following key property of the Hessian.

THEOREM 2.2. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be such that all of the second partials $\frac{\partial^2 f(x)}{\partial x_i \partial x_j}$, $i, j = 1, 2, \dots, n$ exist and are continuous near $x \in \mathbb{R}^n$. Then $\nabla^2 f(x)$ is a real symmetric matrix, i.e., $\nabla^2 f(x) = \nabla^2 f(x)^T$.*

The partial derivative representations of the gradient, Hessian, and Jacobian matrices is a convenient tool for computing these objects. For example, if we have

$$f(x) := 3x_1^2 + x_1 x_2 x_3,$$

then

$$\nabla f(x) = \begin{pmatrix} 6x_1 + x_2 x_3 \\ x_1 x_3 \\ x_1 x_2 \end{pmatrix} \quad \text{and} \quad \nabla^2 f(x) = \begin{bmatrix} 6 & x_3 & x_2 \\ x_3 & 0 & x_1 \\ x_2 & x_1 & 0 \end{bmatrix}.$$

However, the partial derivatives are not the only tool for computing derivatives. In many cases, it is easier to compute the gradient, Hessian, and/or Jacobian directly from the definition using the little-o notation.

3. The Delta Method for Computing Derivatives

Recall that a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is said to be differentiable at a point x if there is a vector $g \in \mathbb{R}^n$ such that

$$(63) \quad f(x + \Delta x) = f(x) + g^T \Delta x + o(\|\Delta x\|).$$

Hence, if we can write $f(x + \Delta x)$ in this form, then $g = \nabla f(x)$. To see how to use this idea, consider the least squares objective function

$$f(x) = \frac{1}{2} \|Ax - b\|_2^2, \quad \text{where } A \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m.$$

Then

$$(64) \quad \begin{aligned} f(x + \Delta x) &= \frac{1}{2} \|A(x + \Delta x) - b\|_2^2 \\ &= \frac{1}{2} \|(Ax - b) + A\Delta x\|_2^2 \\ &= \frac{1}{2} \|Ax - b\|_2^2 + (Ax - b)^T A\Delta x + \frac{1}{2} \|A\Delta x\|_2^2 \\ &= f(x) + (A^T(Ax - b))^T \Delta x + \frac{1}{2} \|A\Delta x\|_2^2. \end{aligned}$$

In this expression, $\frac{1}{2} \|A\Delta x\|_2^2 = o(\|\Delta x\|_2)$ since

$$\frac{\frac{1}{2} \|A\Delta x\|_2^2}{\|\Delta x\|_2} = \frac{1}{2} \|A\Delta x\|_2 \left\| A \frac{\Delta x}{\|\Delta x\|_2} \right\|_2 \rightarrow 0 \quad \text{as } \|\Delta x\|_2 \rightarrow 0.$$

Therefore, by (63), the expression (64) tells us that

$$\nabla f(x) = A^T(Ax - b).$$

This approach to computing the derivative of a function is called the *delta method*. In a similar manner it can be used to compute the Hessian of f by applying the approach to ∇f :

$$\nabla f(x + \Delta x) = A^T(A(x + \Delta x) - b) = A^T(Ax - b) + A^T A\Delta x = \nabla f(x) + A^T A\Delta x,$$

and, hence, $\nabla^2 f(x) = A^T A$.

Let us now apply the delta method to compute the gradient and Hessian of the quadratic function

$$f(x) := \frac{1}{2} x^T Hx + g^T x, \quad \text{where } H \in \mathcal{S}^n, g \in \mathbb{R}^n.$$

Then

$$\begin{aligned} f(x + \Delta x) &= \frac{1}{2} (x + \Delta x)^T H(x + \Delta x) + g^T (x + \Delta x) \\ &= \frac{1}{2} x^T Hx + g^T x + (Hx + g)^T \Delta x + \frac{1}{2} \Delta x^T H\Delta x \\ &= f(x) + (Hx + g)^T \Delta x + \frac{1}{2} \Delta x^T H\Delta x, \end{aligned}$$

where $\frac{1}{2} \Delta x^T H\Delta x = o(\|\Delta x\|_2)$ since

$$\frac{\frac{1}{2} \Delta x^T H\Delta x}{\|\Delta x\|_2} = \frac{1}{2} \Delta x^T H \frac{\Delta x}{\|\Delta x\|_2} \rightarrow 0.$$

Therefore, by (63), we must have

$$\nabla f(x) = Hx + g.$$

Again, we compute the Hessian by applying the delta method to the gradient:

$$\nabla f(x + \Delta x) = H(x + \Delta x) + g = (Hx + g) + H\Delta x = \nabla f(x) + H\Delta x,$$

and so

$$\nabla^2 f(x) = H.$$

4. Differential Calculus

There are many further tools for computing derivatives that do not require a direct appeal to either the partial derivatives or the delta method. These tools allow us to compute new derivatives from derivatives that are already known based on a *calculus* of differentiation. We are familiar with this differential calculus for functions mapping \mathbb{R} to \mathbb{R} . Here we show how a few of these calculus rules extend to mappings from \mathbb{R}^n to \mathbb{R}^m . The most elementary of these are the facts that the derivative of the scalar multiple of a function equals the scalar multiple of the derivative and the derivative of a sum is the sum of derivatives: given $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$, $G : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $\alpha \in \mathbb{R}$,

$$\nabla(\alpha F) = \alpha \nabla F \quad \text{and} \quad \nabla(F + G) = \nabla F + \nabla G .$$

These rules are themselves derivable from the much more powerful *chain rule*.

THEOREM 4.1. *Let $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $H : \mathbb{R}^m \rightarrow \mathbb{R}^k$ be such that F is differentiable at x and H is differentiable at $F(x)$. Then $G := H \circ F$ is differentiable at x with*

$$\nabla G(x) = \nabla H(F(x)) \circ \nabla F(x) .$$

REMARK 4.1. *As noted in Remark 2.1, one must take special care in the interpretation of this chain rule when $k = 1$ due to the presence of an additional transpose. In this case,*

$$\nabla G(x) = \nabla F(x)^T \nabla H(F(x)) .$$

For example, let $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and consider the function

$$f(x) := \frac{1}{2} \|F(x)\|_2^2 = \left(\frac{1}{2} \|\cdot\|_2^2\right) \circ F(x),$$

that is, we are composing half the 2-norm squared with F . Since $\nabla(\frac{1}{2} \|\cdot\|_2^2)(y) = y$, we have

$$\nabla f(x) = \nabla F(x)^T F(x) .$$

This chain rule computation can be verified using the delta method:

$$\begin{aligned} f(x + \Delta x) &= \frac{1}{2} \|F(x + \Delta x)\|_2^2 \\ &= \frac{1}{2} \|F(x) + \nabla F(x)\Delta x + o(\|\Delta x\|_2)\|_2^2 \\ &= \frac{1}{2} \|F(x) + \nabla F(x)\Delta x\|_2^2 + (F(x) + \nabla F(x)\Delta x)^T (o(\|\Delta x\|_2)) + \frac{1}{2} \|o(\|\Delta x\|_2)\|_2^2 \\ &= \frac{1}{2} \|F(x) + \nabla F(x)\Delta x\|_2^2 + o(\|\Delta x\|_2) \\ &= \frac{1}{2} \|F(x)\|_2^2 + (\nabla F(x)^T F(x))^T \Delta x + \frac{1}{2} \|\nabla F(x)\Delta x\|_2^2 + o(\|\Delta x\|_2) \\ &= f(x) + (\nabla F(x)^T F(x))^T \Delta x + o(\|\Delta x\|_2), \end{aligned}$$

where $\lim_{t \rightarrow 0} \frac{o(t)}{t} = 0$ and we have used this notation as described in Remark 4.1. Hence, again $\nabla f(x) = \nabla F(x)^T F(x)$.

5. The Mean Value Theorem

Given $f : \mathbb{R}^n \rightarrow \mathbb{R}$, the defining formula for the derivative,

$$f(y) = f(x) + \nabla f(x)(y - x) + o(\|y - x\|),$$

is a powerful tool for understanding the local behavior of the function f near x . If we drop the little- o term from the right hand side, we obtain the *first-order Taylor expansion* of f at x . This is called a *first-order approximation* to f at x due to the fact that the power of $\|y - x\|$ in the *error term* $o(\|y - x\|)$ is 1. Higher order approximations to f can be obtained using higher order derivatives. But before turning to these approximations, we make a closer study of the first-order expansion. In particular, we wish to extend the *Mean Value Theorem* to functions of many variables.

THEOREM 5.1. [*1-Dimensional Mean Value Theorem*]

Let $\phi : \mathbb{R} \rightarrow \mathbb{R}$ be $k + 1$ times differentiable on an open interval $(a, b) \subset \mathbb{R}$. Then, for every $x, y \in (a, b)$ with $x \neq y$, there exists a $z \in (a, b)$ strictly between x and y such that

$$\phi(y) = \phi(x) + \phi'(x)(y - x) + \cdots + \frac{1}{k!} \phi^{(k)}(x)(y - x)^k + \frac{1}{(k + 1)!} \phi^{(k+1)}(z)(y - x)^{(k+1)} .$$

REMARK 5.1. *Theorem 5.1 is also called Taylor's Theorem with remainder, where we have chosen the Lagrange form for the remainder. Other forms for the remainder include the integral form and Cauchy's form. In most texts the name Mean Value Theorem is reserved for the first-order case alone.*

We use this results to easily obtain the following mean value theorem for function mapping \mathbb{R}^n to \mathbb{R} .

THEOREM 5.2. [*n-Dimensional Mean Value Theorem*]

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be differentiable on an open set containing the two points $x, y \in \mathbb{R}^n$ with $x \neq y$. Define the closed and open line segments connecting x and y by

$$[x, y] := \{(1 - \lambda)x + \lambda y \mid 0 \leq \lambda \leq 1\} \quad \text{and} \quad (x, y) := \{(1 - \lambda)x + \lambda y \mid 0 < \lambda < 1\},$$

respectively. Then there exists a $z, w \in (x, y)$ such that

$$f(y) = f(x) + \nabla f(z)^T(y - x) \quad \text{and} \quad f(y) = f(x) + \nabla f(x)^T(y - x) + \frac{1}{2}(y - x)^T \nabla^2 f(z)(y - x).$$

PROOF. Define the function $\phi : \mathbb{R} \rightarrow \mathbb{R}$ by $\phi(t) := f(x + t(y - x))$. Since f is differentiable, so is ϕ and the chain rule tells us that

$$\phi'(t) = \nabla f(x + t(y - x))^T(y - x) \quad \text{and} \quad \phi'(t) = (y - x)^T \nabla^2 f(x + t(y - x))(y - x).$$

By applying the Mean Value Theorem 5.1 to ϕ we obtain the existence of $t, s \in (0, 1)$ such that

$$f(y) = \phi(1) = \phi(0) + \phi'(t)(1 - 0) = f(x) + \nabla f(x + t(y - x))^T(y - x)$$

and

$$f(y) = \phi(1) = \phi(0) + \phi'(0)(1 - 0) + \frac{1}{2}\phi''(s)(1 - 0)^2 = f(x) + \nabla f(x)^T(y - x) + \frac{1}{2}(y - x)^T \nabla^2 f(x + s(y - x))(y - x).$$

By setting $z := x + t(y - x)$ and $w := x + s(y - x)$ we obtain the result. \square

In a similar manner we can apply the Fundamental Theorem of Calculus to such functions.

THEOREM 5.3. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be differentiable on an open set containing the two points $x, y \in \mathbb{R}^n$ with $x \neq y$. Then

$$f(y) = f(x) + \int_0^1 \nabla f(x + t(y - x))^T(y - x) dt .$$

PROOF. Apply the Fundamental Theorem of Calculus to the function ϕ defined in the proof of Theorem 5.2. \square

Unfortunately, the Mean Value Theorem does not extend to general differentiable function mapping from \mathbb{R}^n to \mathbb{R}^m for $m > 1$. Nonetheless, we have the following approximate result.

THEOREM 5.4. Let $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be differentiable on an open set containing the two points $x, y \in \mathbb{R}^n$ with $x \neq y$. Then

$$(65) \quad \|F(y) - F(x)\|_2 \leq \left[\max_{z \in [x, y]} \|F'(z)\|_F \right] \|y - x\|_2 .$$

PROOF. By the Fundamental Theorem of Calculus, we have

$$F(y) - F(x) = \begin{pmatrix} \int_0^1 \nabla F_1(x + t(y - x))^T(y - x) dt \\ \vdots \\ \int_0^1 \nabla F_m(x + t(y - x))^T(y - x) dt \end{pmatrix} = \int_0^1 \nabla F(x + t(y - x))(y - x) dt .$$

Therefore,

$$\begin{aligned} \|F(y) - F(x)\|_2 &= \left\| \int_0^1 \nabla F(x + t(y - x))(y - x) dt \right\|_2 \\ &\leq \int_0^1 \|\nabla F(x + t(y - x))(y - x)\|_2 dt \\ &\leq \int_0^1 \|\nabla F(x + t(y - x))\|_F \|y - x\|_2 dt \\ &\leq \left[\max_{z \in [x, y]} \|F'(z)\|_F \right] \|y - x\|_2 . \end{aligned}$$

\square

The bound (65) is very useful in many applications. But it can be simplified in cases where ∇F is known to be continuous since in this case the Weierstrass extreme value theorem says that, for every $\beta > 0$,

$$\max_{z \in \beta \mathbb{B}} \|F'(z)\|_F =: K < \infty .$$

Hence, by Theorem 5.4,

$$\|F(x) - F(y)\|_2 \leq K \|x - y\|_2 \quad \forall x, y \in \beta \mathbb{B} .$$

This kind of inequality is extremely useful and leads to the following notion of continuity.

DEFINITION 5.1. [*Lipschitz Continuity*]

We say that $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is Lipschitz continuous on a set $S \subset \mathbb{R}^n$ if there exists a constant $K > 0$ such that

$$\|F(x) - F(y)\| \leq K \|x - y\| \quad \forall x, y \in S .$$

The constant K is called the modulus of Lipschitz continuity for F over S , and depends on the choice of norms for \mathbb{R}^n and \mathbb{R}^m .

As one application of Lipschitz continuity, we give the following lemma concerning the accuracy of the first-order Taylor approximation of a function.

LEMMA 5.1. [*Quadratic Bound Lemma*]

Let $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be such that ∇F is Lipschitz continuous on the set $S \subset \mathbb{R}^n$. If $x, y \in S$ are such that $[x, y] \subset S$, then

$$\|F(y) - (F(x) + \nabla F(x)(y - x))\|_2 \leq \frac{K}{2} \|y - x\|_2^2 ,$$

where K is the modulus of Lipschitz continuity for ∇F on S .

PROOF. Observe that

$$\begin{aligned} F(y) - F(x) - \nabla F(x)(y - x) &= \int_0^1 \nabla F(x + t(y - x))(y - x) dt - \nabla F(x)(y - x) \\ &= \int_0^1 [\nabla F(x + t(y - x)) - \nabla F(x)](y - x) dt . \end{aligned}$$

Hence

$$\begin{aligned} \|F(y) - (F(x) + \nabla F(x)(y - x))\|_2 &= \left\| \int_0^1 [\nabla F(x + t(y - x)) - \nabla F(x)](y - x) dt \right\|_2 \\ &\leq \int_0^1 \|(\nabla F(x + t(y - x)) - \nabla F(x))(y - x)\|_2 dt \\ &\leq \int_0^1 \|\nabla F(x + t(y - x)) - \nabla F(x)\|_F \|y - x\|_2 dt \\ &\leq \int_0^1 K t \|y - x\|_2^2 dt \\ &= \frac{K}{2} \|y - x\|_2^2 . \end{aligned}$$

□

The Mean Value Theorem also allows to obtain the following *second order* approximation.

THEOREM 5.5. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and suppose that $\nabla^2 f(x)$ exists and is continuous at x . Then

$$(66) \quad f(y) = f(x) + \nabla f(x)^T (y - x) + \frac{1}{2} (y - x)^T \nabla^2 f(x) (y - x) + o(\|y - x\|^2) .$$

PROOF. The mean value theorem tells us that for every $y \in x + \epsilon \mathbb{B}$ there is a $z \in (x, y)$ such that

$$\begin{aligned} f(y) &= f(x) + \nabla f(x)^T (y - x) + \frac{1}{2} (y - x)^T \nabla^2 f(z) (y - x) \\ &= f(x) + \nabla f(x)^T (y - x) + \frac{1}{2} (y - x)^T \nabla^2 f(x) (y - x) + \frac{1}{2} (y - x)^T [\nabla^2 f(z) - \nabla^2 f(x)] (y - x) \\ &= f(x) + \nabla f(x)^T (y - x) + \frac{1}{2} (y - x)^T \nabla^2 f(x) (y - x) + o(\|y - x\|^2) . \end{aligned}$$

□

If we drop the $o(\|y - x\|^2)$ in the equation (66), we obtain the *second-order Taylor approximation* to f at x . This is a second-order approximation since the power of $\|y - x\|$ in the little-o term is 2, i.e., $o(\|y - x\|^2)$.

Optimality Conditions for Unconstrained Problems

1. Existence of Optimal Solutions

Consider the problem of minimizing the function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ where f is continuous on all of \mathbb{R}^n :

$$\mathcal{P} \quad \min_{x \in \mathbb{R}^n} f(x).$$

As we have seen, there is no guarantee that f has a minimum value, or if it does, it may not be attained. To clarify this situation, we examine conditions under which a solution is guaranteed to exist. Recall that we already have at our disposal a rudimentary existence result for constrained problems. This is the Weierstrass Extreme Value Theorem.

THEOREM 1.1. (WEIERSTRASS EXTREME VALUE THEOREM) *Every continuous function on a compact set attains its extreme values on that set.*

We now build a basic existence result for unconstrained problems based on this theorem. For this we make use of the notion of a coercive function.

DEFINITION 1.1. *A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is said to be coercive if for every sequence $\{x^\nu\} \subset \mathbb{R}^n$ for which $\|x^\nu\| \rightarrow \infty$ it must be the case that $f(x^\nu) \rightarrow +\infty$ as well.*

Continuous coercive functions can be characterized by an underlying compactness property on their lower level sets.

THEOREM 1.2. (Coercivity and Compactness) *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be continuous on all of \mathbb{R}^n . The function f is coercive if and only if for every $\alpha \in \mathbb{R}$ the set $\{x \mid f(x) \leq \alpha\}$ is compact.*

PROOF. We first show that the coercivity of f implies the compactness of the sets $\{x \mid f(x) \leq \alpha\}$. We begin by noting that the continuity of f implies the closedness of the sets $\{x \mid f(x) \leq \alpha\}$. Thus, it remains only to show that any set of the form $\{x \mid f(x) \leq \alpha\}$ is bounded. We show this by contradiction. Suppose to the contrary that there is an $\alpha \in \mathbb{R}$ such that the set $S = \{x \mid f(x) \leq \alpha\}$ is unbounded. Then there must exist a sequence $\{x^\nu\} \subset S$ with $\|x^\nu\| \rightarrow \infty$. But then, by the coercivity of f , we must also have $f(x^\nu) \rightarrow \infty$. This contradicts the fact that $f(x^\nu) \leq \alpha$ for all $\nu = 1, 2, \dots$. Therefore the set S must be bounded.

Let us now assume that each of the sets $\{x \mid f(x) \leq \alpha\}$ is bounded and let $\{x^\nu\} \subset \mathbb{R}^n$ be such that $\|x^\nu\| \rightarrow \infty$. Let us suppose that there exists a subsequence of the integers $J \subset \mathbb{N}$ such that the set $\{f(x^\nu)\}_J$ is bounded above. Then there exists $\alpha \in \mathbb{R}$ such that $\{x^\nu\}_J \subset \{x \mid f(x) \leq \alpha\}$. But this cannot be the case since each of the sets $\{x \mid f(x) \leq \alpha\}$ is bounded while every subsequence of the sequence $\{x^\nu\}$ is unbounded by definition. Therefore, the set $\{f(x^\nu)\}_J$ cannot be bounded, and so the sequence $\{f(x^\nu)\}$ contains no bounded subsequence, i.e. $f(x^\nu) \rightarrow \infty$. \square

This result in conjunction with Weierstrass's Theorem immediately yields the following existence result for the problem \mathcal{P} .

THEOREM 1.3. (Coercivity implies existence) *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be continuous on all of \mathbb{R}^n . If f is coercive, then f has at least one global minimizer.*

PROOF. Let $\alpha \in \mathbb{R}$ be chosen so that the set $S = \{x \mid f(x) \leq \alpha\}$ is non-empty. By coercivity, this set is compact. By Weierstrass's Theorem, the problem $\min \{f(x) \mid x \in S\}$ has at least one global solution. Obviously, the set of global solutions to the problem $\min \{f(x) \mid x \in S\}$ is a global solution to \mathcal{P} which proves the result. \square

REMARK 1.1. *It should be noted that the coercivity hypothesis is stronger than is strictly required in order to establish the existence of a solution. Indeed, a global minimizer must exist if there exist one non-empty compact lower level set. We do not need all of them to be compact. However, in practice, coercivity is easy to check.*

2. First-Order Optimality Conditions

This existence result can be quite useful, but unfortunately it does not give us a constructive test for optimality. That is, we may know a solution exists, but we still do not have a method for determining whether any given point may or may not be a solution. We now present such a test using the derivatives of the objective function f . For this we will assume that f is twice continuously differentiable on \mathbb{R}^n and develop constructible first- and second-order necessary and sufficient conditions for optimality.

The optimality conditions we consider are built up from those developed in first term calculus for functions mapping from \mathbb{R} to \mathbb{R} . The reduction to the one dimensional case comes about by considering the functions $\phi : \mathbb{R} \rightarrow \mathbb{R}$ given by

$$\phi(t) = f(x + td)$$

for some choice of x and d in \mathbb{R}^n . The key variational object in this context is the directional derivative of f at a point x in the direction d given by

$$f'(x; d) = \lim_{t \downarrow 0} \frac{f(x + td) - f(x)}{t}.$$

When f is differentiable at the point $x \in \mathbb{R}^n$, then

$$f'(x; d) = \nabla f(x)^T d = \phi'(0).$$

Note that if $f'(x; d) < 0$, then there must be a $\bar{t} > 0$ such that

$$\frac{f(x + td) - f(x)}{t} < 0 \quad \text{whenever} \quad 0 < t < \bar{t}.$$

In this case, we must have

$$f(x + td) < f(x) \quad \text{whenever} \quad 0 < t < \bar{t}.$$

That is, we can always reduce the function value at x by moving in the direction d an arbitrarily small amount. In particular, if there is a direction d such that $f'(x; d)$ exists with $f'(x; d) < 0$, then x cannot be a local solution to the problem $\min_{x \in \mathbb{R}^n} f(x)$. Or equivalently, if x is a local to the problem $\min_{x \in \mathbb{R}^n} f(x)$, then $f'(x; d) \geq 0$ whenever $f'(x; d)$ exists. We state this elementary result in the following lemma.

LEMMA 2.1 (Basic First-Order Optimality Result). *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and let $\bar{x} \in \mathbb{R}^n$ be a local solution to the problem $\min_{x \in \mathbb{R}^n} f(x)$. Then*

$$f'(x; d) \geq 0$$

for every direction $d \in \mathbb{R}^n$ for which $f'(x; d)$ exists.

