

ELEC 539A Lecture Notes

Selected Topics in Signal Processing: Machine Learning for Signal Processing

Wu-Sheng Lu

Department of Electrical and Computer Engineering

University of Victoria

January 2015

© Wu-Sheng Lu, 2015
University of Victoria

All rights reserved. This set of notes may not be reproduced in whole or in part by photocopy or other means, without the permission of the author.

COURSE OUTLINE

- **Instructor:**

Dr. W.-S. Lu
Phone: 8692
E-mail: wslu@ece.uvic.ca
URL: www.ece.uvic.ca/~wslu

- **Office Hours:**

Days: Wednesdays
Time: 14:40 – 16:40
Location: EOW 427

- **Lectures:**

Section: A01/CRN 23843
Days: Tuesdays, Wednesdays, and Fridays
Time: 13:30 – 14:20
Location: ELL 161

- **Text:**

Lecture Notes

- **Assessment:**

Assignments:	30%
Project:	30%
Final:	40%

- **Due Dates for Assignments and Project Report:**

Each assignment due date will be given in class as well as posted in the class web site. To receive course credit, you are required to submit the assignment no later than 4:00 pm on the due day. This course does not use drop-box for assignments and project reports. There are three ways to submit your assignment work:

- submit it to the instructor in class.
 - submit by placing it in the instructor's mailbox in EOW 448.
 - submit by sliding it underneath the door of EOW 427.
- Project report is due by 4:00 pm on the same day as the final exam.

The course requires CVX – a software package for convex optimization.

- Download current version of CVX from <http://cvxr.com/cvx/download/> and a manual from <http://cvxr.com/cvx/doc/>