We now apply this result to the case in which $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is differentiable.

THEOREM 2.1. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be differentiable at a point $\bar{x} \in \mathbb{R}^n$. If \bar{x} is a local minimum of f , then $\nabla f(\bar{x}) = 0$.*

PROOF. By Lemma 2.1 we have

$$0 \leq f'(\bar{x}; d) = \nabla f(\bar{x})^T d \quad \text{for all} \quad d \in \mathbb{R}^n.$$

Taking $d = -\nabla f(\bar{x})$ we get

$$0 \leq -\nabla f(\bar{x})^T \nabla f(\bar{x}) = -\|\nabla f(\bar{x})\|^2 \leq 0.$$

Therefore, $\nabla f(\bar{x}) = 0$. □

When $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is differentiable, any point $x \in \mathbb{R}^n$ satisfying $\nabla f(x) = 0$ is said to be a stationary (or, equivalently, a critical) point of f . In our next result we link the notions of coercivity and stationarity.

THEOREM 2.2. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be differentiable on all of \mathbb{R}^n . If f is coercive, then f has at least one global minimizer these global minimizers can be found from among the set of critical points of f .*

PROOF. Since differentiability implies continuity, we already know that f has at least one global minimizer. Differentiability implies that this global minimizer is critical. □

This result indicates that one way to find a global minimizer of a coercive differentiable function is to first find all critical points and then from among these determine those yielding the smallest function value.

3. Second-Order Optimality Conditions

To obtain second-order conditions for optimality we must first recall a few properties of the Hessian matrix $\nabla^2 f(x)$. The calculus tells us that if f is twice continuously differentiable at a point $x \in \mathbb{R}^n$, then the hessian $\nabla^2 f(x)$ is a symmetric matrix. Symmetric matrices are orthogonally diagonalizable. That is, there exists an orthonormal basis of eigenvectors of $\nabla^2 f(x)$, $v^1, v^2, \dots, v^n \in \mathbb{R}^n$ such that

$$\nabla^2 f(x) = V \begin{bmatrix} \lambda_1 & 0 & 0 & \dots & 0 \\ 0 & \lambda_2 & 0 & \dots & 0 \\ \vdots & & \ddots & & \vdots \\ 0 & 0 & \dots & \dots & \lambda_n \end{bmatrix} V^T$$

where $\lambda_1, \lambda_2, \dots, \lambda_n$ are the eigenvalues of $\nabla^2 f(x)$ and V is the matrix whose columns are given by their corresponding vectors v^1, v^2, \dots, v^n :

$$V = [v^1, v^2, \dots, v^n] .$$

It can be shown that $\nabla^2 f(x)$ is positive semi-definite if and only if $\lambda_i \geq 0$, $i = 1, 2, \dots, n$, and it is positive definite if and only if $\lambda_i > 0$, $i = 1, 2, \dots, n$. Thus, in particular, if $\nabla^2 f(x)$ is positive definite, then

$$d^T \nabla^2 f(x) d \geq \lambda_{\min} \|d\|^2 \quad \text{for all } d \in \mathbb{R}^n,$$

where λ_{\min} is the smallest eigenvalue of $\nabla^2 f(x)$.

We now give our main result on second-order necessary and sufficient conditions for optimality in the problem $\min_{x \in \mathbb{R}^n} f(x)$. The key tools in the proof are the notions of positive semi-definiteness and definiteness along with the second-order Taylor series expansion for f at a given point $\bar{x} \in \mathbb{R}^n$:

$$(67) \quad f(x) = f(\bar{x}) + \nabla f(\bar{x})^T (x - \bar{x}) + \frac{1}{2} (x - \bar{x})^T \nabla^2 f(\bar{x}) (x - \bar{x}) + o(\|x - \bar{x}\|^2)$$

where

$$\lim_{x \rightarrow \bar{x}} \frac{o(\|x - \bar{x}\|^2)}{\|x - \bar{x}\|^2} = 0.$$

THEOREM 3.1. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be twice continuously differentiable at the point $\bar{x} \in \mathbb{R}^n$.*

- (1) *(Necessity) If \bar{x} is a local minimum of f , then $\nabla f(\bar{x}) = 0$ and $\nabla^2 f(\bar{x})$ is positive semi-definite.*
- (2) *(Sufficiency) If $\nabla f(\bar{x}) = 0$ and $\nabla^2 f(\bar{x})$ is positive definite, then there is an $\alpha > 0$ such that $f(x) \geq f(\bar{x}) + \alpha \|x - \bar{x}\|^2$ for all x near \bar{x} .*

PROOF. (1) We make use of the second-order Taylor series expansion (67) and the fact that $\nabla f(\bar{x}) = 0$ by Theorem 2.1. Given $d \in \mathbb{R}^n$ and $t > 0$ set $x := \bar{x} + td$, plugging this into (67) we find that

$$0 \leq \frac{f(\bar{x} + td) - f(\bar{x})}{t^2} = \frac{1}{2} d^T \nabla^2 f(\bar{x}) d + \frac{o(t^2)}{t^2}$$

since $\nabla f(\bar{x}) = 0$ by Theorem 2.1. Taking the limit as $t \rightarrow 0$ we get that

$$0 \leq d^T \nabla^2 f(\bar{x}) d.$$

Since d was chosen arbitrarily, $\nabla^2 f(\bar{x})$ is positive semi-definite.

- (2) The Taylor expansion (67) and the hypothesis that $\nabla f(\bar{x}) = 0$ imply that

$$(68) \quad \frac{f(x) - f(\bar{x})}{\|x - \bar{x}\|^2} = \frac{1}{2} \frac{(x - \bar{x})^T}{\|x - \bar{x}\|} \nabla^2 f(\bar{x}) \frac{(x - \bar{x})}{\|x - \bar{x}\|} + \frac{o(\|x - \bar{x}\|^2)}{\|x - \bar{x}\|^2}.$$

If $\lambda_{\min} > 0$ is the smallest eigenvalue of $\nabla^2 f(\bar{x})$, choose $\epsilon > 0$ so that

$$(69) \quad \left| \frac{o(\|x - \bar{x}\|^2)}{\|x - \bar{x}\|^2} \right| \leq \frac{\lambda_{\min}}{4}$$

whenever $\|x - \bar{x}\| < \epsilon$. Then, for all $\|x - \bar{x}\| < \epsilon$, we have from (68) and (69) that

$$\begin{aligned} \frac{f(x) - f(\bar{x})}{\|x - \bar{x}\|^2} &\geq \frac{1}{2} \lambda_{\min} + \frac{o(\|x - \bar{x}\|^2)}{\|x - \bar{x}\|^2} \\ &\geq \frac{1}{4} \lambda_{\min}. \end{aligned}$$

Consequently, if we set $\alpha = \frac{1}{4}\lambda_{\min}$, then

$$f(x) \geq f(\bar{x}) + \alpha\|x - \bar{x}\|^2$$

whenever $\|x - \bar{x}\| < \epsilon$.

□

In order to apply the second-order sufficient condition one must be able to check that a symmetric matrix is positive definite. As we have seen, this can be done by computing the eigenvalues of the matrix and checking that they are all positive. But there is another approach that is often easier to implement using the *principal minors* of the matrix.

THEOREM 3.2. *Let $H \in \mathbb{R}^{n \times n}$ be symmetric. We define the k th principal minor of H , denoted $\Delta_k(H)$, to be the determinant of the upper-left $k \times k$ submatrix of H . Then*

- (1) H is positive definite if and only if $\Delta_k(H) > 0$, $k = 1, 2, \dots, n$.
- (2) H is negative definite if and only if $(-1)^k \Delta_k(H) > 0$, $k = 1, 2, \dots, n$.

DEFINITION 3.1. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be continuously differentiable at \bar{x} . If $\nabla f(\bar{x}) = 0$, but \bar{x} is neither a local maximum or a local minimum, we call \bar{x} a saddle point for f .*

THEOREM 3.3. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be twice continuously differentiable at \bar{x} . If $\nabla f(\bar{x}) = 0$ and $\nabla^2 f(\bar{x})$ has both positive and negative eigenvalues, then \bar{x} is a saddle point of f .*

THEOREM 3.4. *Let $H \in \mathbb{R}^{n \times n}$ be symmetric. If H is neither positive definite or negative definite and all of its principal minors are non-zero, then H has both positive and negative eigenvalues. In this case we say that H is indefinite.*

EXAMPLE 3.1. *Consider the matrix*

$$H = \begin{bmatrix} 1 & 1 & -1 \\ 1 & 5 & 1 \\ -1 & 1 & 4 \end{bmatrix}.$$

We have

$$\Delta_1(H) = 1, \quad \Delta_2(H) = \begin{vmatrix} 1 & 1 \\ 1 & 5 \end{vmatrix} = 4, \quad \text{and} \quad \Delta_3(H) = \det(H) = 8.$$

Therefore, H is positive definite.

4. Convexity

In the previous section we established first- and second-order optimality conditions. These conditions are based on only local information and so only refer to properties of local extrema. In this section we study the notion of convexity which allows us to provide optimality conditions for global solutions.

DEFINITION 4.1. (1) *A set $C \subset \mathbb{R}^n$ is said to be convex if for every $x, y \in C$ and $\lambda \in [0, 1]$ one has*

$$(1 - \lambda)x + \lambda y \in C.$$

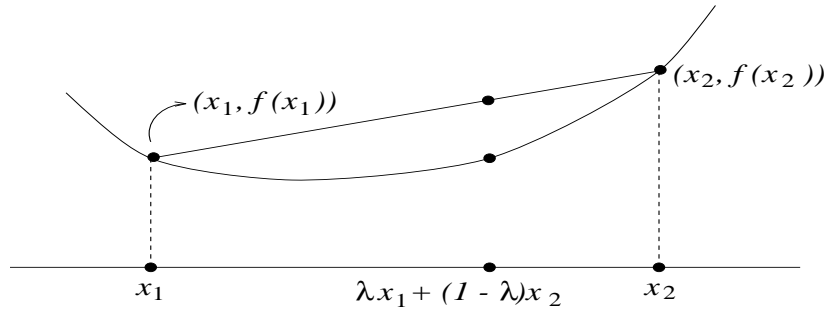
(2) *A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is said to be convex if for every two points $x_1, x_2 \in \mathbb{R}^n$ and $\lambda \in [0, 1]$ we have*

$$(70) \quad f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2).$$

The function f is said to be strictly convex if for every two distinct points $x_1, x_2 \in \mathbb{R}^n$ and $\lambda \in (0, 1)$ we have

$$(71) \quad f(\lambda x_1 + (1 - \lambda)x_2) < \lambda f(x_1) + (1 - \lambda)f(x_2).$$

The inequality (70) is equivalent to the statement that the secant line connecting $(x_1, f(x_1))$ and $(x_2, f(x_2))$ lies above the graph of f on the line segment $\lambda x_1 + (1 - \lambda)x_2$, $\lambda \in [0, 1]$.



That is, the set

$$\text{epi}(f) = \{(x, \mu) : f(x) \leq \mu\},$$

called the *epi-graph* of f is a convex set. Indeed, it can be shown that the convexity of the set $\text{epi}(f)$ is equivalent to the convexity of the function f . This observation allows us to extend the definition of the convexity of a function to functions taking potentially infinite values.

DEFINITION 4.2. A function $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\} = \bar{\mathbb{R}}$ is said to be convex if the set $\text{epi}(f) = \{(x, \mu) : f(x) \leq \mu\}$ is a convex set. We also define the essential domain of f to be the set

$$\text{dom}(f) = \{x : f(x) < +\infty\}.$$

We say that f is strictly convex if the strict inequality (71) holds whenever $x_1, x_2 \in \text{dom}(f)$ are distinct.

EXAMPLE 4.1. $c^T x$, $\|x\|$, e^x , x^2

The role of convexity in linking the global and the local in optimization theory is illustrated by the following result.

THEOREM 4.1. Let $f : \mathbb{R}^n \rightarrow \bar{\mathbb{R}}$ be convex. If $\bar{x} \in \mathbb{R}^n$ is a local minimum for f , then \bar{x} is a global minimum for f .

PROOF. Suppose to the contrary that there is a $\hat{x} \in \mathbb{R}^n$ with $f(\hat{x}) < f(\bar{x})$. Since \bar{x} is a local solution, there is an $\epsilon > 0$ such that

$$f(\bar{x}) \leq f(x) \quad \text{whenever} \quad \|x - \bar{x}\| \leq \epsilon.$$

Taking ϵ smaller if necessary, we may assume that

$$\epsilon < 2\|\bar{x} - \hat{x}\|.$$

Set $\lambda := \epsilon(2\|\bar{x} - \hat{x}\|)^{-1} < 1$ and $x_\lambda := \bar{x} + \lambda(\hat{x} - \bar{x})$. Then $\|x_\lambda - \bar{x}\| \leq \epsilon/2$ and $f(x_\lambda) \leq (1 - \lambda)f(\bar{x}) + \lambda f(\hat{x}) < f(\bar{x})$. This contradicts the choice of ϵ and so no such \hat{x} exists. \square

Strict convexity implies the uniqueness of solutions.

THEOREM 4.2. *Let $f : \mathbb{R}^n \rightarrow \bar{\mathbb{R}}$ be strictly convex. If f has a global minimizer, then it is unique.*

PROOF. Let x^1 and x^2 be distinct global minimizers of f . Then, for $\lambda \in (0, 1)$,

$$f((1 - \lambda)x^1 + \lambda x^2) < (1 - \lambda)f(x^1) + \lambda f(x^2) = f(x^1),$$

which contradicts the assumption that x^1 is a global minimizer. \square

If f is a differentiable convex function, much more can be said. We begin with the following lemma.

LEMMA 4.1. *Let $f : \mathbb{R}^n \rightarrow \bar{\mathbb{R}}$ be convex (not necessarily differentiable).*

(1) *Given $x, d \in \mathbb{R}^n$ the difference quotient*

$$(72) \quad \frac{f(x + td) - f(x)}{t}$$

is a non-decreasing function of t on $(0, +\infty)$.

(2) *For every $x, d \in \mathbb{R}^n$ the directional derivative $f'(x; d)$ always exists and is given by*

$$(73) \quad f'(x; d) := \inf_{t > 0} \frac{f(x + td) - f(x)}{t}.$$

PROOF. We first assume (1) is true and show (2). Recall that

$$(74) \quad f'(x; d) := \lim_{t \downarrow 0} \frac{f(x + td) - f(x)}{t}.$$

Now if the difference quotient (72) is non-decreasing in t on $(0, +\infty)$, then the limit in (74) is necessarily given by the infimum in (73). This infimum always exists and so $f'(x; d)$ always exists and is given by (73).

We now prove (1). Let $x, d \in \mathbb{R}^n$ and let $0 < t_1 < t_2$. Then

$$\begin{aligned} f(x + t_1 d) &= f\left(x + \left(\frac{t_1}{t_2}\right) t_2 d\right) \\ &= f\left[\left(1 - \left(\frac{t_1}{t_2}\right)\right) x + \left(\frac{t_1}{t_2}\right) (x + t_2 d)\right] \\ &\leq \left(1 - \frac{t_1}{t_2}\right) f(x) + \left(\frac{t_1}{t_2}\right) f(x + t_2 d). \end{aligned}$$

Hence

$$\frac{f(x + t_1 d) - f(x)}{t_1} \leq \frac{f(x + t_2 d) - f(x)}{t_2}.$$

\square

A very important consequence of Lemma 4.1 is the *subdifferential inequality*. This inequality is obtained by plugging $t = 1$ and $d = y - x$ into the right hand side of (73) where y is any other point in \mathbb{R}^n . This substitution gives the inequality

$$(75) \quad f(y) \geq f(x) + f'(x; y - x) \quad \text{for all } y \in \mathbb{R}^n \text{ and } x \in \text{dom } f.$$

The subdifferential inequality immediately yields the following result.

THEOREM 4.3 (Convexity and Optimality). *Let $f : \mathbb{R}^n \rightarrow \bar{\mathbb{R}}$ be convex (not necessarily differentiable) and let $\bar{x} \in \text{dom } f$. Then the following three statements are equivalent.*

(i) \bar{x} is a local solution to $\min_{x \in \mathbb{R}^n} f(x)$.

(ii) $f'(\bar{x}; d) \geq 0$ for all $d \in \mathbb{R}^n$.

(iii) \bar{x} is a global solution to $\min_{x \in \mathbb{R}^n} f(x)$.

PROOF. Lemma 2.1 gives the implication (i) \Rightarrow (ii). To see the implication (ii) \Rightarrow (iii) we use the subdifferential inequality and the fact that $f'(\bar{x}; y - \bar{x})$ exists for all $y \in \mathbb{R}^n$ to obtain

$$f(y) \geq f(\bar{x}) + f'(\bar{x}; y - \bar{x}) \geq f(\bar{x}) \quad \text{for all } y \in \mathbb{R}^n.$$

The implication (iii) \Rightarrow (i) is obvious. \square

If it is further assumed that f is differentiable, then we obtain the following elementary consequence of Theorem 4.3.

THEOREM 4.4. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be convex and suppose that $\bar{x} \in \mathbb{R}^n$ is a point at which f is differentiable. Then \bar{x} is a global minimum of f if and only if $\nabla f(\bar{x}) = 0$.*

As Theorems 4.3 and 4.4 demonstrate, convex functions are well suited to optimization theory. Thus, it is important that we are able to recognize when a function is convex. For this reason we give the following result.

THEOREM 4.5. *Let $f : \mathbb{R}^n \rightarrow \bar{\mathbb{R}}$.*

(1) *If f is differentiable on \mathbb{R}^n , then the following statements are equivalent:*

(a) *f is convex,*

(b) *$f(y) \geq f(x) + \nabla f(x)^T(y - x)$ for all $x, y \in \mathbb{R}^n$*

(c) *$(\nabla f(x) - \nabla f(y))^T(x - y) \geq 0$ for all $x, y \in \mathbb{R}^n$.*

(2) *If f is twice differentiable then f is convex if and only if $\nabla^2 f(x)$ is positive semi-definite for all $x \in \mathbb{R}^n$.*

REMARK 4.1. *The condition in Part (c) is called monotonicity.*

PROOF. (a) \Rightarrow (b) If f is convex, then 4.5 holds. By setting $t := 1$ and $d := y - x$ we obtain (b).

(b) \Rightarrow (c) Let $x, y \in \mathbb{R}^n$. From (b) we have

$$f(y) \geq f(x) + \nabla f(x)^T(y - x)$$

and

$$f(x) \geq f(y) + \nabla f(y)^T(x - y).$$

By adding these two inequalities we obtain (c).

(c) \Rightarrow (b) Let $x, y \in \mathbb{R}^n$. By the Mean Value Theorem there exists $0 < \lambda < 1$ such that

$$f(y) - f(x) = \nabla f(x_\lambda)^T(y - x)$$

where $x_\lambda := \lambda y + (1 - \lambda)x$. By hypothesis,

$$\begin{aligned} 0 &\leq [\nabla f(x_\lambda) - \nabla f(x)]^T(x_\lambda - x) \\ &= \lambda[\nabla f(x_\lambda) - \nabla f(x)]^T(y - x) \\ &= \lambda[f(y) - f(x) - \nabla f(x)^T(y - x)]. \end{aligned}$$

Hence $f(y) \geq f(x) + \nabla f(x)^T(y - x)$.

(b) \Rightarrow (a) Let $x, y \in \mathbb{R}^n$ and set

$$\alpha := \max_{\lambda \in [0,1]} \varphi(\lambda) := [f(\lambda y + (1 - \lambda)x) - (\lambda f(y) + (1 - \lambda)f(x))].$$

We need to show that $\alpha \leq 0$. Since $[0, 1]$ is compact and φ is continuous, there is a $\lambda \in [0, 1]$ such that $\varphi(\lambda) = \alpha$. If λ equals zero or one, we are done. Hence we may as well assume that $0 < \lambda < 1$ in which case

$$0 = \varphi'(\lambda) = \nabla f(x_\lambda)^T(y - x) + f(x) - f(y)$$

where $x_\lambda = x + \lambda(y - x)$, or equivalently

$$\lambda f(y) = \lambda f(x) - \nabla f(x_\lambda)^T(x - x_\lambda).$$

But then

$$\begin{aligned} \alpha &= f(x_\lambda) - (f(x) + \lambda(f(y) - f(x))) \\ &= f(x_\lambda) + \nabla f(x_\lambda)^T(x - x_\lambda) - f(x) \\ &\leq 0 \end{aligned}$$

by (b).

2) Suppose f is convex and let $x, d \in \mathbb{R}^n$, then by (b) of Part (1),

$$f(x + td) \geq f(x) + t\nabla f(x)^T d$$

for all $t \in \mathbb{R}$. Replacing the left hand side of this inequality with its second-order Taylor expansion yields the inequality

$$f(x) + t\nabla f(x)^T d + \frac{t^2}{2} d^T \nabla^2 f(x) d + o(t^2) \geq f(x) + t\nabla f(x)^T d,$$

or equivalently,

$$\frac{1}{2}d^T \nabla^2 f(x) d + \frac{o(t^2)}{t^2} \geq 0.$$

Letting $t \rightarrow 0$ yields the inequality

$$d^T \nabla^2 f(x) d \geq 0.$$

Since d was arbitrary, $\nabla^2 f(x)$ is positive semi-definite.

Conversely, if $x, y \in \mathbb{R}^n$, then by the Mean Value Theorem there is a $\lambda \in (0, 1)$ such that

$$f(y) = f(x) + \nabla f(x)^T (y - x) + \frac{1}{2} (y - x)^T \nabla^2 f(x_\lambda) (y - x)$$

where $x_\lambda = \lambda y + (1 - \lambda)x$. Hence

$$f(y) \geq f(x) + \nabla f(x)^T (y - x)$$

since $\nabla^2 f(x_\lambda)$ is positive semi-definite. Therefore, f is convex by (b) of Part (1). □

Convexity is also preserved by certain operations on convex functions. A few of these are given below.

THEOREM 4.6. *Let $f : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$, $h : \mathbb{R}^s \times \mathbb{R}^k \rightarrow \overline{\mathbb{R}}$ and $f_\nu : \mathbb{R}^n \rightarrow \overline{\mathbb{R}}$ be convex functions for $\nu \in N$ where N is an arbitrary index set, and let $\nu_i \in N$ and $\lambda_i \geq 0$, $i = 1, \dots, m$. Then the following functions are also convex.*

- (1) $\phi \circ f$, where $\phi : \mathbb{R} \rightarrow \mathbb{R}$ is any non-decreasing function on \mathbb{R} .
- (2) $f(x) := \sum_{i=1}^m \lambda_i f_{\nu_i}(x)$ (Non-negative linear combinations)
- (3) $f(x) := \max_{\nu \in N} f_\nu(x)$ (pointwise max)
- (4) $f(x) := \sup \left\{ \sum_{i=1}^m f_{\nu_i}(x^i) \mid x = \sum_{i=1}^m x^i \right\}$ (infimal convolution)
- (5) $f^*(y) := \sup_{x \in \mathbb{R}^n} [y^T x - f(x)]$ (convex conjugation)
- (6) $\psi(y) = \inf_{x \in \mathbb{R}^s} h(x, y)$ (infimal projection)

4.0.1. More on the Directional Derivative. It is a powerful fact that convex functions are directionally differentiable at every point of their domain in every direction. But this is just the beginning of the story. The directional derivative of a convex function possesses several other important and surprising properties. We now develop a few of these.

DEFINITION 4.3. *Let $h : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$. We say that h is positively homogeneous if*

$$h(\lambda x) = \lambda h(x) \quad \text{for all } x \in \mathbb{R} \text{ and } \lambda > 0.$$

We say that h is subadditive if

$$h(x + y) \leq h(x) + h(y) \quad \text{for all } x, y \in \mathbb{R}.$$

Finally, we say that h is sublinear if it is both positively homogeneous and subadditive.

There are numerous important examples of sublinear functions (as we shall soon see), but perhaps the most familiar of these is the norm $\|x\|$. Positive homogeneity is obvious and subadditivity is simply the triangle inequality. In a certain sense the class of sublinear functions is a generalization of norms. It is also important to note that sublinear functions are always convex functions. Indeed, given $x, y \in \text{dom } h$ and $0 \leq \lambda \leq 1$,

$$\begin{aligned} h(\lambda x + (1 - \lambda)y) &\leq h(\lambda x) + h((1 - \lambda)y) \\ &= \lambda h(x) + (1 - \lambda)h(y). \end{aligned}$$

THEOREM 4.7. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ be a convex function. Then at every point $x \in \text{dom } f$ the directional derivative $f'(x; d)$ is a sublinear function of the d argument, that is, the function $f'(x; \cdot) : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is sublinear. Thus, in particular, the function $f'(x; \cdot)$ is a convex function.*

REMARK 4.2. *Since f is convex and $x \in \text{dom } f$, $f'(x; d)$ exists for all $d \in \mathbb{R}^n$.*

PROOF. Let $x \in \text{dom } f$, $d \in \mathbb{R}^n$, and $\lambda > 0$. Then

$$\begin{aligned}
 f'(x; \lambda d) &= \lim_{t \downarrow 0} \frac{f(x + t\lambda d) - f(x)}{t} \\
 &= \lim_{t \downarrow 0} \lambda \frac{f(x + t\lambda d) - f(x)}{\lambda t} \\
 &= \lambda \lim_{(\lambda t) \downarrow 0} \frac{f(x + (t\lambda)d) - f(x)}{(\lambda t)} \\
 &= \lambda f'(x; d),
 \end{aligned}$$

showing that $f'(x; \cdot)$ is positively homogeneous.

Next let $d_1, d_2 \in \mathbb{R}^n$, Then

$$\begin{aligned}
 f'(x; d_1 + d_2) &= \lim_{t \downarrow 0} \frac{f(x + t(d_1 + d_2)) - f(x)}{t} \\
 &= \lim_{t \downarrow 0} \frac{f(\frac{1}{2}(x + 2td_1) + \frac{1}{2}(x + 2td_2)) - f(x)}{t} \\
 &\leq \lim_{t \downarrow 0} \frac{\frac{1}{2}f(x + 2td_1) + \frac{1}{2}f(x + 2td_2) - f(x)}{t} \\
 &\leq \lim_{t \downarrow 0} \frac{\frac{1}{2}(f(x + 2td_1) - f(x)) + \frac{1}{2}(f(x + 2td_2) - f(x))}{t} \\
 &= \lim_{t \downarrow 0} \frac{f(x + 2td_1) - f(x)}{2t} + \lim_{t \downarrow 0} \frac{f(x + 2td_2) - f(x)}{2t} \\
 &= f'(x; d_1) + f'(x; d_2),
 \end{aligned}$$

showing that $f'(x; \cdot)$ is subadditive and completing the proof. \square

Optimality Conditions for Constrained Optimization

1. First-Order Conditions

In this section we consider first-order optimality conditions for the constrained problem

$$\mathcal{P} : \begin{array}{ll} \text{minimize} & f_0(x) \\ \text{subject to} & x \in \Omega, \end{array}$$

where $f_0 : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuously differentiable and $\Omega \subset \mathbb{R}^n$ is closed and non-empty. The first step in the analysis of the problem \mathcal{P} is to derive conditions that allow us to recognize when a particular vector \bar{x} is a solution, or local solution, to the problem. For example, when we minimize a function of one variable we first take the derivative and see if it is zero. If it is, then we take the second derivative and check that it is positive. If this is also true, then we know that the point under consideration is a local minimizer of the function. Of course, the presence of constraints complicates this kind of test.

To understand how an optimality test might be derived in the constrained case, let us first suppose that we are at a feasible point x and we wish to find a better point \tilde{x} . That is, we wish to find a point \tilde{x} such that $\tilde{x} \in \Omega$ and $f(\tilde{x}) < f(x)$. As in the unconstrained case, one way to do this is to find a direction d in which the directional derivative of f in the direction d is negative: $f'(x; d) < 0$. We know that for such directions we can reduce the value of the function by moving away from the point x in the direction d . However, moving in such a direction may violate feasibility. That is, it may happen that $x + td \notin \Omega$ regardless of how small we take $t > 0$. To avoid this problem, we consider the notion of a *feasible direction*.

DEFINITION 1.1. [FEASIBLE DIRECTIONS]

Given a subset Ω of \mathbb{R}^n and a point $x \in \Omega$, we say that a direction $d \in \mathbb{R}^n$ is a *feasible direction* for Ω at x if there is a $\bar{t} > 0$ such that $x + td \in \Omega$ for all $t \in [0, \bar{t}]$.

THEOREM 1.1. If \bar{x} is a local solution to the problem \mathcal{P} , then $f'(\bar{x}; d) \geq 0$ for all feasible directions d for Ω at \bar{x} for which $f'(\bar{x}; d)$ exists.

PROOF. The proof is a straightforward application of the definitions. If the result were false, then there would be a direction of descent for f at \bar{x} that is also a feasible direction for Ω at \bar{x} . But then moving a little bit in this direction both keeps us in Ω and strictly reduces the value of f . This contradicts the assumption that \bar{x} is a local solution. Therefore, the result must be true. \square

Unfortunately, this result is insufficient in many important cases. The insufficiency comes from the dependence on the notion of *feasible direction*. For example, if

$$\Omega = \{(x_1, x_2)^T : x_1^2 + x_2^2 = 1\},$$

then the only feasible direction at any point of Ω is the zero direction. Hence, regardless of the objective function f and the point $\bar{x} \in \Omega$, we have that $f'(\bar{x}; d) \geq 0$ for every feasible direction to Ω at \bar{x} . In this case, Theorem 1.1 has no content.

To overcome this deficiency we introduce a general notion of *tangency* that considers all directions d pointing into Ω at $x \in \Omega$ in a limiting sense. Define the *tangent cone* to Ω at a point $x \in \Omega$ to be the set of limiting directions obtained from sequences in Ω that converge to x . Specifically, the tangent cone is given by

$$T(x | \Omega) := \{d : \exists \tau_i \searrow 0 \text{ and } \{x_i\} \subset \Omega, \text{ with } x_i \rightarrow x, \text{ such that } \tau_i^{-1}(x_i - x) \rightarrow d\}.$$

EXAMPLE 1.1. (1) If $\Omega = \{x : Ax = b\}$, where $A \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$, then $T(x | \Omega) = \text{Nul}(A)$ for every $x \in \Omega$.

Reason: Let $x \in \Omega$. Note that if $d \in \text{Nul}(A)$, then for every $t \geq 0$ we have $A(x + td) = Ax + tAd = Ax = b$ so that $d \in T(x|\Omega)$. Since $d \in \text{Nul}(A)$ was chosen arbitrarily, this implies that $\text{Nul}(A) \subset T(x|\Omega)$. Hence we only need to establish the reverse inclusion to verify the equivalence of these sets.

Let $d \in T(x|\Omega)$. Then, by definition, there are sequences $t_i \downarrow 0$ and $\{x^i\} \subset \Omega$ with $x^i \rightarrow x$ such that $d^i \rightarrow d$ where $d^i = t_i^{-1}(x^i - x)$, $i = 1, 2, \dots$. Note that

$$Ad^i = t_i^{-1}A(x^i - x) = t_i^{-1}[Ax^i - Ax] = t_i^{-1}[b - b] = 0 \quad \forall i = 1, 2, \dots$$

Therefore, $Ad = \lim_{i \rightarrow \infty} Ad^i = 0$ so that $d \in \text{Nul}(A)$. Since d was chosen arbitrarily from $T(x|\Omega)$, we have $T(x|\Omega) \subset \text{Nul}(A)$ which proves the equivalence.

(2) If $\Omega = \{(x_1, x_2)^T : x_1^2 + x_2^2 = 1\}$, then $T(x|\Omega) = \{(y_1, y_2) : x_1 y_1 + x_2 y_2 = 0\}$.

(3) A convex set is said to be polyhedral if it can be represented as the solution set of a finite number of linear equality and/or inequality constraints. Thus, for example the constraint region for an LPs is a convex polyhedron. If it is assumed that Ω is a convex polyhedron, then

$$T(x|\Omega) = \bigcup_{\lambda \geq 0} \lambda(\Omega - x) = \{\lambda(y - x) \mid \lambda \geq 0, y \in \Omega\}.$$

(4) If Ω is a convex subset of \mathbb{R}^n , then

$$T(x|\Omega) = \overline{\bigcup_{\lambda \geq 0} \lambda(\Omega - x)} = \text{cl} \{\lambda(y - x) \mid \lambda \geq 0, y \in \Omega\}.$$

THEOREM 1.2. [BASIC CONSTRAINED FIRST-ORDER NECESSARY CONDITIONS]

Suppose that the function $f_0 : \mathbb{R}^n \rightarrow \mathbb{R}$ in \mathcal{P} is continuously differentiable near the point $\bar{x} \in \Omega$. If \bar{x} is a local solution to \mathcal{P} , then

$$f'_0(\bar{x}; d) \geq 0 \quad \text{for all } d \in T(\bar{x}|\Omega).$$

PROOF. Note that the MVT (Mean Value Theorem) implies that

$$f'_0(\bar{x}; d) = \lim_{\tau \searrow 0} \frac{f_0(\bar{x} + \tau d) - f_0(\bar{x})}{\tau} = \lim_{\substack{s \rightarrow d \\ \tau \searrow 0}} \frac{f_0(\bar{x} + \tau s) - f_0(\bar{x})}{\tau}$$

since f_0 is continuously differentiable.

Suppose \bar{x} is a local solution to \mathcal{P} and let $d \in T(\bar{x}|\Omega)$. Since $d \in T(\bar{x}|\Omega)$, there is a sequence $\{x_i\} \subset \Omega$ and $t_i \downarrow 0$ such that $x_i \rightarrow \bar{x}$ and $s_i = t_i^{-1}(x_i - \bar{x}) \rightarrow d$. Note that $\bar{x} + t_i d \approx \bar{x} + t_i s_i = x_i$, and so $f(\bar{x} + t_i s_i) = f(x_i) \geq f(\bar{x})$. Using the representation of the directional derivative given above, we obtain

$$f'_0(\bar{x}; d) = \lim_{\substack{s \rightarrow d \\ \tau \searrow 0}} \frac{f_0(\bar{x} + \tau s) - f_0(\bar{x})}{\tau} = \lim_{i \rightarrow \infty} \frac{f_0(\bar{x} + t_i s_i) - f_0(\bar{x})}{t_i} = \lim_{i \rightarrow \infty} \frac{f_0(x_i) - f_0(\bar{x})}{t_i} \geq 0.$$

□

This general result is not particularly useful on its own since it refers to the abstract notion of tangent cone. In order to make it useful, we need to be able to compute the tangent cone once a representation for Ω is given. We now show how this can be done.

We begin by assuming that Ω has the form

$$(76) \quad \Omega := \{x : f_i(x) \leq 0, i = 1, \dots, s, f_i(x) = 0, i = s + 1, \dots, m\},$$

where each $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuously differentiable on \mathbb{R}^n . Observe that if $x \in \Omega$ and $d \in T(x|\Omega)$ then there are sequences $\{x_k\} \subset \Omega$ and $\tau_k \searrow 0$ with $x_k \rightarrow x$ such that $\tau_k^{-1}(x_k - x) \rightarrow d$. Setting $d_k = \tau_k^{-1}(x_k - x)$ for all k we have that

$$f'_i(x; d) = \lim_{k \rightarrow \infty} \frac{f_i(x + \tau_k d_k) - f_i(x)}{\tau_k}$$

equals 0 for $i \in \{s + 1, \dots, m\}$ and is less than or equal to 0 for $i \in I(x)$ where

$$I(x) := \{i : i \in \{1, \dots, s\}, f_i(x) = 0\}.$$

Consequently,

$$T(x|\Omega) \subset \{d : \nabla f_i(x)^T d \leq 0, i \in I(x), \nabla f_i(x)^T d = 0, i = s + 1, \dots, m\}.$$

The set on the right hand side of this inclusion is a computationally tractable. Moreover, in a certain sense, the cases where these two sets do not coincide are exceptional. For this reason we make the following definition.

DEFINITION 1.2. [REGULARITY]

We say that the set Ω is regular at $x \in \Omega$ if

$$T(x|\Omega) = \{d \in \mathbb{R}^n : f'_i(x; d) \leq 0, i \in I(x), f'_i(x; d) = 0 \text{ } i = s+1, \dots, m\}.$$

But it is important to note that not every set is regular.

EXERCISE 1.1. Graph the set

$$\Omega := \{x \in \mathbb{R}^2 \mid -x_1^3 \leq x_2 \leq x_1^3\},$$

and show that it is not regular at the origin. This is done by first showing that

$$T_\Omega(0) = \{(d_1, d_2)^T \mid d_1 \geq 0, d_2 = 0\}.$$

Then set

$$f_1(x_1, x_2) = -x_1^3 - x_2 \quad \text{and} \quad f_2(x_1, x_2) = -x_1^3 + x_2,$$

so that $\Omega = \{(x_1, x_2)^T \mid f_1(x_1, x_2) \leq 0, f_2(x_1, x_2) \leq 0\}$. Finally, show that

$$\{d \mid \nabla f_1(0, 0)^T d \leq 0, \nabla f_2(0, 0)^T d \leq 0\} = \{(d_1, d_2)^T \mid d_2 = 0\} \neq T_\Omega(0).$$

Next let us suppose we are at a given point $x \in \Omega$ and that we wish to obtain a new point $x_+ = x + td$ for which $f(x_+) < f(x)$ for some direction $d \in \mathbb{R}^n$ and steplength $t > 0$. A good candidate for a search direction d is one that minimizes $f'(x; d)$ over all directions that point into Ω up to first-order. That is, we should minimize $\nabla f(x)^T d$ over the set of tangent directions. Remarkably, this search for a *feasible direction of steepest descent* can be posed as the following linear program (assuming regularity):

$$(77) \quad \begin{array}{ll} \max & (-\nabla f_0(\bar{x}))^T d \\ \text{subject to} & \nabla f_i(\bar{x})^T d \leq 0 \quad i \in I(\bar{x}) \\ & \nabla f_i(\bar{x})^T d = 0 \quad i = s+1, \dots, m. \end{array}$$

The dual of (77) is the linear program

$$(78) \quad \begin{array}{ll} \min & 0 \\ \text{subject to} & \sum_{i \in I(\bar{x})} u_i \nabla f_i(\bar{x}) + \sum_{i=s+1}^m u_i \nabla f_i(\bar{x}) = -\nabla f_0(\bar{x}) \\ & 0 \leq u_i, \quad i \in I(\bar{x}). \end{array}$$

If we assume that \bar{x} is a local solution to \mathcal{P} , Theorem 1.2 tells us that the maximum in (77) is less than or equal to zero. But $d = 0$ is feasible for (77), hence the maximum value in (77) is zero. Therefore, by the Strong Duality Theorem for Linear Programming, the linear program (78) is feasible, that is, there exist scalars u_i , $i \in I(\bar{x}) \cup \{s+1, \dots, m\}$ with $u_i \geq 0$ for $i \in I(\bar{x})$ such that

$$(79) \quad 0 = \nabla f_0(\bar{x}) + \sum_{i \in I(\bar{x})} u_i \nabla f_i(\bar{x}) + \sum_{i=s+1}^m u_i \nabla f_i(\bar{x}).$$

This observation yields the following result.

THEOREM 1.3. [CONSTRAINED FIRST-ORDER OPTIMALITY CONDITIONS]

Let $\bar{x} \in \Omega$ be a local solution to \mathcal{P} at which Ω is regular. Then there exist $u \in \mathbb{R}^m$ such that

- (1) $0 = \nabla_x L(\bar{x}, u)$,
- (2) $0 = u_i f_i(\bar{x})$ for $i = 1, \dots, s$, and
- (3) $0 \leq u_i$, $i = 1, \dots, s$,

where the mapping $L : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$ is defined by

$$L(x, u) := f_0(x) + \sum_{i=1}^m u_i f_i(x)$$

and is called the Lagrangian for the problem \mathcal{P} .

PROOF. For $i \in I(\bar{x}) \cup \{s+1, \dots, m\}$ let u_i be as given in (79) and for $i \in \{1, \dots, s\} \setminus I(\bar{x})$ set $u_i = 0$. Then this choice of $u \in \mathbb{R}^m$ satisfies (1)–(3) above. \square

DEFINITION 1.3. [KKT CONDITIONS]

Let $x \in \mathbb{R}^n$ and $u \in \mathbb{R}^m$. We say that (x, u) is a Karush-Kuhn-Tucker (KKT) pair for \mathcal{P} if

- (1) $f_i(x) \leq 0$ $i = 1, \dots, s$, $f_i(x) = 0$ $i = s + 1, \dots, m$ (Primal feasibility),
- (2) $u_i \geq 0$ for $i = 1, \dots, s$ (Dual feasibility),
- (3) $0 = u_i f_i(x)$ for $i = 1, \dots, s$ (complementarity), and
- (4) $0 = \nabla_x L(x, u)$ (stationarity of the Lagrangian).

Given $x \in \mathbb{R}^n$, if there is a $u \in \mathbb{R}^m$ such that (x, u) is a Karush-Kuhn-Tucker pair for \mathcal{P} , then we say that x is a KKT point for \mathcal{P} (we also refer to such an x as a stationary point for \mathcal{P}). \square

2. Regularity and Constraint Qualifications

We now briefly discuss conditions that yield the regularity of Ω at a point $x \in \Omega$. These conditions should be testable in the sense that there is a finitely terminating algorithm that can determine whether they are satisfied or not satisfied. The condition that we will concentrate on is the so called *Mangasarian-Fromovitz constraint qualification* (MFCQ).

DEFINITION 2.1. [MFCQ]

We say that a point $x \in \Omega$ satisfies the Mangasarian-Fromovitz constraint qualification (or MFCQ) at x if

- (1) there is a $d \in \mathbb{R}^n$ such that

$$\begin{aligned} \nabla f_i(x)^T d &< 0 \text{ for } i \in I(x), \\ \nabla f_i(x)^T d &= 0 \text{ for } i = s + 1, \dots, m, \end{aligned}$$

and

- (2) the gradients $\{\nabla f_i(x) \mid i = s + 1, \dots, m\}$ are linearly independent.

We have the following key result which we shall not prove.

THEOREM 2.1. [MFCQ \rightarrow REGULARITY] Let $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$, $i = 1, 2, \dots, m$ be C^1 near $\bar{x} \in \Omega$. If the MFCQ holds at \bar{x} , then Ω is regular at \bar{x} .

The MFCQ is algorithmically verifiable. This is seen by considering the LP

$$(80) \quad \begin{aligned} \min \quad & 0 \\ \text{subject to} \quad & \nabla f_i(\bar{x})^T d \leq -1 \quad i \in I(\bar{x}) \\ & \nabla f_i(\bar{x})^T d = 0 \quad i = s + 1, \dots, m. \end{aligned}$$

Clearly, the MFCQ is satisfied at \bar{x} if and only if the above LP is feasible and the gradients $\{\nabla f_i(\bar{x}) \mid i = s + 1, \dots, m\}$ are linearly independent. This observation also leads to a dual characterization of the MFCQ by considering the dual of the LP (80).

LEMMA 2.1. [DUAL MFCQ]

The MFCQ is satisfied at a point $\bar{x} \in \Omega$ if and only if the only solution to the system

$$\begin{aligned} \sum_{i=1}^m u_i \nabla f_i(\bar{x}) &= 0, \\ u_i f_i(\bar{x}) &= 0 \quad i = 1, 2, \dots, s, \text{ and} \\ u_i &\geq 0 \quad i = 1, 2, \dots, s, \end{aligned}$$

is $u_i = 0$, $i = 1, 2, \dots, m$.

PROOF. The dual the LP (80) is the LP

$$(81) \quad \begin{aligned} \min \quad & \sum_{i \in I(\bar{x})} u_i \\ \text{subject to} \quad & \sum_{i \in I(\bar{x})} u_i \nabla f_i(\bar{x}) + \sum_{i=s+1}^m u_i \nabla f_i(\bar{x}) = 0 \\ & 0 \leq u_i, \quad i \in I(\bar{x}). \end{aligned}$$

This LP is always feasible, simply take all u_i 's equal to zero. Hence, by the Strong Duality Theorem of Linear Programming, the LP (80) is feasible if and only if the LP (81) is finite valued in which case the optimal value in both is zero. That is, the MFCQ holds at \bar{x} if and only if the optimal value in (81) is zero and the gradients

$\{\nabla f_i(\bar{x}) \mid i = s+1, \dots, m\}$ are linearly independent. The latter statement is equivalent to the statement that the only solution to the system

$$\begin{aligned} \sum_{i=1}^m u_i \nabla f_i(\bar{x}) &= 0, \\ u_i f_i(\bar{x}) &= 0 \quad i = 1, 2, \dots, s, \text{ and} \\ u_i &\geq 0 \quad i = 1, 2, \dots, s, \end{aligned}$$

is $u_i = 0, i = 1, 2, \dots, m$. □

Techniques similar to these show that the MFCQ is a local property. That is, if it is satisfied at a point then it must be satisfied on a neighborhood of that point. The MFCQ is a powerful tool in the analysis of constraint systems as it implies many useful properties. One such property is established in the following result.

THEOREM 2.2. [MFCQ \rightarrow COMPACT MULTIPLIER SET]

Let $\bar{x} \in \Omega$ be a local solution to \mathcal{P} at which the set of Karush-Kuhn-Tucker multipliers

$$(82) \quad KKT(\bar{x}) := \left\{ u \in \mathbb{R}^m \left| \begin{array}{l} \nabla_x L(\bar{x}, u) = 0 \\ u_i f_i(\bar{x}) = 0, \quad i = 1, 2, \dots, s, \\ 0 \leq u_i, \quad i = 1, 2, \dots, s \end{array} \right. \right\}$$

is non-empty. Then $KKT(\bar{x})$ is a compact set if and only if the MFCQ is satisfied at \bar{x} .

PROOF. (\Rightarrow) If MFCQ is not satisfied at \bar{x} , then from the Strong Duality Theorem for linear programming, Lemma 2.1, and the LP (81) guarantees the existence of a non-zero vector $\bar{u} \in \mathbb{R}^m$ satisfying

$$\sum_{i=1}^m u_i \nabla f_i(\bar{x}) = 0 \text{ and } 0 \leq u_i \text{ with } 0 = u_i f_i(\bar{x}) \text{ for } i = 1, 2, \dots, s.$$

Then for each $u \in KKT(\bar{x})$ we have that $u + t\bar{u} \in KKT(\bar{x})$ for all $t > 0$. Consequently, $KKT(\bar{x})$ cannot be compact.

(\Leftarrow) If $KKT(\bar{x})$ is not compact, there is a sequence $\{u^j\} \subset KKT(\bar{x})$ with $\|u^j\| \uparrow +\infty$. With no loss of generality, we may assume that

$$\frac{u^j}{\|u^j\|} \rightarrow u.$$

But then

$$\begin{aligned} u_i &\geq 0, \quad i = 1, 2, \dots, s, \\ u_i f_i(\bar{x}) &= \lim_{i \rightarrow \infty} \frac{u^j}{\|u^j\|} f_i(\bar{x}) = 0, \quad i = 1, 2, \dots, s, \text{ and} \\ \sum_{i=1}^m u_i f_i(\bar{x}) &= \lim_{i \rightarrow \infty} \frac{\nabla_x L(\bar{x}, u^j)}{\|u^j\|} = 0. \end{aligned}$$

Hence, by Lemma 2.1, the MFCQ cannot be satisfied at \bar{x} . □

Before closing this section we introduce one more constraint qualification. This is the so called *LI* condition and is associated with the uniqueness of the multipliers..

DEFINITION 2.2 (LINEAR INDEPENDENCE CONDITION). The *LI condition* is said to be satisfied at the point $x \in \Omega$ if the constraint gradients

$$\{\nabla f_i(x) \mid i \in I(x) \cup \{s+1, \dots, m\}\}$$

are linearly independent.

Clearly, the LI condition implies the MFCQ. However, it is a much stronger condition in the presence of inequality constraints. In particular, the LI condition implies the uniqueness of the multipliers at a local solution to \mathcal{P} .

3. Second-Order Conditions

Second-order conditions are introduced by way of the Lagrangian. As is illustrated in the following result, the multipliers provide a natural way to incorporate the curvature of the constraints.

THEOREM 3.1. [CONSTRAINED SECOND-ORDER SUFFICIENCY]

Let Ω have representation (76) and suppose that each of the functions f_i , $i = 0, 1, 2, \dots, m$ are C^2 . Let $\bar{x} \in \Omega$. If $(\bar{x}, \bar{u}) \in \mathbb{R}^n \times \mathbb{R}^m$ is a Karush-Kuhn-Tucker pair for \mathcal{P} such that

$$d^T \nabla_x^2 L(\bar{x}, \bar{u}) d > 0$$

for all $d \in T_\Omega(\bar{x})$, $d \neq 0$, with $\nabla f_0(\bar{x})^T d = 0$, then there is an $\epsilon > 0$ and $\nu > 0$ such that

$$f_0(x) \geq f_0(\bar{x}) + \nu \|x - \bar{x}\|^2$$

for every $x \in \Omega$ with $\|x - \bar{x}\| \leq \epsilon$, in particular \bar{x} is a strict local solution to \mathcal{P} .

PROOF. Suppose to the contrary that no such $\epsilon > 0$ and $\nu > 0$ exist, then there exist sequences $\{x_k\} \subset \Omega$, $\{\nu_k\} \subset \mathbb{R}_+$ such that $x_k \rightarrow \bar{x}$, $\nu_k \downarrow 0$, and

$$f_0(x_k) \leq f_0(\bar{x}) + \nu_k \|x_k - \bar{x}\|^2$$

for all $k = 1, 2, \dots$. For every $x \in \Omega$ we know that $\bar{u}^T f(x) \leq 0$ and $0 = \bar{u}^T f(\bar{x})$ where the i th component of $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is f_i . Hence

$$\begin{aligned} L(x_k, \bar{u}) \leq f_0(x_k) &\leq f_0(\bar{x}) + \nu_k \|x_k - \bar{x}\|^2 \\ &= L(\bar{x}, \bar{u}) + \nu_k \|x_k - \bar{x}\|^2. \end{aligned}$$

Therefore,

$$(83) \quad f_0(\bar{x}) + \nabla f_0(\bar{x})^T (x_k - \bar{x}) + o(\|x_k - \bar{x}\|) \leq f_0(\bar{x}) + \nu_k \|x_k - \bar{x}\|^2$$

and

$$(84) \quad \begin{aligned} L(\bar{x}, \bar{u}) &+ \nabla_x L(\bar{x}, \bar{u})^T (x_k - \bar{x}) \\ &+ \frac{1}{2} (x_k - \bar{x})^T \nabla_x^2 L(\bar{x}, \bar{u}) (x_k - \bar{x}) + o(\|x_k - \bar{x}\|^2) \\ &\leq L(\bar{x}, \bar{u}) + \nu_k \|x_k - \bar{x}\|^2. \end{aligned}$$

With no loss of generality, we can assume that

$$d_k := \frac{x_k - \bar{x}}{\|x_k - \bar{x}\|} \rightarrow \bar{d} \in T_\Omega(\bar{x}).$$

Dividing (83) through by $\|x_k - \bar{x}\|$ and taking the limit we find that $\nabla f_0(\bar{x})^T \bar{d} \leq 0$. Since

$$T_\Omega(\bar{x}) \subset \{d : \nabla f_i(\bar{x})^T d \leq 0, i \in I(\bar{x}), \nabla f_i(\bar{x})^T d = 0, i = s+1, \dots, m\},$$

we have $\nabla f_i(\bar{x})^T \bar{d} \leq 0$, $i \in I(\bar{x}) \cup \{0\}$ and $\nabla f_i(\bar{x})^T \bar{d} = 0$ for $i = s+1, \dots, m$. On the other hand, (\bar{x}, \bar{u}) is a Karush-Kuhn-Tucker point so

$$\nabla f_0(\bar{x})^T \bar{d} = - \sum_{i \in I(\bar{x})} \bar{u}_i \nabla f_i(\bar{x})^T \bar{d} \geq 0.$$

Hence $\nabla f_0(\bar{x})^T \bar{d} = 0$, so that

$$\bar{d}^T \nabla_x^2 L(\bar{x}, \bar{u}) \bar{d} > 0.$$

But if we divide (84) by $\|x_k - \bar{x}\|^2$ and take the limit, we arrive at the contradiction

$$\frac{1}{2} \bar{d}^T \nabla_x^2 L(\bar{x}, \bar{u}) \bar{d} \leq 0,$$

whereby the result is established. \square

The assumptions required to establish Theorem 3.1 are somewhat strong but they do lead to a very practical and, in many cases, satisfactory second-order sufficiency result. In order to improve on this result one requires a much more sophisticated mathematical machinery. We do not take the time to develop this machinery. Instead we simply state a very general result. The statement of this result employs the entire set of Karush-Kuhn-Tucker multipliers $KKT(\bar{x})$.

THEOREM 3.2 (GENERAL CONSTRAINED SECOND-ORDER NECESSITY AND SUFFICIENCY). Let $\bar{x} \in \Omega$ be a point at which Ω is regular.

(1) If \bar{x} is a local solution to \mathcal{P} , then $KKT(\bar{x}) \neq \emptyset$, and for every $d \in T(\bar{x} | \Omega)$ there is a $u \in KKT(\bar{x})$ such that

$$d^T \nabla_x^2 L(\bar{x}, u) d \geq 0.$$

(2) If $KKT(\bar{x}) \neq \emptyset$, and for every $d \in T(\bar{x} | \Omega)$, $d \neq 0$, for which $\nabla f_0(\bar{x})^T d = 0$ there is a $u \in KKT(\bar{x})$ such that

$$d^T \nabla_x^2 L(\bar{x}, u) d > 0,$$

then there is an $\epsilon > 0$ and $\nu > 0$ such that

$$f_0(x) \geq f_0(\bar{x}) + \nu \|x - \bar{x}\|^2$$

for every $x \in \Omega$ with $\|x - \bar{x}\| \leq \epsilon$, in particular \bar{x} is a strict local solution to \mathcal{P} .

4. Optimality Conditions in the Presence of Convexity

As we saw in the unconstrained case, convexity can have profound implications for optimality and optimality conditions. To begin with, we have the following very powerful result whose proof is identical to the proof in the unconstrained case.

THEOREM 4.1. [CONVEXITY+LOCAL OPTIMALITY→GLOBAL OPTIMALITY]

Suppose that $f_0 : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex and that $\Omega \subset \mathbb{R}^n$ is a convex set. If $\bar{x} \in \mathbb{R}^n$ is a local solution to \mathcal{P} , then \bar{x} is a global solution to \mathcal{P} .

PROOF. Suppose there is a $\hat{x} \in \Omega$ with $f_0(\hat{x}) < f_0(\bar{x})$. Let $\epsilon > 0$ be such that

$$f_0(\bar{x}) \leq f_0(x) \quad \text{whenever} \quad \|x - \bar{x}\| \leq \epsilon \text{ and } x \in \Omega,$$

and

$$\epsilon < 2\|\bar{x} - \hat{x}\|.$$

Set $\lambda := \epsilon(2\|\bar{x} - \hat{x}\|)^{-1} < 1$ and $x_\lambda := \bar{x} + \lambda(\hat{x} - \bar{x}) \in \Omega$. Then $\|x_\lambda - \bar{x}\| \leq \epsilon/2$ and $f_0(x_\lambda) \leq (1 - \lambda)f_0(\bar{x}) + \lambda f_0(\hat{x}) < f_0(\bar{x})$. This contradicts the choice of ϵ and so no such \hat{x} exists. \square

We also have the following first-order necessary conditions for optimality. The proof of this result again follows that for the unconstrained case.

THEOREM 4.2. [1ST-ORDER NECESSITY AND SUFFICIENCY]

Suppose that $f_0 : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex and that $\Omega \subset \mathbb{R}^n$ is a convex set, and let $\bar{x} \in \Omega$. Then the following statements are equivalent.

- (i) \bar{x} is a local solution to \mathcal{P} .
- (ii) $f'_0(\bar{x}; y - \bar{x}) \geq 0$ for all $y \in \Omega$.
- (iii) \bar{x} is a global solution to \mathcal{P} .

PROOF. The implication (i)⇒(ii) follows from Theorem 1.1 since each of the directions $d = y - \bar{x}$, $y \in \Omega$ is a feasible direction for Ω at \bar{x} due to the convexity of Ω . To see the implication (ii)⇒(iii), we again resort to the subdifferential inequality. Let y be any other point in Ω . Then $d = y - \bar{x} \in T_\Omega(\bar{x})$ and so by the subdifferential inequality we have

$$f_0(y) \geq f_0(\bar{x}) + f'_0(\bar{x}; y - \bar{x}) \geq f_0(\bar{x}).$$

Since $y \in \Omega$ was arbitrary the implication (ii)⇒(iii) follows. The implication (iii)⇒(i) is trivial. \square

The utility of this result again depends on our ability to represent the tangent cone $T_\Omega(\bar{x})$ in a computationally tractable manner. Following the general case, we assume that the set Ω has the representation (76):

$$(85) \quad \Omega := \{x : f_i(x) \leq 0, i = 1, \dots, s, f_i(x) = 0, i = s + 1, \dots, m\}.$$

The first issue we must address is to determine reasonable conditions on the functions f_i that guarantee that the set Ω is convex. We begin with the following elementary facts about convex functions and convex sets whose proofs we leave to the reader.

LEMMA 4.1. If $C_i \subset \mathbb{R}^n$, $i = 1, 2, \dots, N$, are convex sets, then so is the set $C = \bigcap_{i=1}^N C_i$.

LEMMA 4.2. If $h : \mathbb{R}^n \rightarrow \bar{\mathbb{R}}$ is a convex function, then for every $\alpha \in \mathbb{R}$ the set

$$\text{lev}_h(\alpha) = \{x \mid h(x) \leq \alpha\}$$

is a convex set.

These facts combine to give the following result.

LEMMA 4.3. *If the functions f_i , $i = 1, 2, \dots, s$ are convex and the functions f_i , $i = s + 1, \dots, m$ are linear, then the set Ω given by (85) is a convex set.*

REMARK 4.1. *Recall that a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is said to be linear if there exists $c \in \mathbb{R}^n$ and $\alpha \in \mathbb{R}$ such that $f(x) = c^T x + \alpha$.*

PROOF. Note that

$$\Omega = \left(\bigcap_{i=1}^m \text{lev}_{f_i}(0) \right) \cap \left(\bigcap_{i=s+1}^m \text{lev}_{-f_i}(0) \right),$$

where each of the functions $f_i, i = 1, \dots, m$ and $-f_i, i = s + 1, \dots, m$ is convex. Therefore, the convexity of Ω follows from Lemmas 4.2 and 4.1. \square

In order to make the link to the KKT condition in the presence of convexity, we still require the regularity of the set Ω at the point of interest \bar{x} . If the set Ω is a polyhedral convex set, i.e.

$$\Omega = \{x \mid Ax \leq a, Bx = b\}$$

for some $A \in \mathbb{R}^{s \times n}, a \in \mathbb{R}^s, B \in \mathbb{R}^{(m-s) \times n}$, and $b \in \mathbb{R}^{(m-s)}$, then the set Ω is everywhere regular (Why?). In the general convex case this may not be true. However, convexity can be used to derive a much simpler test for the regularity of non-polyhedral convex sets.

DEFINITION 4.1 (THE SLATER CONSTRAINT QUALIFICATION). *Let $\Omega \subset \mathbb{R}^n$ be as given in (85) with f_i , $i = 1, \dots, s$ convex and f_i , $i = s + 1, \dots, m$ linear. We say that Ω satisfies the Slater constraint qualification if there exists $\tilde{x} \in \Omega$ such that $f_i(\tilde{x}) < 0$ for $i = 1, \dots, s$.*

THEOREM 4.3 (CONVEXITY AND REGULARITY). *Suppose $\Omega \subset \mathbb{R}^n$ is as given in (85) with f_i , $i = 1, \dots, s$ convex and f_i , $i = s + 1, \dots, m$ linear. If either Ω is polyhedral convex or satisfies the Slater constraint qualification, then Ω is regular at every point $\bar{x} \in \Omega$ at which the function f_i , $i = 1, \dots, s$ are differentiable.*

We do not present the proof of this result as it takes us too far afield of our study. Nonetheless, we make use of this fact in the following result of the KKT conditions.

THEOREM 4.4 (CONVEXITY+REGULARITY \rightarrow (OPTIMALITY \Leftrightarrow KKT CONDITIONS)). *Let $f_0 : \mathbb{R}^n \rightarrow \mathbb{R}$ be a differentiable convex function and let Ω be as given in Lemma 4.3 where each of the function f_i , $i = 1, \dots, s$ is differentiable.*

- (i) *If $\bar{x} \in \Omega$ is a KKT point for \mathcal{P} , then \bar{x} is a global solution to \mathcal{P} .*
- (ii) *Suppose the functions f_i , $i = 0, 1, \dots, s$ are continuously differentiable. If \bar{x} is a solution to \mathcal{P} at which Ω is regular, then \bar{x} is a KKT point for \mathcal{P} .*

PROOF. Part (ii) of this theorem is just a restatement of Theorem 1.3 and so we need only prove Part (i).

Since \bar{x} is a KKT point there exists $\bar{y} \in \mathbb{R}^m$ such that (\bar{x}, \bar{y}) is a KKT pair for \mathcal{P} . Consider the function $h : \mathbb{R}^n \rightarrow \mathbb{R}$ given by

$$h(x) = L(x, \bar{y}) = f_0(x) + \sum_{i=1}^m \bar{y}_i f_i(x).$$

By construction, the function h is convex with $0 = \nabla h(\bar{x}) = \nabla_x L(\bar{x}, \bar{y})$. Therefore, \bar{x} is a global solution to the problem $\min_{x \in \mathbb{R}^n} h(x)$. Also note that for every $x \in \Omega$ we have

$$\sum_{i=1}^m \bar{y}_i f_i(x) \leq 0,$$

since $\bar{y}_i f_i(x) \leq 0$ $i = 1, \dots, s$ and $\bar{y}_i f_i(x) = 0$ $i = s + 1, \dots, m$. Consequently,

$$\begin{aligned} f_0(\bar{x}) &= h(\bar{x}) \leq h(x) = L(x, \bar{y}) \\ &= f_0(x) + \sum_{i=1}^m \bar{y}_i f_i(x) \\ &\leq f_0(x) \end{aligned}$$

for all $x \in \Omega$. This establishes Part (i). \square

If all of the functions f_i $i = 0, 1, \dots, m$ are twice continuously differentiable, then the second-order sufficiency conditions stated in Theorem 3.1 apply. However, in the presence of convexity another kind of second-order condition is possible that does not directly incorporate curvature information about the functions f_i $i = 1, \dots, m$. These second-order conditions are most appropriate when Ω is polyhedral convex.

THEOREM 4.5. [2ND-ORDER OPTIMALITY CONDITIONS FOR POLYHEDRAL CONSTRAINTS]

Let $f_0 : \mathbb{R}^n \rightarrow \mathbb{R}$ be C^2 and \bar{x} be an element of the convex set Ω .

- (1) (necessity) If $\bar{x} \in \mathbb{R}^n$ is a local solution to \mathcal{P} with Ω a polyhedral convex set, then $\nabla f_0(\bar{x})^T d \geq 0$ for all $d \in T_\Omega(\bar{x})$ and

$$d^T \nabla^2 f_0(\bar{x}) d \geq 0$$

for all $d \in T_\Omega(\bar{x})$ with $\nabla f_0(\bar{x})^T d = 0$.

- (2) (sufficiency) If $\bar{x} \in \mathbb{R}^n$ is such that $\nabla f_0(\bar{x})^T (y - \bar{x}) \geq 0$ for all $d \in T_\Omega(\bar{x})$ and

$$d^T \nabla^2 f_0(\bar{x}) d > 0$$

for all $d \in T_\Omega(\bar{x}) \setminus \{0\}$ with $\nabla f_0(\bar{x})^T d = 0$, then there exist $\epsilon, \nu > 0$ such that

$$f_0(x) \geq f_0(\bar{x}) + \nu \|x - \bar{x}\|^2$$

for all $x \in \Omega$ with $\|x - \bar{x}\| \leq \epsilon$.

PROOF. (1) Since Ω is polyhedral convex, we have $T_\Omega(\bar{x}) = \bigcup_{\lambda \geq 0} (\Omega - \bar{x})$. Therefore, the fact that $\nabla f_0(\bar{x})^T d \geq 0$ for all $d \in T_\Omega(\bar{x})$ follows from Theorem 4.2. Next let $d \in T_\Omega(\bar{x}) = \bigcup_{\lambda \geq 0} (\Omega - \bar{x})$ be such that $\nabla f_0(\bar{x})^T d = 0$. Then there is a $y \in \Omega$, $y \neq \bar{x}$, and a $\lambda_0 > 0$ such that $d = \lambda_0 (y - \bar{x})$. Let $\epsilon > 0$ be such that $f_0(\bar{x}) \leq f_0(x)$ for all $x \in \Omega$ with $\|x - \bar{x}\| \leq \epsilon$. Set $\bar{\lambda} = \min\{\lambda_0, \epsilon(\lambda_0 \|y - \bar{x}\|)^{-1}\} > 0$ so that $\bar{x} + \lambda d \in \Omega$ and $\|\bar{x} - (\bar{x} + \lambda d)\| \leq \epsilon$ for all $\lambda \in [0, \bar{\lambda}]$. By hypothesis, we now have

$$\begin{aligned} f_0(\bar{x}) &\leq f_0(\bar{x} + \lambda d) \\ &= f_0(\bar{x}) + \lambda \nabla f_0(\bar{x})^T (y - \bar{x}) + \frac{\lambda^2}{2} d^T \nabla^2 f_0(\bar{x}) d + o(\lambda^2) \\ &= f_0(\bar{x}) + \frac{\lambda^2}{2} d^T \nabla^2 f_0(\bar{x}) d + o(\lambda^2), \end{aligned}$$

where the second equality follows from the choice of d ($\nabla f_0(\bar{x})^T d = 0$). Therefore $d^T \nabla^2 f_0(\bar{x}) d \geq 0$.

(2) We show that $f_0(\bar{x}) \leq f_0(x) - \nu \|x - \bar{x}\|^2$ for some $\nu > 0$ for all $x \in \Omega$ near \bar{x} . Indeed, if this were not the case there would exist sequences $\{x_k\} \subset \Omega$, $\{\nu_k\} \subset \mathbb{R}_+$ with $x_k \rightarrow \bar{x}$, $\nu_k \downarrow 0$, and

$$f_0(x_k) < f_0(\bar{x}) + \nu_k \|x_k - \bar{x}\|^2$$

for all $k = 1, 2, \dots$ where, with no loss of generality, $\frac{x_k - \bar{x}}{\|x_k - \bar{x}\|} \rightarrow d$. Clearly, $d \in T_\Omega(\bar{x})$. Moreover,

$$\begin{aligned} f_0(\bar{x}) + \nabla f_0(\bar{x})^T (x_k - \bar{x}) &+ o(\|x_k - \bar{x}\|) \\ &= f_0(x_k) \\ &\leq f_0(\bar{x}) + \nu_k \|x_k - \bar{x}\|^2 \end{aligned}$$

so that $\nabla f_0(\bar{x})^T d = 0$.

Now, since $\nabla f_0(\bar{x})^T (x_k - \bar{x}) \geq 0$ for all $k = 1, 2, \dots$,

$$\begin{aligned} f_0(\bar{x}) + \frac{1}{2} (x_k - \bar{x})^T \nabla^2 f_0(\bar{x}) (x_k - \bar{x}) &+ o(\|x_k - \bar{x}\|^2) \\ &\leq f_0(\bar{x}) + \nabla f_0(\bar{x})^T (x_k - \bar{x}) + \frac{1}{2} (x_k - \bar{x})^T \nabla^2 f_0(\bar{x}) (x_k - \bar{x}) \\ &\quad + o(\|x_k - \bar{x}\|^2) \\ &= f_0(x_k) \\ &< f_0(\bar{x}) + \nu_k \|x_k - \bar{x}\|^2. \end{aligned}$$

Hence,

$$\left(\frac{x_k - \bar{x}}{\|x_k - \bar{x}\|} \right)^T \nabla^2 f_0(\bar{x}) \left(\frac{x_k - \bar{x}}{\|x_k - \bar{x}\|} \right) \leq \nu_k + \frac{o(\|x_k - \bar{x}\|^2)}{\|x_k - \bar{x}\|^2}$$

Taking the limit in k we obtain the contradiction

$$0 < d^T \nabla^2 f_0(\bar{x}) d \leq 0,$$

whereby the result is established. \square

Although it is possible to weaken the assumption of polyhedrality in Part 1, such weakenings are somewhat artificial as they essentially imply that $T_\Omega(x) = \bigcup_{\lambda \geq 0} (\Omega - x)$. The following example illustrates what can go wrong when the assumption of polyhedrality is dropped.

EXAMPLE 4.1. *Consider the problem*

$$\begin{aligned} \min \quad & \frac{1}{2}(x_2 - x_1^2) \\ \text{subject to} \quad & 0 \leq x_2, \quad x_1^3 \leq x_2^2. \end{aligned}$$

Observe that the constraint region in this problem can be written as $\Omega := \{(x_1, x_2)^T : |x_1|^{\frac{3}{2}} \leq x_2\}$, therefore

$$\begin{aligned} f_0(x) &= \frac{1}{2}(x_2 - x_1^2) \\ &\geq \frac{1}{2}(|x_1|^{\frac{3}{2}} - |x_1|^2) \\ &= \frac{1}{2}|x_1|^{\frac{3}{2}}(1 - |x_1|^{\frac{1}{2}}) > 0 \end{aligned}$$

whenever $0 < |x_1| \leq 1$. Consequently, the origin is a strict local solution for this problem. Nonetheless,

$$T_\Omega(0) \cap [\nabla f_0(0)]^\perp = \{(\delta, 0)^T : \delta \in \mathbb{R}\},$$

while

$$\nabla^2 f_0(0) = \begin{bmatrix} -1 & 0 \\ 0 & 0 \end{bmatrix}.$$

That is, even though the origin is a strict local solution, the Hessian of f_0 is not positive semidefinite on $T_\Omega(0)$.

When using the second-order conditions given above, one needs to be careful about the relationship between the Hessian of f_0 and the set $K := T_\Omega(x) \cap [\nabla f_0(x)]^\perp$. In particular, the positive definiteness (or semidefiniteness) of the Hessian of f_0 on the cone K does not necessarily imply the positive definiteness (or semidefiniteness) of the Hessian of f_0 on the subspace spanned by K . This is illustrated by the following example.

EXAMPLE 4.2. *Consider the problem*

$$\begin{aligned} \min \quad & (x_1^2 - \frac{1}{2}x_2^2) \\ \text{subject to} \quad & -x_1 \leq x_2 \leq x_1. \end{aligned}$$

Clearly, the origin is the unique global solution for this problem. Moreover, the constraint region for this problem, Ω , satisfies

$$T_\Omega(0) \cap [\nabla f(0)]^\perp = T_\Omega(0) = \Omega,$$

with the span of Ω being all of \mathbb{R}^2 . Now, while the Hessian of f_0 is positive definite on Ω , it is not positive definite on all of \mathbb{R}^2 .

In the polyhedral case it is easy to see that the sufficiency result in Theorem 4.5 is equivalent to the sufficiency result of Theorem 3.1. However, in the nonpolyhedral case, these results are not comparable. It is easy to see that Theorem 4.5 can handle situations where Theorem 3.1 does not apply even if Ω is given in the form (76). Just let one of the active constraint functions be nondifferentiable at the solution. Similarly, Theorem 3.1 can provide information when Theorem 4.5 does not. This is illustrated by the following example.

EXAMPLE 4.3. *Consider the problem*

$$\begin{aligned} \min \quad & x_2 \\ \text{subject to} \quad & x_1^2 \leq x_2. \end{aligned}$$

Clearly, $\bar{x} = 0$ is the unique global solution to this convex program. Moreover,

$$\begin{aligned} f_0(\bar{x}) + \frac{1}{2} \|x - \bar{x}\|^2 &= \frac{1}{2}(x_1^2 + x_2^2) \\ &\leq \frac{1}{2}(x_2 + x_2^2) \\ &\leq x_2 = f_0(x) \end{aligned}$$

for all x in the constraint region Ω with $\|x - \bar{x}\| \leq 1$. It is easily verified that this growth property is predicted by Theorem 4.5.

5. Convex Optimization, Saddle Point Theory, and Lagrangian Duality

In this section we extend the duality theory for linear programming to general problems of convex optimization. This is accomplished using the saddle point properties of the Lagrangian in convex optimization. Again, consider the problem

$$\begin{aligned} \mathcal{P} \quad & \text{minimize} && f_0(x) \\ & \text{subject to} && f_i \leq 0, \quad i = 1, 2, \dots, s \\ & && f_i(x) = 0, \quad i = s + 1, \dots, m, \end{aligned}$$

where it is assumed that the functions f_0, f_1, \dots, f_s are convex functions mapping \mathbb{R}^n to $\overline{\mathbb{R}}$, and f_{s+1}, \dots, f_m are affine mappings from \mathbb{R}^n to \mathbb{R} . We denote the constraint region for \mathcal{P} by Ω .

The Lagrangian for \mathcal{P} is the function

$$L(x, y) = f_0(x) + y_1 f_1(x) + y_2 f_2(x) + \dots + y_m f_m(x),$$

where it is always assumed that $0 \leq y_i$, $i = 1, 2, \dots, s$. Set $K = \mathbb{R}_+^s \times \mathbb{R}^{m-s} \subset \mathbb{R}^m$. A pair $(\bar{x}, \bar{y}) \in \mathbb{R}^n \times K$ is said to be a saddle point for L if

$$L(\bar{x}, y) \leq L(\bar{x}, \bar{y}) \leq L(x, \bar{y}) \quad \forall (x, y) \in \mathbb{R}^n \times K.$$

We have the following basic saddle point theorem for L .

THEOREM 5.1 (SADDLE POINT THEOREM). *Let $\bar{x} \in \mathbb{R}^n$. If there exists $\bar{y} \in K$ such that (\bar{x}, \bar{y}) is a saddle point for the Lagrangian L , then \bar{x} solves \mathcal{P} . Conversely, if \bar{x} is a solution to \mathcal{P} at which the Slater C.Q. is satisfied, then there is a $\bar{y} \in K$ such that (\bar{x}, \bar{y}) is a saddle point for L .*

PROOF. If $(\bar{x}, \bar{y}) \in \mathbb{R}^n \times K$ is a saddle point for \mathcal{P} then

$$\sup_{y \in K} L(\bar{x}, y) = \sup_{y \in K} f_0(\bar{x}) + y_1 f_1(\bar{x}) + y_2 f_2(\bar{x}) + \dots + y_m f_m(\bar{x}) \leq L(\bar{x}, \bar{y}).$$

If for some $i \in \{1, \dots, s\}$ such that $f_i(\bar{x}) > 0$, then we could send $y_i \uparrow +\infty$ to find that the supremum on the left is $+\infty$ which is a contradiction, so we must have $f_i(\bar{x}) \leq 0$, $i = 1, \dots, s$. Moreover, if $f_i(\bar{x}) \neq 0$ for some $i \in \{s + 1, \dots, m\}$, then we could send $y_i \uparrow -\text{sign}(f_i(\bar{x}))\infty$ to again find that the supremum on the left is $+\infty$ again a contradiction, so we must have $f_i(\bar{x}) = 0$, $i = s + 1, \dots, m$. That is, we must have $\bar{x} \in \Omega$. Since $L(\bar{x}, \bar{y}) = \sup_{y \in K} L(\bar{x}, y)$, we must have $\sum_{i=1}^m \bar{y}_i f_i(\bar{x}) = 0$. Therefore the right half of the saddle point condition implies that

$$f_0(\bar{x}) = L(\bar{x}, \bar{y}) \leq \inf_x L(x, \bar{y}) \leq \inf_{x \in \Omega} L(x, \bar{y}) \leq \inf_{x \in \Omega} f_0(x) \leq f_0(\bar{x}),$$

and so \bar{x} solves \mathcal{P} .

Conversely, if \bar{x} is a solution to \mathcal{P} at which the Slater C.Q. is satisfied, then there is a vector \bar{y} such that (\bar{x}, \bar{y}) is a KKT pair for \mathcal{P} . Primal feasibility ($\bar{x} \in \Omega$), dual feasibility ($\bar{y} \in K$), and complementarity ($\bar{y}_i f_i(\bar{x}) = 0$, $i = 1, \dots, s$) imply that

$$L(\bar{x}, y) \leq f_0(\bar{x}) = L(\bar{x}, \bar{y}) \quad \forall y \in K.$$

On the other hand, dual feasibility and convexity imply the convexity of the function $L(x, \bar{y})$ in x . Hence the condition $0 = \nabla_x L(\bar{x}, \bar{y})$ implies that \bar{x} is a global minimizer for the function $x \rightarrow L(x, \bar{y})$, that is

$$L(\bar{x}, \bar{y}) \leq L(x, \bar{y}) \quad \forall x \in \mathbb{R}^n.$$

Therefore, (\bar{x}, \bar{y}) is a saddle point for L . □

Note that it is always the case that

$$\sup_{y \in K} \inf_{x \in \mathbb{R}^n} L(x, y) \leq \inf_{x \in \mathbb{R}^n} \sup_{y \in K} L(x, y)$$

since the largest minimum is always smaller than the smallest maximum. On the other hand, if (\bar{x}, \bar{y}) is a saddle point for L , then

$$\inf_{x \in \mathbb{R}^n} \sup_{y \in K} L(x, y) \leq \sup_{y \in K} L(\bar{x}, y) \leq L(\bar{x}, \bar{y}) \leq \inf_{x \in \mathbb{R}^n} L(x, \bar{y}) \leq \sup_{y \in K} \inf_{x \in \mathbb{R}^n} L(x, y).$$

Hence, if a saddle point for L exists on $\mathbb{R}^n \times K$, then

$$\sup_{y \in K} \inf_{x \in \mathbb{R}^n} L(x, y) = \inf_{x \in \mathbb{R}^n} \sup_{y \in K} L(x, y).$$

Such a result is called a mini-max theorem and provides conditions under which one can exchange and inf-sup for a sup-inf. This mini-max result can be used as a basis for convex duality theory.

Observe that we have already shown that

$$\sup_{y \in K} L(x, y) = \begin{cases} +\infty & \text{if } x \notin \Omega, \\ f_0(x) & \text{if } x \in \Omega. \end{cases}$$

Therefore,

$$\inf_{x \in \mathbb{R}^n} \sup_{y \in K} L(x, y) = \inf_{x \in \Omega} f_0(x).$$

We will call this the *primal* problem. This is the inf-sup side of the saddle point problem. The other side, the sup-inf problem, we will call the *dual* problem with dual objective function

$$g(y) = \inf_{x \in \mathbb{R}^n} L(x, y).$$

The Saddle Point Theorem says that if (\bar{x}, \bar{y}) is a saddle point for L , then \bar{x} solves the primal problem, \bar{y} solves the dual problem, and the optimal values in the primal and dual problems coincide. This is a *Weak Duality Theorem*. The Strong Duality Theorem follows from the second half of the Saddle Point Theorem and requires the use of the Slater Constraint Qualification.

5.1. Linear Programming Duality. We now show how the Lagrangian Duality Theory described above gives linear programming duality as a special case. Consider the following LP:

$$\begin{aligned} \mathcal{P} \quad & \text{minimize} && b^T x \\ & \text{subject to} && A^T x \geq c, \quad 0 \leq x. \end{aligned}$$

The Lagrangian is

$$L(x, y, v) = b^T x + y^T (c - A^T x) - v^T x, \quad \text{where } 0 \leq y, \quad 0 \leq v.$$

The dual objective function is

$$g(y, u) = \min_{x \in \mathbb{R}^n} L(x, y, v) = \min_{x \in \mathbb{R}^n} b^T x + y^T (c - A^T x) - v^T x.$$

Our first goal is to obtain a closed form expression for $g(y, u)$. This is accomplished by using the optimality conditions for minimizing $L(x, y, u)$ to eliminate x from the definition of L . Since $L(x, y, v)$ is a convex function in x , the global solution to $\min_{x \in \mathbb{R}^n} L(x, y, v)$ is obtained by solving the equation $0 = \nabla_x L(x, y, u) = b - Ay - v$ with $0 \leq y, \quad 0 \leq v$. Using this condition in the definition of L we get

$$L(x, y, u) = b^T x + y^T (c - A^T x) - v^T x = (b - Ay - v)^T x + c^T y = c^T y,$$

subject to $b - A^T y = v$ and $0 \leq y, \quad 0 \leq v$. Hence the Lagrangian dual problem

$$\begin{aligned} & \text{maximize} && g(y, v) \\ & \text{subject to} && 0 \leq y, \quad 0 \leq v \end{aligned}$$

can be written as

$$\begin{aligned} \mathcal{D} \quad & \text{maximize} && c^T y \\ & \text{subject to} && b - Ay = v, \quad 0 \leq y, \quad 0 \leq v. \end{aligned}$$

Note that we can treat the variable v as a slack variable in this LP and write

$$\begin{aligned} \mathcal{D} \quad & \text{maximize} && c^T y \\ & \text{subject to} && Ay \leq b, \quad 0 \leq y. \end{aligned}$$

The linear program \mathcal{D} is the dual to the linear program \mathcal{P} .

5.2. Convex Quadratic Programming Duality. One can also apply the Lagrangian Duality Theory in the context of Convex Quadratic Programming. To see how this is done let $Q \in \mathbb{R}^{n \times n}$ be symmetric and positive definite, and let $c \in \mathbb{R}^n$. Consider the convex quadratic program

$$\begin{aligned} \mathcal{D} \quad & \text{minimize} && \frac{1}{2}x^T Qx + c^T x \\ & \text{subject to} && Ax \leq b, \quad 0 \leq x. \end{aligned}$$

The Lagrangian is given by

$$L(x, y, v) = \frac{1}{2}x^T Qx + c^T x + y^T (A^T x - b) - v^T x \quad \text{where } 0 \leq y, \quad 0 \leq v.$$

The dual objective function is

$$g(y, v) = \min_{x \in \mathbb{R}^n} L(x, y, v).$$

The goal is to obtain a closed form expression for g with the variable x removed by using the first-order optimality condition $0 = \nabla_x L(x, y, v)$. This optimality condition completely identifies the solution since L is convex in x . We have

$$0 = \nabla_x L(x, y, v) = Qx + c + A^T y - v.$$

Since Q is invertible, we have

$$x = Q^{-1}(v - A^T y - c).$$

Plugging this expression for x into $L(x, y, v)$ gives

$$\begin{aligned} g(y, v) &= L(Q^{-1}(v - A^T y - c), y, v) \\ &= \frac{1}{2}(v - A^T y - c)^T Q^{-1}(v - A^T y - c) \\ &\quad + c^T Q^{-1}(v - A^T y - c) + y^T (AQ^{-1}(v - A^T y - c) - b) - v^T Q^{-1}(v - A^T y - c) \\ &= \frac{1}{2}(v - A^T y - c)^T Q^{-1}(v - A^T y - c) - (v - A^T y - c)^T Q^{-1}(v - A^T y - c) - b^T y \\ &= -\frac{1}{2}(v - A^T y - c)^T Q^{-1}(v - A^T y - c) - b^T y. \end{aligned}$$

Hence the dual problem is

$$\begin{aligned} & \text{maximize} && -\frac{1}{2}(v - A^T y - c)^T Q^{-1}(v - A^T y - c) - b^T y \\ & \text{subject to} && 0 \leq y, \quad 0 \leq v. \end{aligned}$$

Moreover, (\bar{y}, \bar{v}) solve the dual problem if and only if $\bar{x} = Q^{-1}(\bar{v} - A^T \bar{y} - c)$ solves the primal problem with the primal and dual optimal values coinciding.

Exercises

- (1) Locate all of the KKT points for the following problems. Can you show that these points are local solutions? Global solutions?

(a)

$$\begin{aligned} & \text{minimize} && e^{(x_1-x_2)} \\ & \text{subject to} && e^{x_1} + e^{x_2} \leq 20 \\ & && 0 \leq x_1 \end{aligned}$$

(b)

$$\begin{aligned} & \text{minimize} && e^{(-x_1+x_2)} \\ & \text{subject to} && e^{x_1} + e^{x_2} \leq 20 \\ & && 0 \leq x_1 \end{aligned}$$

(c)

$$\begin{aligned} & \text{minimize} && x_1^2 + x_2^2 - 4x_1 - 4x_2 \\ & \text{subject to} && x_1^2 \leq x_2 \\ & && x_1 + x_2 \leq 2 \end{aligned}$$

(d)

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \|x\|^2 \\ & \text{subject to} && Ax = b \end{aligned}$$

where $b \in \mathbb{R}^m$ and $A \in \mathbb{R}^{m \times n}$ satisfies $\text{Nul}(A^T) = \{0\}$.

- (2) Show that the set

$$\Omega := \{x \in \mathbb{R}^2 \mid -x_1^3 \leq x_2 \leq x_1^3\}$$

is not regular at the origin. Graph the set Ω .

- (3) Construct an example of a constraint region of the form (76) at which the MFCQ is satisfied, but the LI condition is not satisfied.

- (4) Suppose $\Omega = \{x; Ax \leq b, Ex = h\}$ where $A \in \mathbb{R}^{m \times n}$, $E \in \mathbb{R}^{k \times n}$, $b \in \mathbb{R}^m$, and $h \in \mathbb{R}^k$.

(a) Given $x \in \Omega$, show that

$$T(x|\Omega) = \{d : A_i d \leq 0 \text{ for } i \in I(x), Ed = 0\},$$

where A_i denotes the i th row of the matrix A and $I(x) = \{i \mid A_i x = b_i\}$.

(b) Given $x \in \Omega$, show that every $d \in T(x|\Omega)$ is a feasible direction for Ω at x .

(c) Note that parts (a) and (b) above show that

$$T(x|\Omega) = \bigcup_{\lambda > 0} \lambda(\Omega - x)$$

whenever Ω is a convex polyhedral set. Why?

- (5) Let $C \subset \mathbb{R}^n$ be non-empty, closed and convex. For any $x \in \mathbb{R}^n$ consider the problem of finding the closest point in C to x using the 2-norm:

$$\mathcal{D} \quad \begin{aligned} & \text{minimize} && \frac{1}{2} \|x - z\|_2^2 \\ & \text{subject to} && z \in C \end{aligned}$$

Show that $\bar{z} \in C$ solves this problem if and only if

$$\langle x - \bar{z}, z - \bar{z} \rangle \leq 0 \quad \text{for all } z \in C.$$

- (6) Let Ω be a non-empty closed convex subset of \mathbb{R}^n . The geometric object *dual* to the tangent cone is called the *normal cone*:

$$N(x|\Omega) = \{z; \langle z, d \rangle \leq 0, \text{ for all } d \in T(x|\Omega)\}.$$

(a) Show that if \bar{x} solves the problem $\min\{f(x) : x \in \Omega\}$ then

$$-\nabla f(\bar{x}) \in N(\bar{x}|\Omega).$$

(b) Show that

$$N(\bar{x}|\Omega) = \{z : \langle z, x - \bar{x} \rangle \leq 0, \text{ for all } x \in \Omega\}.$$

(c) Let $\bar{x} \in \Omega$. Show that \bar{x} solves the problem $\min\{\frac{1}{2} \|x - y\|_2^2 : x \in \Omega\}$ for every $y \in \bar{x} + N(\bar{x}|\Omega)$.

(7) Consider the functions

$$f(x) = \frac{1}{2}x^T Qx - c^T x$$

and

$$f_t(x) = \frac{1}{2}x^T Qx - c^T x + t\phi(x),$$

where $t > 0$, $Q \in \mathbb{R}^{n \times n}$ is positive semi-definite, $c \in \mathbb{R}^n$, and $\phi : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ is given by

$$\phi(x) = \begin{cases} -\sum_{i=1}^n \ln x_i & , \text{ if } x_i > 0, i = 1, 2, \dots, n, \\ +\infty & , \text{ otherwise.} \end{cases}$$

- (a) Show that ϕ is a convex function.
- (b) Show that both f and f_t are convex functions.
- (c) Show that the solution to the problem $\min f_t(x)$ always exists and is unique.
- (d) Let $\{t_i\}$ be a decreasing sequence of positive real scalars with $t_i \downarrow 0$, and let x^i be the solution to the problem $\min f_{t_i}(x)$. Show that if the sequence $\{x^i\}$ has a cluster point \bar{x} , then \bar{x} must be a solution to the problem $\min\{f(x) : 0 \leq x\}$.

Hint: Use the KKT conditions for the QP $\min\{f(x) : 0 \leq x\}$.

Line Search Methods

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be given and suppose that x_c is our current best estimate of a solution to

$$\mathcal{P} \quad \min_{x \in \mathbb{R}^n} f(x) .$$

A standard method for improving the estimate x_c is to choose a direction of search $d \in \mathbb{R}^n$ and the compute a step length $t^* \in \mathbb{R}$ so that $x_c + t^*d$ approximately optimizes f along the line $\{x + td \mid t \in \mathbb{R}\}$. The new estimate for the solution to \mathcal{P} is then $x_n = x_c + t^*d$. The procedure for choosing t^* is called a *line search method*. If t^* is taken to be the global solution to the problem

$$\min_{t \in \mathbb{R}} f(x_c + td) ,$$

then t^* is called the *Curry* step length. However, except in certain very special cases, the Curry step length is far too costly to compute. For this reason we focus on a few easily computed step lengths. We begin the simplest and the most commonly used line search method called backtracking.

1. The Basic Backtracking Algorithm

In the backtracking line search we assume that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is differentiable and that we are given a direction d of strict descent at the current point x_c , that is $f'(x_c; d) < 0$.

INITIALIZATION: Choose $\gamma \in (0, 1)$ and $c \in (0, 1)$.

Having x_c obtain x_n as follows:

STEP 1: Compute the backtracking stepsize

$$\begin{aligned} t^* &:= \max \gamma^\nu \\ &\text{s.t. } \nu \in \{0, 1, 2, \dots\} \text{ and} \\ &f(x_c + \gamma^\nu d) \leq f(x_c) + c\gamma^\nu f'(x_c; d). \end{aligned}$$

STEP 2: Set $x_n = x_c + t^*d$.

The backtracking line search method forms the basic structure upon which most line search methods are built. Due to the importance of this method, we take a moment to emphasize its key features.

(1) The update to x_c has the form

$$(86) \quad x_n = x_c + t^*d .$$

Here d is called the *search direction* while t^* is called the *step length* or *stepsize*.

(2) The search direction d must satisfy

$$f'(x_c; d) < 0 .$$

Any direction satisfying this strict inequality is called a *direction of strict descent* for f at x_c . If $\nabla f(x_c) \neq 0$, then a direction of strict descent always exists. Just take $d = -\nabla f(x_c)$. As we have already seen

$$f'(x_c; -\nabla f(x_c)) = -\|\nabla f(x_c)\|^2 .$$

It is important to note that if d is a direction of strict descent for f at x_c , then there is a $\bar{t} > 0$ such that

$$f(x_c + td) < f(x_c) \quad \forall t \in (0, \bar{t}) .$$

In order to see this recall that

$$f'(x_c; d) = \lim_{t \downarrow 0} \frac{f(x_c + td) - f(x_c)}{t} .$$

Hence, if $f'(x_c; d) < 0$, there is a $\bar{t} > 0$ such that

$$\frac{f(x_c + td) - f(x_c)}{t} < 0 \quad \forall t \in (0, \bar{t}),$$

that is

$$f(x_c + td) < f(x_c) \quad \forall t \in (0, \bar{t}).$$

(3) In Step 1 of the algorithm, we require that the step length t^* be chosen so that

$$(87) \quad f(x_c + t^*d) \leq f(x_c) + c\gamma^\nu f'(x_c; d).$$

This inequality is called the Armijo-Goldstein inequality. It is named after the two researchers to first use it in the design of line search routines (Allen Goldstein is a Professor Emeritus here at the University of Washington). Observe that this inequality guarantees that

$$f(x_c + t^*d) < f(x_c).$$

For this reason, the algorithm described above is called a *descent algorithm*. It was observed in point (2) above that it is always possible to choose t^* so that $f(x_c + t^*d) < f(x_c)$. But the Armijo-Goldstein inequality is a somewhat stronger statement. To see that it too can be satisfied observe that since $f'(x_c; d) < 0$,

$$\lim_{t \downarrow 0} \frac{f(x_c + td) - f(x_c)}{t} = f'(x_c; d) < cf'(x_c; d) < 0.$$

Hence, there is a $\bar{t} > 0$ such that

$$\frac{f(x_c + td) - f(x_c)}{t} \leq cf'(x_c; d) \quad \forall t \in (0, \bar{t}),$$

that is

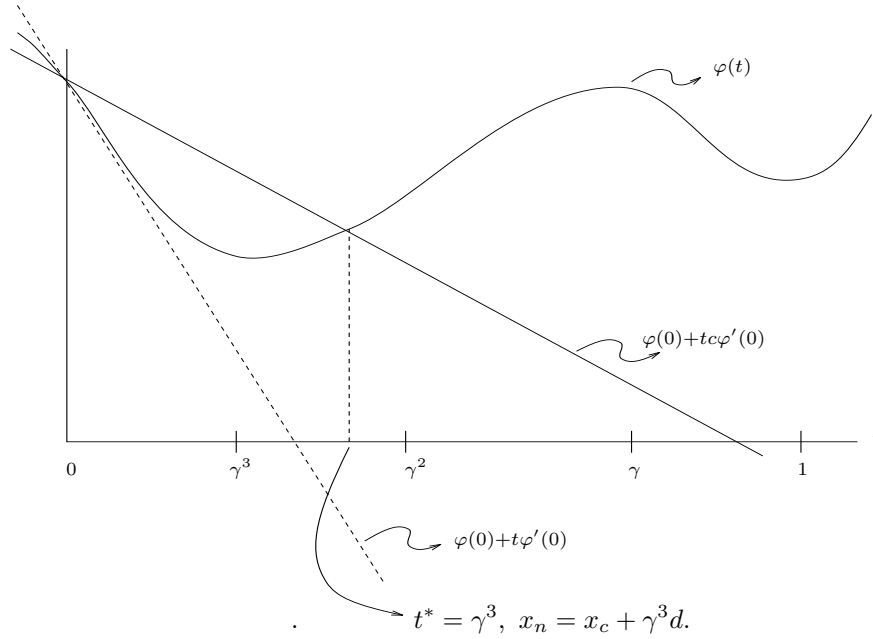
$$f(x_c + td) \leq f(x_c) + tc f'(x_c; d) \quad \forall t \in (0, \bar{t}).$$

- (4) The Armijo-Goldstein inequality is known as a condition of *sufficient decrease*. It is essential that we do not choose t^* too small. This is the reason for setting t^* equal to the first (largest) member of the geometric sequence $\{\gamma^\nu\}$ for which the Armijo-Goldstein inequality is satisfied. In general, we always wish to choose t^* as large as possible since it is often the case that some effort was put into the selection of the search direction d . Indeed, as we will see, for Newton's method we must take $t^* = 1$ in order to achieve rapid local convergence.
- (5) There is a balance that must be struck between taking t^* as large as possible and not having to evaluating the function at many points. Such a balance is obtained with an appropriate selection of the parameters γ and c . Typically one takes $\gamma \in [.5, .8]$ while $c \in [.001, .1]$ with adjustments depending on the cost of function evaluation and degree of nonlinearity.
- (6) The backtracking procedure of Step 1 is easy to program. A pseudo-Matlab code follows:

$$\left[\begin{array}{l} f_c = f(x_c) \\ \Delta f = cf'(x_c; d) \\ \text{new } f = f(x_c + d) \\ t = 1 \\ \text{while } \text{new } f > f_c + t\Delta f \\ \quad t = \gamma t \\ \quad \text{new } f = f(x_c + td) \\ \text{endwhile} \end{array} \right.$$

Point (3) above guarantees that this procedure is finitely terminating.

- (7) The backtracking procedure has a nice graphical illustration. Set $\varphi(t) = f(x_c + td)$ so that $\varphi'(0) = f'(x_c; d)$.



Before proceeding to a convergence result for the backtracking algorithm, we consider some possible choices for the search directions d . There are essentially three directions of interest:

(1) Steepest Descent (or Cauchy Direction):

$$d = -\nabla f(x_c) / \|\nabla f(x_c)\| .$$

(2) Newton Direction:

$$d = -\nabla^2 f(x_c)^{-1} \nabla f(x_c) .$$

(3) Newton-Like Direction:

$$d = -H \nabla f(x_c),$$

where $H \in \mathbb{R}^{n \times n}$ is symmetric and constructed to approximate the inverse of $\nabla^2 f(x_c)$.

In order to base a descent method on these directions we must have

$$f'(x_c; d) < 0.$$

For the Cauchy direction $-\nabla f(x_c) / \|\nabla f(x_c)\|$, this inequality always holds when $\nabla f(x_c) \neq 0$;

$$f'(x_c; -\nabla f(x_c) / \|\nabla f(x_c)\|) = -\|\nabla f(x_c)\| < 0.$$

On the other hand the Newton and Newton-like directions do not always satisfy this property:

$$f'(x_c; -H \nabla f(x_c)) = -\nabla f(x_c)^T H \nabla f(x_c).$$

These directions are directions of strict descent if and only if

$$0 < \nabla f(x_c)^T H \nabla f(x_c) .$$

This condition is related to second-order sufficiency conditions for optimality when H is an approximation to the inverse of the Hessian.

The advantage of the Cauchy direction is that it always provides a direction of strict descent. However, once the iterates get “close” to a stationary point, the procedure takes a very long time to obtain a moderately accurate estimate of the stationary point. Most often numerical error takes over due to very small stepsizes and the iterates behave chaotically.

On the other hand, Newton’s method (and its approximation, the secant method), may not define directions of strict descent until one is very close to a stationary point satisfying the second-order sufficiency condition. However, once one is near such a stationary point, then Newton’s method (and some Newton-Like methods) zoom in on the stationary point very rapidly. This behavior will be made precise when we establish our convergence result from Newton’s method.

Let us now consider the basic convergence result for the backtracking algorithm.

THEOREM 1.1. (CONVERGENCE FOR BACKTRACKING) *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $x_0 \in \mathbb{R}^n$ be such that f is differentiable on \mathbb{R}^n with ∇f Lipschitz continuous on an open convex set containing the set $\{x : f(x) \leq f(x_0)\}$. Let $\{x^k\}$ be the sequence satisfying $x^{k+1} = x^k$ if $\nabla f(x^k) = 0$; otherwise,*

$$x^{k+1} = x^k + t_k d^k, \quad \text{where } d^k \text{ satisfies } f'(x^k; d^k) < 0,$$

and t_k is chosen by the backtracking stepsize selection method. Then one of the following statements must be true:

- (i) *There is a k_0 such that $\nabla f'(x^{k_0}) = 0$.*
- (ii) *$f(x^k) \searrow -\infty$*
- (iii) *The sequence $\{\|d^k\|\}$ diverges ($\|d^k\| \rightarrow \infty$).*
- (iv) *For every subsequence $J \subset \mathbb{N}$ for which $\{d^k : k \in J\}$ is bounded, we have*

$$\lim_{k \in J} f'(x^k; d^k) = 0.$$

REMARK 1.1. *It is important to note that this theorem says nothing about the convergence of the sequence $\{x^k\}$. Indeed, this sequence may diverge. The theorem only concerns the function values and the first-order necessary condition for optimality.*

Before proving this Theorem, we first consider some important corollaries concerning the Cauchy and Newton search directions. Each corollary assumes that the hypotheses of Theorem 1.1 hold.

COROLLARY 1.1.1. *If the sequences $\{d^k\}$ and $\{f(x^k)\}$ are bounded, then*

$$\lim_{k \rightarrow \infty} f'(x^k; d^k) = 0.$$

PROOF. The hypotheses imply that either (i) or (iv) with $J = \mathbb{N}$ occurs in Theorem 1.1. Hence, $\lim_{k \rightarrow \infty} f'(x^k; d^k) = 0$. □

COROLLARY 1.1.2. *If $d^k = -\nabla f'(x^k) / \|\nabla f(x^k)\|$ is the Cauchy direction for all k , then every accumulation point, \bar{x} , of the sequence $\{x^k\}$ satisfies $\nabla f(\bar{x}) = 0$.*

PROOF. The sequence $\{f(x^k)\}$ is decreasing. If \bar{x} is any accumulation point of the sequence $\{x^k\}$, then we claim that $f(\bar{x})$ is a lower bound for the sequence $\{f(x^k)\}$. Indeed, if this were not the case, then for some k_0 and $\epsilon > 0$

$$f(x^k) + \epsilon < f(\bar{x})$$

for all $k > k_0$ since $\{f(x^k)\}$ is decreasing. But \bar{x} is a cluster point of $\{x^k\}$ and f is continuous. Hence, there is a $\hat{k} > k_0$ such that

$$|f(\bar{x}) - f(x^{\hat{k}})| < \epsilon/2.$$

But then

$$f(\bar{x}) < \frac{\epsilon}{2} + f(x^{\hat{k}}) \quad \text{and} \quad f(x^{\hat{k}}) + \epsilon < f(\bar{x}).$$

Hence,

$$f(x^{\hat{k}}) + \epsilon < \frac{\epsilon}{2} + f(x^{\hat{k}}), \quad \text{or} \quad \frac{\epsilon}{2} < 0.$$

This contradiction implies that $\{f(x^k)\}$ is bounded below by $f(\bar{x})$. But then the sequence $\{f(x^k)\}$ is bounded so that Corollary 1.1.1 applies. That is,

$$0 = \lim_{k \rightarrow \infty} f' \left(x^k; \frac{-\nabla f(x^k)}{\|\nabla f(x^k)\|} \right) = \lim_{k \rightarrow \infty} -\|\nabla f(x^k)\|.$$

Since ∇f is continuous, $\nabla f(\bar{x}) = 0$. □

COROLLARY 1.1.3. *Let us further assume that f is twice continuously differentiable and that there is a $\beta > 0$ such that, for all $u \in \mathbb{R}^n$, $\beta \|u\|^2 < u^T \nabla^2 f(x) u$ on $\{x : f(x) \leq f(x^0)\}$. If the Basic Backtracking algorithm is implemented using the Newton search directions,*

$$d^k = -\nabla^2 f(x^k)^{-1} \nabla f(x^k),$$

then every accumulation point, \bar{x} , of the sequence $\{x^k\}$ satisfies $\nabla f(\bar{x}) = 0$.

PROOF. Let \bar{x} be an accumulation point of the sequence $\{x^k\}$ and let $J \subset \mathbb{N}$ be such that $x^k \xrightarrow{J} \bar{x}$. Clearly, $\{x^k : k \in J\}$ is bounded. Hence, the continuity of ∇f and $\nabla^2 f$, along with the Weierstrass Compactness Theorem, imply that the sets $\{\nabla f(x^k) : k \in J\}$ and $\{\nabla^2 f(x^k) : k \in J\}$ are also bounded. Let M_1 be a bound on the values $\{\|\nabla f(x^k)\| : k \in J\}$ and let M_2 be an upper bound on the values $\{\|\nabla^2 f(x^k)\| : k \in J\}$. Recall that by hypotheses $\beta \|u\|^2$ is a uniform lower bound on the values $\{u^T \nabla^2 f(x^k) u\}$ for every $u \in \mathbb{R}^n$. Take $u = d^k$ to obtain the bound

$$\beta \|d^k\|^2 \leq \nabla f(x^k)^T \nabla^2 f(x^k)^{-1} \nabla f(x^k) \leq \|d^k\| \|\nabla f(x^k)\|,$$

and so

$$\|d^k\| \leq \beta^{-1} M_1 \quad \forall k \in J.$$

Therefore, the sequence $\{d^k : k \in J\}$ is bounded. Moreover, as in the proof of Corollary 1.1.2, the sequence $\{f(x^k)\}$ is also bounded. On the other hand,

$$\|\nabla f(x^k)\| = \|\nabla^2 f(x^k) d^k\| \leq M_2 \|d^k\| \quad \forall k \in J.$$

Therefore,

$$M_2^{-1} \|\nabla f(x^k)\| \leq \|d^k\| \quad \forall k \in J.$$

Consequently, Theorem 1.1 Part (iv) implies that

$$\begin{aligned} 0 &= \lim_{k \in J} |f'(x^k; d^k)| \\ &= \lim_{k \in J} |\nabla f(x^k)^T \nabla^2 f(x^k)^{-1} \nabla f(x^k)| \\ &\geq \lim_{k \in J} \beta \|d^k\|^2 \\ &\geq \lim_{k \in J} \beta M_2^{-2} \|\nabla f(x^k)\|^2 \\ &= \beta M_2^{-2} \|\nabla f(\bar{x})\|^2. \end{aligned}$$

Therefore, $\nabla f(\bar{x}) = 0$. □

PROOF OF THEOREM 1.1: We assume that none of (i), (ii), (iii), and (iv) hold and establish a contradiction.

Since (i) does not occur, $\nabla f(x^k) \neq 0$ for all $k = 1, 2, \dots$. Since (ii) does not occur, the sequence $\{f(x^k)\}$ is bounded below. Since $\{f(x^k)\}$ is a bounded decreasing sequence in \mathbb{R} , we have $f(x^k) \searrow \bar{f}$ for some \bar{f} . In particular, $(f(x^{k+1}) - f(x^k)) \rightarrow 0$. Next, since (iii) and (iv) do not occur, there is a subsequence $J \subset \mathbb{N}$ and a vector \bar{d} such that $d^k \xrightarrow{J} \bar{d}$ and

$$\sup_{k \in J} f'(x^k; d^k) =: \beta < 0.$$

The Armijo-Goldstein inequality combined with the fact that $(f(x^{k+1}) - f(x^k)) \rightarrow 0$, imply that

$$t_k f'(x^k; d^k) \rightarrow 0.$$

Since $f'(x^k; d^k) \leq \beta < 0$ for $k \in J$, we must have $t_k \xrightarrow{J} 0$. With no loss in generality, we assume that $t_k < 1$ for all $k \in J$. Hence,

$$(88) \quad c\gamma^{-1} t_k f'(x^k; d^k) < f(x^k + t_k \gamma^{-1} d^k) - f(x^k)$$

for all $k \in J$ due to Step 1 of the line search and the fact that $\tau_k < 1$. By the Mean Value Theorem, there exists for each $k \in J$ a $\theta_k \in (0, 1)$ such that

$$f(x^k + t_k \gamma^{-1} d^k) - f(x^k) = t_k \gamma^{-1} f'(\hat{x}^k; d^k)$$

where

$$\begin{aligned}\widehat{x}^n &:= (1 - \theta_k)x^k + \theta_k(x^k + t_k\gamma^{-1}d^k) \\ &= x^k + \theta_k t_k \gamma^{-1} d^k.\end{aligned}$$

Now, since ∇f is Lipschitz continuous, we have

$$\begin{aligned}f(x^k + t_k\gamma^{-1}d^k) - f(x^k) &= t_k\gamma^{-1}f'(\widehat{x}^k; d^k) \\ &= t_k\gamma^{-1}f'(x^k; d^k) + t_k\gamma^{-1}[f'(\widehat{x}^k; d^k) - f'(x^k; d^k)] \\ &= t_k\gamma^{-1}f'(x^k; d^k) + t_k\gamma^{-1}[\nabla f(\widehat{x}^k) - \nabla f(x^k)]^T d^k \\ &\leq t_k\gamma^{-1}f'(x^k; d^k) + t_k\gamma^{-1}L\|\widehat{x}^k - x^k\|\|d^k\| \\ &= t_k\gamma^{-1}f'(x^k; d^k) + L(t_k\gamma^{-1})^2\theta_k\|d^k\|^2.\end{aligned}$$

Combining this inequality with inequality (88) yields the inequality

$$ct_k\gamma^{-1}f'(x^k; d^k) < t_k\gamma^{-1}f'(x^k; d^k) + L(t_k\gamma^{-1})^2\theta_k\|d^k\|^2.$$

By rearranging and then substituting β for $f'(x^k; d^k)$ we obtain

$$0 < (1 - c)\beta + (t_k\gamma^{-1})L\|\delta_k\|^2 \quad \forall k \in J.$$

Now taking the limit over $k \in J$, we obtain the contradiction

$$0 \leq (1 - c)\beta < 0.$$

□

2. The Wolfe Conditions

We now consider a couple of modifications to the basic backtracking line search that attempt to better approximate an exact line-search (Curry line search), i.e. the stepsize t_k is chosen to satisfy

$$f(x^k + t_k d^k) = \min_{t \in \mathbb{R}} f(x^k + t d^k).$$

In this case, the first-order optimality conditions tell us that $0 = \nabla f(x^k + t_k d^k)^T d^k$. The Wolfe conditions try to combine the Armijo-Goldstein sufficient decrease condition with a condition that tries to push $\nabla f(x^k + t_k d^k)^T d^k$ either toward zero, or at least to a point where the search direction d^k is less of a direction of descent. To describe these line search conditions, we take parameters $0 < c_1 < c_2 < 1$.

Weak Wolfe Conditions

$$(89) \quad f(x^k + t_k d^k) \leq f(x^k) + c_1 t_k f'(x^k; d^k)$$

$$(90) \quad c_2 f'(x^k; d^k) \leq f'(x^k + t_k d^k; d^k).$$

Strong Wolfe Conditions

$$(91) \quad f(x^k + t_k d^k) \leq f(x^k) + c_1 t_k f'(x^k; d^k)$$

$$(92) \quad |f'(x^k + t_k d^k; d^k)| \leq c_2 |f'(x^k; d^k)|.$$

The weak Wolfe condition (90) tries to make d^k less of a direction of descent (and possibly a direction of ascent) at the new point, while the strong Wolfe condition tries to push the directional derivative in the direction d^k closer to zero at the new point. Imposing one or the other of the Wolfe conditions on a line search procedure has become standard practice for optimization software based on line search methods.

We now give a result showing that there exists stepsizes satisfying the weak Wolfe conditions. A similar result (with a similar proof) holds for the strong Wolfe conditions.

LEMMA 2.1. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be continuously differentiable and suppose that $x, d \in \mathbb{R}^n$ are such that the set $\{f(x + td) : t \geq 0\}$ is bounded below and $f'(x; d) < 0$, then for each $0 < c_1 < c_2 < 1$ the set*

$$\left\{ t \mid \begin{array}{l} t > 0, f'(x + td; d) \geq c_2 f'(x; d), \text{ and} \\ f(x + td) \leq f(x) + c_1 t f'(x; d) \end{array} \right\}$$

has non-empty interior.

PROOF. Set $\phi(t) = f(x + td) - (f(x) + c_1 t f'(x; d))$. Then $\phi(0) = 0$ and $\phi'(0) = (1 - c_1)f'(x; d) < 0$. So there is a $\bar{t} > 0$ such that $\phi(t) < 0$ for $t \in (0, \bar{t})$. Moreover, since $f'(x; d) < 0$ and $\{f(x + td) : t \geq 0\}$ is bounded below, we have $\phi(t) \rightarrow +\infty$ as $t \uparrow \infty$. Hence, by the continuity of f , there exists $\hat{t} > 0$ such that $\phi(\hat{t}) = 0$. Let $t^* = \inf \{\hat{t} \mid 0 \leq t, \phi(\hat{t}) = 0\}$. Since $\phi(t) < 0$ for $t \in (0, \bar{t})$, $t^* > 0$ and by continuity $\phi(t^*) = 0$. By Rolle's theorem (or the mean value theorem) there must exist $\tilde{t} \in (0, t^*)$ with $\phi'(\tilde{t}) = 0$. That is,

$$\nabla f(x + \tilde{t}d)^T d = c_1 \nabla f(x)^T d > c_2 \nabla f(x)^T d.$$

From the definition of t^* and the fact that $\tilde{t} \in (0, t^*)$, we also have

$$f(x + td) - (f(x) + c_1 \tilde{t} \nabla f(x)^T d) < 0.$$

The result now follows from the continuity of f and ∇f . \square

We now describe a bisection method that either computes a stepsize satisfying the weak Wolfe conditions or sends the function values to $-\infty$. Let x and d in \mathbb{R}^n be such that $f'(x; d) < 0$.

A Bisection Method for the Weak Wolfe Conditions

INITIALIZATION: Choose $0 < c_1 < c_2 < 1$, and set $\alpha = 0$, $t = 1$, and $\beta = +\infty$.

REPEAT

If $f(x + td) > f(x) + c_1 t f'(x; d)$,
 set $\beta = t$ and reset $t = \frac{1}{2}(\alpha + \beta)$.
 Else if $f'(x + td; d) < c_2 f'(x; d)$,
 set $\alpha = t$ and reset

$$t = \begin{cases} 2\alpha, & \text{if } \beta = +\infty \\ \frac{1}{2}(\alpha + \beta), & \text{otherwise.} \end{cases}$$

Else, STOP.

END REPEAT

LEMMA 2.2. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be continuously differentiable and suppose that $x, d \in \mathbb{R}^n$ are such that $f'(x; d) < 0$. Then one of the following two possibilities must occur in the Bisection Method for the Weak Wolfe Condition described above.

- (i) The procedure terminates finitely at a value of t for which the weak Wolfe conditions are satisfied.
- (ii) The procedure does not terminate finitely, the parameter β is never set to a finite value, the parameter α becomes positive on the first iteration and is doubled in magnitude at every iteration thereafter, and $f(x + td) \downarrow -\infty$.

PROOF. Let us suppose that the procedure does not terminate finitely. If the parameter β is never set to a finite value, then it must be the case that α becomes positive on the first iteration (since we did not terminate) and is doubled on each subsequent iteration with

$$f(x + \alpha d) \leq f(x) + c_1 \alpha f'(x; d).$$

But then $f(x + td) \downarrow -\infty$ since $f'(x; d) < 0$. That is, option (ii) above occurs. Hence, we may as well assume that β is eventually finite and the procedure is not finitely terminating. For the sake of clarity, let us index the bounds and trial steps by iteration as follows: $\alpha_k < t_k < \beta_k$, $k = 1, 2, \dots$. Since β is eventually finite, the bisection procedure guarantees that there is a $\bar{t} > 0$ such that

$$(93) \quad \alpha_k \uparrow \bar{t}, \quad t_k \rightarrow \bar{t}, \quad \text{and} \quad \beta_k \downarrow \bar{t}.$$

If $\alpha_k = 0$ for all k , then $\bar{t} = 0$ and

$$\frac{f(x + t_k d) - f(x)}{t_k} - c_1 f'(x; d) > 0 \quad \forall k.$$

But then, taking the limit in k , we obtain $f'(x; d) \geq c_1 f'(x; d)$, or equivalently, $0 > (1 - c_1)f'(x; d) \geq 0$ which is a contradiction. Hence, we can assume that eventually $\alpha_k > 0$.

We now have that the sequences $\{\alpha_k\}$, $\{t_k\}$, and $\{\beta_k\}$ are infinite with (93) satisfied, and there is a k_0 such that $0 < \alpha_k < t_k < \beta_k < \infty$ for all $k \geq k_0$. By construction, we know that for all $k > k_0$

$$(94) \quad f(x + \alpha_k d) \leq f(x) + c_1 \alpha_k f'(x; d)$$

$$(95) \quad f(x) + c_1 \beta_k f'(x; d) < f(x + \beta_k d)$$

$$(96) \quad f'(x + \alpha_k d; d) < c_2 f'(x; d) .$$

Taking the limit in k in (96) tells us that

$$(97) \quad f'(x + \bar{t}d; d) \leq c_2 f'(x; d) .$$

Adding (94) and (95) together and using the Mean Value Theorem gives

$$c_1(\beta_k - \alpha_k)f'(x; d) \leq f(x + \beta_k d) - f(x + \alpha_k d) = (\beta_k - \alpha_k)f'(x + \hat{t}_k d; d) \quad \forall k > k_0,$$

where $\alpha_k \leq \hat{t}_k \leq \beta_k$. Dividing by $(\beta_k - \alpha_k) > 0$ and taking the limit in k gives $c_1 f'(x; d) \leq f'(x + \bar{t}d; d)$ which combined with (97) yields the contradiction $f'(x + \bar{t}d; d) \leq c_2 f'(x; d) < c_1 f'(x; d) \leq f'(x + \bar{t}d; d)$. Consequently, option (i) above must occur if (ii) does not. \square

A global convergence result for a line search routine based on the Weak Wolfe conditions now follows.

THEOREM 2.1. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $x^0 \in \mathbb{R}^n$, and $0 < c_1 < c_2 < 1$. Assume that $\nabla f(x)$ exists and is Lipschitz continuous on an open set containing the set $\{x \mid f(x) \leq f(x^0)\}$. Let $\{x^\nu\}$ be a sequence initiated at x^0 and generated by the following algorithm:*

Step 0: Set $k = 0$.

Step 1: Choose $d^k \in \mathbb{R}^n$ such that $f'(x^k; d^k) < 0$.

If no such d^k exists, then STOP.

First-order necessary conditions for optimality are satisfied at x^k .

Step 2: Let t^k be a stepsize satisfying the Weak Wolfe conditions (89) and (90).

If no such t^k exists, then STOP.

The function f is unbounded below.

Step 3: Set $x^{k+1} = x^k + t_k d^k$, reset $k = k + 1$, and return to Step 1.

One of the following must occur:

(i) The algorithm terminates finitely at a first-order stationary point for f .

(ii) For some k the stepsize selection procedure generates a sequence of trial stepsizes $t_{k\nu} \uparrow +\infty$ such that $f(x^k + t_{k\nu} d^k) \rightarrow -\infty$.

(iii) $f(x^k) \downarrow -\infty$.

(iv) $\sum_{k=0}^{\infty} \|\nabla f(x^k)\|^2 \cos^2 \theta_k < +\infty$, where $\cos \theta_k = \frac{\nabla f(x^k)^T d^k}{\|\nabla f(x^k)\| \|d^k\|}$ for all $k = 1, 2, \dots$

PROOF. We assume that (i), (ii), and (iii) do not occur and show that (iv) occurs. Since (i) and (ii) do not occur the sequence $\{x^\nu\}$ is infinite and $f'(x^k; d^k) < 0$ for all $k = 1, 2, \dots$. Since (ii) does not occur, the weak Wolfe conditions are satisfied at every iteration. The condition (89) implies that the sequence $\{f(x^k)\}$ is strictly decreasing. In particular, this implies that $\{x^\nu\} \subset \{x \mid f(x) \leq f(x^0)\}$. The condition (90) implies that

$$(c_2 - 1)\nabla f(x^k)^T d^k \leq (\nabla f(x^{k+1}) - \nabla f(x^k))^T d^k$$

for all k . Combining this with the Lipschitz continuity of ∇f on an open neighborhood of $\{x \mid f(x) \leq f(x^0)\}$, gives

$$(c_2 - 1)\nabla f(x^k)^T d^k \leq (\nabla f(x^{k+1}) - \nabla f(x^k))^T d^k \leq Lt_k \|d^k\|^2 .$$

Hence

$$t_k \geq \frac{c_2 - 1}{L} \frac{\nabla f(x^k)^T d^k}{\|d^k\|^2} > 0.$$

Plugging this into (89) give the inequality

$$f(x^{k+1}) \leq f(x^k) - c_1 \frac{1 - c_2}{L} \frac{(\nabla f(x^k)^T d^k)^2}{\|d^k\|^2} = f(x^k) - c_1 \frac{1 - c_2}{L} \|\nabla f(x^k)\|^2 \cos^2 \theta_k.$$

Setting $c = c_1 \frac{1-c_2}{L}$ and summing over k gives

$$f(x^{k+1}) \leq f(x^0) - c \sum_{\nu=0}^k \|\nabla f(x^\nu)\|^2 \cos^2 \theta_\nu .$$

Since (iii) does not occur, we can take the limit in k and obtain

$$\sum_{\nu=0}^{\infty} \|\nabla f(x^\nu)\|^2 \cos^2 \theta_\nu < +\infty .$$

□

If the function f is bounded below and the algorithm does not terminate finitely, then Part (iv) of this theorem states that

$$\|\nabla f(x^k)\| \cos^2 \theta_k \rightarrow 0 .$$

Hence, if the search directions d^k are chosen so that there is a $\delta > 0$, independent of the iteration k , such that $\cos \theta_k < -\delta$ for all k , then it must be the case that $\|\nabla f(x^k)\| \rightarrow 0$ so that every cluster point of the sequence $\{x^k\}$ is a first-order stationary point for f . For example, we have the following corollary to the theorem.

COROLLARY 2.1.1. *Let f and $\{x^k\}$ be as in the theorem, and let $\{B_k\}$ be a sequence of symmetric positive definite matrices for which there exists $\underline{\lambda} > \underline{\lambda} > 0$ such that*

$$(98) \quad \underline{\lambda} \|u\|^2 \leq u^T B_k u \leq \bar{\lambda} \|u\|^2 \quad \forall u \in \mathbb{R}^n \text{ and } k = 1, 2, \dots .$$

Let us further assume that f is bounded below. If the search directions d^k are given by

$$d^k = -B_k \nabla f(x^k) \quad \forall k = 1, 2, \dots ,$$

then $\nabla f(x^k) \rightarrow 0$.

PROOF. It is easily shown (see exercises) that the condition (98) implies that the eigenvalues of the sequence $\{B_k\}$ are uniformly lower bounded by $\underline{\lambda}$ and uniformly upper bounded by $\bar{\lambda}$. In particular, this implies that

$$\underline{\lambda} \|u\| \leq \|B_k u\| \leq \bar{\lambda} \|u\| \quad \forall u \in \mathbb{R}^n \text{ and } k = 1, 2, \dots$$

(see exercises). Hence for all k

$$\begin{aligned} \cos \theta_k &= \frac{\nabla f(x^k)^T d^k}{\|\nabla f(x^k)\| \|d^k\|} \\ &= -\frac{\nabla f(x^k)^T B_k \nabla f(x^k)}{\|\nabla f(x^k)\| \|B_k \nabla f(x^k)\|} \\ &\leq -\frac{\underline{\lambda} \|\nabla f(x^k)\|^2}{\|\nabla f(x^k)\| \|B_k \nabla f(x^k)\|} \\ &\leq -\frac{\underline{\lambda} \|\nabla f(x^k)\|^2}{\|\nabla f(x^k)\| \bar{\lambda} \|\nabla f(x^k)\|} \\ &= -\underline{\lambda} / \bar{\lambda} \\ &< 0 . \end{aligned}$$

Therefore $\nabla f(x^k) \rightarrow 0$. □

A possible choice for the matrices B_k in the above result is $B_k = I$ for all k . This essentially gives the method of steepest descent.

Search Directions for Unconstrained Optimization

In this chapter we study the choice of search directions used in our basic updating scheme

$$x^{k+1} = x^k + t_k d^k .$$

for solving

$$\mathcal{P} \quad \min_{x \in \mathbb{R}^n} f(x).$$

All of the search directions considered can be classified as *Newton-like* since they are all of the form

$$d^k = -H_k \nabla f(x^k),$$

where H_k is a symmetric $n \times n$ matrix. If $H_k = \mu_k I$ for all k , the resulting search directions are a scaled steepest descent direction with scale factors μ_k . More generally, we choose H_k to approximate $\nabla^2 f(x^k)^{-1}$ in order to approximate Newton's method for optimization. The Newton's method is important since it possesses rapid local convergence properties, and can be shown to be *scale independent*. We precede our discussion of search directions by making precise a useful notion of speed or *rate of convergence*.

1. Rate of Convergence

We focus on notions of quotient rates convergence, or Q-convergence rates. Let $\{x^\nu\} \subset \mathbb{R}^n$ and $\bar{x} \in \mathbb{R}^n$ be such that $\bar{x}^\nu \rightarrow \bar{x}$. We say that $\bar{x}^\nu \rightarrow \bar{x}$ at a *linear* rate if

$$\limsup_{\nu \rightarrow \infty} \frac{\|x^{\nu+1} - \bar{x}\|}{\|x^\nu - \bar{x}\|} < 1 .$$

The convergence is said to be *superlinear* if this limsup is 0. The convergence is said to be *quadratic* if

$$\limsup_{\nu \rightarrow \infty} \frac{\|x^{\nu+1} - \bar{x}\|}{\|x^\nu - \bar{x}\|^2} < \infty .$$

For example, given $\gamma \in (0, 1)$ the sequence $\{\gamma^\nu\}$ converges linearly to zero, but not superlinearly. The sequence $\{\gamma^{\nu^2}\}$ converges superlinearly to 0, but not quadratically. Finally, the sequence $\{\gamma^{\nu^2}\}$ converges quadratically to zero. Superlinear convergence is much faster than linear convergences, but quadratic convergence is much, much faster than superlinear convergence.

2. Newton's Method for Solving Equations

Newton's method is an iterative scheme designed to solve nonlinear equations of the form

$$(99) \quad g(x) = 0,$$

where $g : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is assumed to be continuously differentiable. Many problems of importance can be posed in this way. In the context of the optimization problem \mathcal{P} , we wish to locate critical points, that is, points at which $\nabla f(x) = 0$. We begin our discussion of Newton's method in the usual context of equation solving.

Assume that the function g in (99) is continuously differentiable and that we have an approximate solution $x^0 \in \mathbb{R}^n$. We now wish to improve on this approximation. If \bar{x} is a solution to (99), then

$$0 = g(\bar{x}) = g(x^0) + g'(x^0)(\bar{x} - x^0) + o\|\bar{x} - x^0\|.$$

Thus, if x^0 is "close" to \bar{x} , it is reasonable to suppose that the solution to the linearized system

$$(100) \quad 0 = g(x^0) + g'(x^0)(x - x^0)$$

is even closer. This is Newton's method for finding the roots of the equation $g(x) = 0$. It has one obvious pitfall. Equation (100) may not be consistent. That is, there may not exist a solution to (100).

For the sake of the present argument, we assume that (3) holds, i.e. $g'(x^0)^{-1}$ exists. Under this assumption (100) defines the iteration scheme,

$$(101) \quad x^{k+1} := x^k - [g'(x^k)]^{-1}g(x^k),$$

called the Newton iteration. The associated direction

$$(102) \quad d^k := -[g'(x^k)]^{-1}g(x^k).$$

is called the Newton direction. We analyze the convergence behavior of this scheme under the additional assumption that only an approximation to $g'(x^k)^{-1}$ is available. We denote this approximation by J_k . The resulting iteration scheme is

$$(103) \quad x^{k+1} := x^k - J_k g(x^k).$$

Methods of this type are called *Newton-Like methods*.

THEOREM 2.1. *Let $g : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be differentiable, $x^0 \in \mathbb{R}^n$, and $J_0 \in \mathbb{R}^{n \times n}$. Suppose that there exists \bar{x} , $x^0 \in \mathbb{R}^n$, and $\epsilon > 0$ with $\|x^0 - \bar{x}\| < \epsilon$ such that*

- (1) $g(\bar{x}) = 0$,
- (2) $g'(x)^{-1}$ exists for $x \in B(\bar{x}; \epsilon) := \{x \in \mathbb{R}^n : \|x - \bar{x}\| < \epsilon\}$ with

$$\sup\{\|g'(x)^{-1}\| : x \in B(\bar{x}; \epsilon)\} \leq M_1$$

- (3) g' is Lipschitz continuous on $\text{cl}B(\bar{x}; \epsilon)$ with Lipschitz constant L , and
- (4) $\theta_0 := \frac{LM_1}{2}\|x^0 - \bar{x}\| + M_0K < 1$ where $K \geq \|(g'(x^0)^{-1} - J_0)y^0\|$, $y^0 := g(x^0)/\|g(x^0)\|$, and $M_0 = \max\{\|g'(x)\| : x \in B(\bar{x}; \epsilon)\}$.

Further suppose that iteration (103) is initiated at x^0 where the J_k 's are chosen to satisfy one of the following conditions;

- (i) $\|(g'(x^k)^{-1} - J_k)y^k\| \leq K$,
- (ii) $\|(g'(x^k)^{-1} - J_k)y^k\| \leq \theta_1^k K$ for some $\theta_1 \in (0, 1)$,
- (iii) $\|(g'(x^k)^{-1} - J_k)y^k\| \leq \min\{M_3\|x^k - x^{k-1}\|, K\}$, for some $M_2 > 0$, or
- (iv) $\|(g'(x^k)^{-1} - J_k)y^k\| \leq \min\{M_2\|g(x^k)\|, K\}$, for some $M_3 > 0$,

where for each $k = 1, 2, \dots$, $y^k := g(x^k)/\|g(x^k)\|$.

These hypotheses on the accuracy of the approximations J_k yield the following conclusions about the rate of convergence of the iterates x^k .

- (a) If (i) holds, then $x^k \rightarrow \bar{x}$ linearly.
- (b) If (ii) holds, then $x^k \rightarrow \bar{x}$ superlinearly.
- (c) If (iii) holds, then $x^k \rightarrow \bar{x}$ two step quadratically.
- (d) If (iv) holds, then $x^k \rightarrow \bar{x}$ quadratically.

PROOF. We begin by inductively establishing the basic inequalities

$$(104) \quad \|x^{k+1} - \bar{x}\| \leq \frac{LM_1}{2}\|x^k - \bar{x}\|^2 + \|(g'(x^k)^{-1} - J_k)g(x^k)\|,$$

and

$$(105) \quad \|x^{k+1} - \bar{x}\| \leq \theta_0\|x^k - \bar{x}\|$$

as well as the inclusion

$$(106) \quad x^{k+1} \in B(\bar{x}; \epsilon)$$

for $k = 0, 1, 2, \dots$. For $k = 0$ we have

$$\begin{aligned} x^1 - \bar{x} &= x^0 - \bar{x} - g'(x^0)^{-1}g(x^0) + [g'(x^0)^{-1} - J_0]g(x^0) \\ &= g'(x^0)^{-1}[g(\bar{x}) - (g(x^0) + g'(x^0)(\bar{x} - x^0))] \\ &\quad + [g'(x^0)^{-1} - J_0]g(x^0), \end{aligned}$$

since $g'(x^0)^{-1}$ exists by the hypotheses. Consequently, the hypothese (1)–(4) plus the quadratic bound lemma imply that

$$\begin{aligned} \|x^{k+1} - \bar{x}\| &\leq \|g'(x^0)^{-1}\| \|g(\bar{x}) - (g(x^0) + g'(x^0)(\bar{x} - x^0))\| \\ &\quad + \|(g'(x^0)^{-1} - J_0)g(x^0)\| \\ &\leq \frac{M_1 L}{2} \|x^0 - \bar{x}\|^2 + K \|g(x^0) - g(\bar{x})\| \\ &\leq \frac{M_1 L}{2} \|x^0 - \bar{x}\|^2 + M_0 K \|x^0 - \bar{x}\| \\ &\leq \theta_0 \|x^0 - \bar{x}\| < \epsilon, \end{aligned}$$

whereby (104) – (106) are established for $k = 0$.

Next suppose that (104) – (106) hold for $k = 0, 1, \dots, s-1$. We show that (104) – (106) hold at $k = s$. Since $x^s \in B(\bar{x}, \epsilon)$, hypotheses (2)–(4) hold at x^s , one can proceed exactly as in the case $k = 0$ to obtain (104). Now if any one of (i)–(iv) holds, then (i) holds. Thus, by (104), we find that

$$\begin{aligned} \|x^{s+1} - \bar{x}\| &\leq \frac{M_1 L}{2} \|x^s - \bar{x}\|^2 + \|(g'(x^s)^{-1} - J_s)g(x^s)\| \\ &\leq \left[\frac{M_1 L}{2} \theta_0^s \|x^0 - \bar{x}\| + M_0 K \right] \|x^s - \bar{x}\| \\ &\leq \left[\frac{M_1 L}{2} \|x^0 - \bar{x}\| + M_0 K \right] \|x^s - \bar{x}\| \\ &= \theta_0 \|x^s - \bar{x}\|. \end{aligned}$$

Hence $\|x^{s+1} - \bar{x}\| \leq \theta_0 \|x^s - \bar{x}\| \leq \theta_0 \epsilon < \epsilon$ and so $x^{s+1} \in B(\bar{x}, \epsilon)$. We now proceed to establish (a)–(d).

(a) This clearly holds since the induction above established that

$$\|x^{k+1} - \bar{x}\| \leq \theta_0 \|x^k - \bar{x}\|.$$

(b) From (104), we have

$$\begin{aligned} \|x^{k+1} - \bar{x}\| &\leq \frac{LM_1}{2} \|x^k - \bar{x}\|^2 + \|(g'(x^k)^{-1} - J_k)g(x^k)\| \\ &\leq \frac{LM_1}{2} \|x^k - \bar{x}\|^2 + \theta_1^k K \|g(x^k)\| \\ &\leq \left[\frac{LM_1}{2} \theta_0^k \|x^0 - \bar{x}\| + \theta_1^k M_0 K \right] \|x^k - \bar{x}\| \end{aligned}$$

Hence $x^k \rightarrow \bar{x}$ superlinearly.

(c) From (104) and the fact that $x^k \rightarrow \bar{x}$, we eventually have

$$\begin{aligned} \|x^{k+1} - \bar{x}\| &\leq \frac{LM_1}{2} \|x^k - \bar{x}\|^2 + \|(g'(x^k)^{-1} - J_k)g(x^k)\| \\ &\leq \frac{LM_1}{2} \|x^k - \bar{x}\|^2 + M_2 \|x^k - x^{k-1}\| \|g(x^k)\| \\ &\leq \left[\frac{LM_1}{2} \|x^k - \bar{x}\| + M_0 M_2 [\|x^{k-1} - \bar{x}\| + \|x^k - \bar{x}\|] \right] \|x^k - \bar{x}\| \\ &\leq \left[\frac{LM_1}{2} \theta_0 \|x^{k-1} - \bar{x}\| + M_0 M_2 (1 + \theta_0) \|x^{k-1} - \bar{x}\| \right] \\ &\quad \times \theta_0 \|x^{k-1} - \bar{x}\| \\ &= \left[\frac{LM_1}{2} \theta_0 + M_0 M_2 (1 + \theta_0) \right] \theta_0 \|x^{k-1} - \bar{x}\|^2. \end{aligned}$$

Hence $x^k \rightarrow \bar{x}$ two step quadratically.

(d) Again by (104) and the fact that $x^k \rightarrow \bar{x}$, we eventually have

$$\begin{aligned} \|x^{k+1} - \bar{x}\| &\leq \frac{LM_1}{2} \|x^k - \bar{x}\|^2 + \|(g'(x^k)^{-1} - J_k)g(x^k)\| \\ &\leq \frac{LM_1}{2} \|x^k - \bar{x}\|^2 + M_2 \|g(x^k)\|^2 \\ &\leq \left[\frac{LM_1}{2} + M_2 M_0^2 \right] \|x^k - \bar{x}\|^2 . \end{aligned}$$

□

Note that the conditions required for the approximations to the Jacobian matrices $g'(x^k)^{-1}$ given in (i)–(ii) do not imply that $J_k \rightarrow g'(\bar{x})^{-1}$. The stronger conditions

- (i)' $\|g'(x^k)^{-1} - J_k\| \leq \|g'(x^0)^{-1} - J_0\|$,
- (ii)' $\|g'(x^{k+1})^{-1} - J_{k+1}\| \leq \theta_1 \|g'(x^k)^{-1} - J_k\|$ for some $\theta_1 \in (0, 1)$,
- (iii)' $\|g'(x^k)^{-1} - J_k\| \leq \min\{M_2 \|x^{k+1} - x^k\|, \|g'(x^0)^{-1} - J_0\|\}$ for some $M_2 > 0$, or
- (iv)' $g'(x^k)^{-1} = J_k$,

which imply the conditions (i) through (iv) of Theorem 2.1 respectively, all imply the convergence of the inverse Jacobian approximates to $g'(\bar{x})^{-1}$. The conditions (i)'–(iv)' are less desirable since they require greater expense and care in the construction of the inverse Jacobian approximates.

3. Newton's Method for Minimization

We now translate the results of previous section to the optimization setting. The underlying problem is

$$\mathcal{P} \quad \min_{x \in \mathbb{R}^n} f(x) .$$

The Newton-like iterations take the form

$$x^{k+1} = x^k - H_k \nabla f(x^k),$$

where H_k is an approximation to the inverse of the Hessian matrix $\nabla^2 f(x^k)$.

THEOREM 3.1. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be twice continuously differentiable, $x^0 \in \mathbb{R}^n$, and $H_0 \in \mathbb{R}^{n \times n}$. Suppose that*

- (1) *there exists $\bar{x} \in \mathbb{R}^n$ and $\epsilon > 0$ such that $f(\bar{x}) \leq f(x)$ whenever $\|x - \bar{x}\| \leq \epsilon$,*
- (2) *there is a $\delta > 0$ such that $\delta \|z\|_2^2 \leq z^T \nabla^2 f(x) z$ for all $x \in B(\bar{x}, \epsilon)$,*
- (3) *$\nabla^2 f$ is Lipschitz continuous on $\text{cl}(B(\bar{x}, \epsilon))$ with Lipschitz constant L , and*
- (4) *$\theta_0 := \frac{L}{2\delta} \|x^0 - \bar{x}\| + M_0 K < 1$ where $M_0 > 0$ satisfies $z^T \nabla^2 f(x) z \leq M_0 \|z\|_2^2$ for all $x \in B(\bar{x}, \epsilon)$ and $K \geq \|(\nabla^2 f(x^0))^{-1} - H_0\|$ with $y^0 = \nabla f(x^0) / \|\nabla f(x^0)\|$.*

Further, suppose that the iteration

$$(107) \quad x^{k+1} := x^k - H_k \nabla f(x^k)$$

is initiated at x^0 where the H_k 's are chosen to satisfy one of the following conditions:

- (i) $\|(\nabla^2 f(x^k))^{-1} - H_k\| \leq K$,
- (ii) $\|(\nabla^2 f(x^k))^{-1} - H_k\| \leq \theta_1^k K$ for some $\theta_1 \in (0, 1)$,
- (iii) $\|(\nabla^2 f(x^k))^{-1} - H_k\| \leq \min\{M_2 \|x^k - x^{k-1}\|, K\}$, for some $M_2 > 0$, or
- (iv) $\|(\nabla^2 f(x^k))^{-1} - H_k\| \leq \min\{M_3 \|\nabla f(x^k)\|, K\}$, for some $M_3 > 0$,

where for each $k = 1, 2, \dots$ $y^k := \nabla f(x^k) / \|\nabla f(x^k)\|$.

These hypotheses on the accuracy of the approximations H_k yield the following conclusions about the rate of convergence of the iterates x^k .

- (a) If (i) holds, then $x^k \rightarrow \bar{x}$ linearly.
- (b) If (ii) holds, then $x^k \rightarrow \bar{x}$ superlinearly.
- (c) If (iii) holds, then $x^k \rightarrow \bar{x}$ two step quadratically.
- (d) If (iv) holds, then $x^k \rightarrow \bar{x}$ quadratically.

To more fully understand the convergence behavior described in this theorem, let us examine the nature of the controlling parameters L , M_0 , and M_1 . Since L is a Lipschitz constant for $\nabla^2 f$ it loosely corresponds to a bound on the third-order behavior of f . Thus the assumptions for convergence make implicit demands on the third derivative. The constant δ is a local lower bound on the eigenvalues of $\nabla^2 f$ near \bar{x} . That is, f behaves locally as if it were a *strongly convex function* (see exercises) with modulus δ . Finally, M_0 can be interpreted as a local Lipschitz constant for ∇f and only plays a role when $\nabla^2 f$ is approximated inexactly by H_k 's.

We now illustrate the performance differences between the method of steepest descent and Newton's method on a simple one dimensional problem. Let $f(x) = x^2 + e^x$. Clearly, f is a strongly convex function with

$$\begin{aligned} f(x) &= x^2 + e^x \\ f'(x) &= 2x + e^x \\ f''(x) &= 2 + e^x > 2 \\ f'''(x) &= e^x. \end{aligned}$$

If we apply the steepest descent algorithm with backtracking ($\gamma = 1/2$, $c = 0.01$) initiated at $x^0 = 1$, we get the following table

k	x^k	$f(x^k)$	$f'(x^k)$	s
0	1	.37182818	4.7182818	0
1	0	1	1	0
2	-.5	.8565307	-0.3934693	1
3	-.25	.8413008	0.2788008	2
4	-.375	.8279143	-.0627107	3
5	-.34075	.8273473	.0297367	5
6	-.356375	.8272131	-.01254	6
7	-.3485625	.8271976	.0085768	7
8	-.3524688	.8271848	-.001987	8
9	-.3514922	.8271841	.0006528	10
10	-.3517364	.827184	-.0000072	12

If we apply Newton's method from the same starting point and take a unit step at each iteration, we obtain a dramatically different table.

x	$f'(x)$
1	4.7182818
0	1
-1/3	.0498646
-.3516893	.00012
-.3517337	.0000000064

In addition, one more iteration gives $|f'(x^5)| \leq 10^{-20}$. This is a stunning improvement in performance and shows why one always uses Newton's method (or an approximation to it) whenever possible.

Our next objective is to develop numerically viable methods for approximating Jacobians and Hessians in Newton-like methods.

4. Matrix Secant Methods

Let us return to the problem of finding $\bar{x} \in \mathbb{R}^n$ such that $g(\bar{x}) = 0$ where $g : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is continuously differentiable. In this section we consider Newton-Like methods of a special type. Recall that in a Newton-Like method the iteration scheme takes the form

$$(108) \quad x^{k+1} := x^k - J_k g(x^k),$$

where J_k is meant to approximate the inverse of $g'(x^k)$. In the one dimensional case, a method proposed by the Babylonians 3700 years ago is of particular significance. Today we call it the *secant method*:

$$(109) \quad J_k = \frac{x^k - x^{k-1}}{g(x^k) - g(x^{k-1})}.$$

With this approximation one has

$$g'(x^k)^{-1} - J_k = \frac{g(x^{k-1}) - [g(x^k) + g'(x^k)(x^{k-1} - x^k)]}{g'(x^k)[g(x^{k-1}) - g(x^k)]}.$$

Near a point x^* at which $g'(x^*) \neq 0$ one can use the MVT to show there exists an $\alpha > 0$ such that

$$\alpha \|x - y\| \leq \|g(x) - g(y)\|.$$

Consequently, by the Quadratic Bound Lemma,

$$\|g'(x^k)^{-1} - J_k\| \leq \frac{\frac{L}{2} \|x^{k-1} - x^k\|^2}{\alpha \|g'(x^k)\| \|x^{k-1} - x^k\|} \leq K \|x^{k-1} - x^k\|$$

for some constant $K > 0$ whenever x^k and x^{k-1} are sufficiently close to x^* . Therefore, by our convergence Theorem for Newton Like methods, the secant method is locally two step quadratically convergent to a non-singular solution of the equation $g(x) = 0$. An additional advantage of this approach is that no extra function evaluations are required to obtain the approximation J_k .

4.0.1. *Matrix Secant Methods for Equations.* Unfortunately, the secant approximation (109) is meaningless if the dimension n is greater than 1 since division by vectors is undefined. But this can be rectified by multiplying (109) on the right by $(g(x^{k-1}) - g(x^k))$ and writing

$$(110) \quad J_k(g(x^k) - g(x^{k-1})) = x^k - x^{k-1}.$$

Equation (110) is called the *Quasi-Newton equation* (QNE), or *matrix secant equation* (MSE), at x^k . Here the matrix J_k is unknown, but is required to satisfy the n linear equations of the MSE. These equations determine an n dimensional affine manifold in $\mathbb{R}^{n \times n}$. Since J_k contains n^2 unknowns, the n linear equations in (110) are not sufficient to uniquely determine J_k . To nail down a specific J_k further conditions on the update J_k must be given. What conditions should these be?

To develop sensible conditions on J_k , let us consider an overall iteration scheme based on (108). For convenience, let us denote J_k^{-1} by B_k (i.e. $B_k = J_k^{-1}$). Using the B_k 's, the MSE (110) becomes

$$(111) \quad B_k(x^k - x^{k-1}) = g(x^k) - g(x^{k-1}).$$

At every iteration we have (x^k, B_k) and compute x^{k+1} by (108). Then B_{k+1} is constructed to satisfy (111). If B_k is close to $g'(x^k)$ and x^{k+1} is close to x^k , then B_{k+1} should be chosen not only to satisfy (111) but also to be as “close” to B_k as possible. With this in mind, we must now decide what we mean by “close”. From a computational perspective, we prefer “close” to mean *easy to compute*. That is, B_{k+1} should be *algebraically close* to B_k in the sense that B_{k+1} is only a rank 1 modification of B_k . Since we are assuming that B_{k+1} is a rank 1 modification to B_k , there are vectors $u, v \in \mathbb{R}^n$ such that

$$(112) \quad B_{k+1} = B_k + uv^T.$$

We now use the matrix secant equation (111) to derive conditions on the choice of u and v . In this setting, the MSE becomes

$$B_{k+1}s^k = y^k,$$

where

$$s^k := x^{k+1} - x^k \quad \text{and} \quad y^k := g(x^{k+1}) - g(x^k).$$

Multiplying (??) by s^k gives

$$y^k = B_{k+1}s^k = B_k s^k + uv^T s^k.$$

Hence, if $v^T s^k \neq 0$, we obtain

$$u = \frac{y^k - B_k s^k}{v^T s^k}$$

and

$$(113) \quad B_{k+1} = B_k + \frac{(y^k - B_k s^k) v^T}{v^T s^k}.$$

Equation (113) determines a whole class of rank one updates that satisfy the MSE where one is allowed to choose $v \in \mathbb{R}^n$ as long as $v^T s^k \neq 0$. If $s^k \neq 0$, then an obvious choice for v is s^k yielding the update

$$(114) \quad B_{k+1} = B_k = \frac{(y^k - B_k s^k) s^{kT}}{s^{kT} s^k}.$$

This is known as Broyden's update. It turns out that the Broyden update is also analytically close.

THEOREM 4.1. *Let $A \in \mathbb{R}^{n \times n}$, $s, y \in \mathbb{R}^n$, $s \neq 0$. Then for any matrix norms $\|\cdot\|$ and $\|\cdot\|_2$ such that*

$$\|AB\| \leq \|A\| \|B\|_2$$

and

$$\left\| \frac{vv^T}{v^T v} \right\|_2 \leq 1,$$

the solution to

$$(115) \quad \min\{\|B - A\| : Bs = y\}$$

is

$$(116) \quad A_+ = A + \frac{(y - As)s^T}{s^T s}.$$

In particular, (116) solves (115) when $\|\cdot\|$ is the ℓ_2 matrix norm, and (116) solves (115) uniquely when $\|\cdot\|$ is the Frobenius norm.

PROOF. Let $B \in \{B \in \mathbb{R}^{n \times n} : Bs = y\}$, then

$$\begin{aligned} \|A_+ - A\| &= \left\| \frac{(y - As)s^T}{s^T s} \right\| = \left\| (B - A) \frac{ss^T}{s^T s} \right\| \\ &\leq \|B - A\| \left\| \frac{ss^T}{s^T s} \right\|_2 \leq \|B - A\|. \end{aligned}$$

Note that if $\|\cdot\|_2 = \|\cdot\|_2$, then

$$\begin{aligned} \left\| \frac{vv^T}{v^T v} \right\|_2 &= \sup \left\{ \left\| \frac{vv^T}{v^T v} x \right\|_2 \mid \|x\|_2 = 1 \right\} \\ &= \sup \left\{ \sqrt{\frac{(v^T x)^2}{\|v\|^2}} \mid \|x\|_2 = 1 \right\} \\ &= 1, \end{aligned}$$

so that the conclusion of the result is not vacuous. For uniqueness observe that the Frobenius norm is strictly convex and $\|A \cdot B\|_F \leq \|A\|_F \|B\|_2$. \square

Therefore, the Broyden update (114) is both algebraically and analytically close to B_k . These properties indicate that it should perform well in practice and indeed it does.

Algorithm: Broyden's Method

Initialization: $x^0 \in \mathbb{R}^n$, $B_0 \in \mathbb{R}^{n \times n}$

Having (x^k, B_k) compute (x^{k+1}, B_{k+1}) as follows:

Solve $B_k s^k = -g(x^k)$ for s^k and set

$$\begin{aligned} x^{k+1} &: = x^k + s^k \\ y^k &: = g(x^{k+1}) - g(x^k) \\ B_{k+1} &: = B_k + \frac{(y^k - B_k s^k)s^{kT}}{s^{kT} s^k}. \end{aligned}$$

We would prefer to write the Broyden update in terms of the matrices $J_k = B_k^{-1}$ so that we can write the step computation as $s^k = -J_k g(x^k)$ avoiding the need to solve an equation. To obtain the formula for J_k we use the following important lemma for matrix inversion.

LEMMA 4.1. (Sherman-Morrison-Woodbury) *Suppose $A \in \mathbb{R}^{n \times n}$, $U \in \mathbb{R}^{n \times k}$, $V \in \mathbb{R}^{n \times k}$ are such that both A^{-1} and $(I + V^T A^{-1} U)^{-1}$ exist, then*

$$(A + UV^T)^{-1} = A^{-1} - A^{-1}U(I + V^T A^{-1}U)^{-1}V^T A^{-1}$$

The above lemma verifies that if $B_k^{-1} = J_k$ exists and $s^{kT} J_k y^k = s^{kT} B_k^{-1} y^k \neq 0$, then

$$(117) \quad J_{k+1} = \left[B_k + \frac{(y^k - B_k s^k) s^{kT}}{s^{kT} s^k} \right]^{-1} = B_k^{-1} + \frac{(s^k - B_k^{-1} y^k) s^{kT} B_k^{-1}}{s^{kT} B_k^{-1} y^k} = J_k + \frac{(s^k - J_k y^k) s^{kT} J_k}{s^{kT} J_k y^k}.$$

In this case, it is possible to directly update the inverses J_k . It should be cautioned though that this process can become numerically unstable if $|s^{kT} J_k y^k|$ is small. Therefore, in practise, the value $|s^{kT} J_k y^k|$ must be monitored to avoid numerical instability.

Although we do not pause to establish the convergence rates here, we do give the following result due to Dennis and Moré (1974).

THEOREM 4.2. *Let $g : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be continuously differentiable in an open convex set $D \subset \mathbb{R}^n$. Assume that there exists $x^* \in \mathbb{R}^n$ and $r, \beta > 0$ such that $x^* + r\mathbb{B} \subset D$, $g(x^*) = 0$, $g'(x^*)^{-1}$ exists with $\|g'(x^*)^{-1}\| \leq \beta$, and g' is Lipschitz continuous on $x^* + r\mathbb{B}$ with Lipschitz constant $\gamma > 0$. Then there exist positive constants ϵ and δ such that if $\|x^0 - x^*\|_2 \leq \epsilon$ and $\|B_0 - g'(x^0)\| \leq \delta$, then the sequence $\{x^k\}$ generated by the iteration*

$$\begin{cases} x^{k+1} & := x^k + s^k \text{ where } s^k \text{ solves } 0 = g(x^k) + B_k s \\ B_{k+1} & := B_k + \frac{(y^k - B_k s^k) s_k^T}{s_k^T s^k} \text{ where } y^k = g(x^{k+1}) - g(x^k) \end{cases}$$

is well-defined with $x^k \rightarrow x^*$ superlinearly.

4.0.2. Matrix Secant Methods for Minimization. We now extend these matrix secant ideas to optimization, specifically minimization. The underlying problem we consider is

$$\mathcal{P} : \underset{x \in \mathbb{R}^n}{\text{minimize}} f(x),$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is assumed to be twice continuously differentiable. In this setting, we wish to solve the equation $\nabla f(x) = 0$ and the MSE (111) becomes

$$(118) \quad H_{k+1} y^k = s^k,$$

where $s^k := x^{k+1} - x^k$ and

$$y^k := \nabla f(x^{k+1}) - \nabla f(x^k).$$

Here the matrix H_k is intended to be an approximation to the inverse of the hessian matrix $\nabla^2 f(x^k)$. Writing $M_k = H_k^{-1}$, a straightforward application of Broyden's method gives the update

$$M_{k+1} = M_k + \frac{(y^k - M_k s^k) s^{kT}}{s^{kT} s^k}.$$

However, this is unsatisfactory for two reasons:

- (1) Since M_k approximates $\nabla^2 f(x^k)$ it must be symmetric.
- (2) Since we are minimizing, then M_k must be positive definite to insure that $s^k = -M_k^{-1} \nabla f(x^k)$ is a direction of descent for f at x^k .

To address problem 1 above, one could return to equation (113) and find an update that preserves symmetry. Such an update is uniquely obtained by setting

$$v = (y^k - M_k s^k).$$

This is called the symmetric rank 1 update or SR1. Although this update can on occasion exhibit problems with numerical stability, it has recently received a great deal of renewed interest. The stability problems occur whenever

$$(119) \quad v^T s^k = (y^k - M_k s^k)^T s^k$$

has small magnitude. The inverse SR1 update is given by

$$H_{k+1} = H_k + \frac{(s^k - H_k y^k)(s^k - H_k y^k)^T}{(s^k - H_k y^k)^T y^k}$$

which exists whenever $(s^k - H_k y^k)^T y^k \neq 0$.

We now approach the question of how to update M_k in a way that addresses both the issue of symmetry and positive definiteness while still using the Broyden updating ideas. Given a symmetric positive definite matrix M and two vectors s and y , our goal is to find a symmetric positive definite matrix \bar{M} such that $\bar{M}s = y$. Since M is symmetric and positive definite, there is a non-singular $n \times n$ matrix L such that $M = LL^T$. Indeed, L can be

chosen to be the lower triangular Cholesky factor of M . If \overline{M} is also symmetric and positive definite then there is a matrix $J \in \mathbb{R}^{n \times n}$ such that $\overline{M} = JJ^T$. The MSE (??) implies that if

$$(120) \quad J^T s = v$$

then

$$(121) \quad Jv = y.$$

Let us apply the Broyden update technique to (121), J , and L . That is, suppose that

$$(122) \quad J = L + \frac{(y - Lv)v^T}{v^T v}.$$

Then by (120)

$$(123) \quad v = J^T s = L^T s + \frac{v(y - Lv)^T s}{v^T v}.$$

This expression implies that v must have the form

$$v = \alpha L^T s$$

for some $\alpha \in \mathbb{R}$. Substituting this back into (123) we get

$$\alpha L^T s = L^T s + \frac{\alpha L^T s (y - \alpha L L^T s)^T s}{\alpha^2 s^T L L^T s}.$$

Hence

$$(124) \quad \alpha^2 = \left[\frac{s^T y}{s^T M s} \right].$$

Consequently, such a matrix J satisfying (123) exists only if $s^T y > 0$ in which case

$$J = L + \frac{(y - \alpha M s) s^T L}{\alpha s^T M s},$$

with

$$\alpha = \left[\frac{s^T y}{s^T M s} \right]^{1/2},$$

yielding

$$(125) \quad \overline{M} = M + \frac{yy^T}{y^T s} - \frac{M s s^T M}{s^T M s}.$$

Moreover, the Cholesky factorization for \overline{M} can be obtained directly from the matrices J . Specifically, if the QR factorization of J^T is $J^T = QR$, we can set $\overline{L} = R$ yielding

$$\overline{M} = JJ^T = R^T Q^T QR = \overline{L}\overline{L}^T.$$

The formula for updating the inverses is again given by applying the Sherman-Morrison-Woodbury formula to obtain

$$(126) \quad \overline{H} = H + \frac{(s + Hy)^T y s s^T}{(s^T y)^2} - \frac{H y s^T + s y^T H}{s^T y},$$

where $H = M^{-1}$. The update (125) is called the BFGS update and (126) the inverse BFGS update. The letter BFGS stand for Broyden, Fletcher, Goldfarb, and Shanno.

We have shown that beginning with a symmetric positive definite matrix M_k we can obtain a symmetric and positive definite update M_{k+1} that satisfies the MSE $M_{k+1} s_k = y_k$ by applying the formula (125) whenever $s_k^T y_k > 0$. We must now address the question of how to choose x^{k+1} so that $s_k^T y_k > 0$. Recall that

$$y = y^k = \nabla f(x^{k+1}) - \nabla f(x^k)$$

and

$$s^k = x^{k+1} - x^k = t_k d^k,$$

where

$$d^k = -t_k H_k \nabla f(x^k)$$

is the matrix secant search direction and t_k is the stepsize. Hence

$$\begin{aligned} y^{kT} s^k &= \nabla f(x^{k+1})^T s^k - \nabla f(x^k)^T s^k \\ &= t_k (\nabla f(x^k + t_k d_k)^T d^k - \nabla f(x^k)^T d^k), \end{aligned}$$

where $d^k := -H_k \nabla f(x^k)$. Since H_k is positive definite the direction d^k is a descent direction for f at x^k and so $t_k > 0$. Therefore, to insure that $s^{kT} y^k > 0$ we need only show that $t_k > 0$ can be chosen so that

$$(127) \quad \nabla f(x^k + t_k d^k)^T d^k \geq \beta \nabla f(x^k)^T d^k$$

for some $\beta \in (0, 1)$ since in this case

$$\nabla f(x^k + t_k d_k)^T d^k - \nabla f(x^k)^T d^k \geq (\beta - 1) \nabla f(x^k)^T d^k > 0.$$

But this precisely the second condition in the weak Wolfe conditions with $\beta = c_2$. Hence a successful BFGS update can always be obtained. The BFGS update and is currently considered the best matrix secant update for minimization.

BFGS Updating

$$\begin{aligned} \sigma &:= \sqrt{s^{kT} y^k} \\ \hat{s}^k &:= s^k / \sigma \\ \hat{y}^k &:= y^k / \sigma \\ H_{k+1} &:= H_k + (\hat{s}^k - H_k \hat{y}^k)(\hat{s}^k)^T + \hat{s}^k (\hat{s}^k - H_k \hat{y}^k)^T - (\hat{s}^k - H_k \hat{y}^k)^T \hat{y}^k \hat{s}^k (\hat{s}^k)^T \end{aligned}$$

Index

- \mathcal{S}^n , 39
- \mathcal{S}_{++}^n , 39
- \mathcal{S}_{+}^n , 39
- adjoint, 38
- Affine Sets, 42
- basis, 19
- Block Matrix Multiplication, 11
- block structure, 7
- boundary point, 54
- bounded set, 54
- chain rule, 59
- characteristic polynomial, 37
- Cholesky diagonalization procedure, 45
- Cholesky Factorization, 46
- closed set, 54
- closed unit ball, 54
- cluster point, 54
- compact set, 54
- complementary subspaces, 28
- condensed QR Factorization, 32
- Conjugacy, 49
- Conjugate Gradient Algorithm, 48, 50
- constraint region, 5
- Continuous Functions, 53
- delta method, 58
- Differentiable Functions, 56
- dimension, 19
- directional derivative, 55
- Eigenvalue Decomposition for Symmetric Matrices, 38
- Elementary Matrices, 10
- Elementary permutation matrices, 10
- elementary row operations, 10
- elementary unit coordinate matrices, 10
- finite impulse response, 24
- first-order approximation, 59
- first-order Taylor expansion, 59
- Four Fundamental Subspaces, 19
- Frobenius inner product, 55
- Frobenius norm, 55
- full QR factorization, 32
- Fundamental Theorem of the Alternative, 20
- Gauss-Jordan Elimination Matrices, 13
- Gaussian-Jordan Elimination Matrix, 14
- generalized Cholesky factorization, 46
- Gram-Schmidt orthogonalization, 31
- Hessian, 57
- Householder reflection matrices, 34
- Identity Matrix, 15
- indefinite, 39
- Inverse Matrices, 16
- Jacobian, 56
- Kalman Smoothing, 24
- Lagrange multiplier, 43
- linear combination, 8
- linear least squares, 6
- linear programming, 5
- Linear Regression, 23
- linear span, 8
- linearly independent, 19
- Lipschitz Continuity, 61
- LU Factorization, 17
- matrix free, 48
- matrix multiplication, 9
- Matrix vector Multiplication, 8
- maximum likelihood estimation, 23
- Mean Value Theorem, 59
- Modified Gram-Schmidt, 31
- modulus of Lipschitz continuity, 61
- negative definite, 39
- nonlinear programming, 5
- normal equations, 26
- objective function, 5
- open set, 54
- orthogonal projection, 28
- Orthogonal Projections, 29
- partial derivative, 56
- Permutation Matrices, 16
- permutation matrix, 16
- Polynomial Fitting, 21
- positive definite, 39
- positive semi-definite, 39
- positively homogeneous, 56
- Principal Minor, 45
- projection, 28
- QR factorization, 31
- quadratic functions, 6
- Rayleigh-Ritz Theorem, 39
- Reduced Echelon Form, 13

second-order Taylor approximation, 61
self-adjoint, 38
Skew Symmetric Matrices, 16
Square Matrices, 15
Square Roots of Positive Semi-Definite Matrices, 40
Subspace Projection Theorem, 29
Symmetric Matrices, 16
symmetric part, 6, 37
System Identification in Signal Processing, 24

Taylor's Theorem with remainder, 60
Toeplitz matrices, 24

unconstrained optimization problem, 5
uniformly continuous, 53
unit vectors, 54
Unitary Matrices, 16
unitary matrix, 38

Vandermonde matrix, 22
Vector Norm, 53

Weierstrass Compactness Theorem, 54
Weierstrass Extreme Value Theorem, 54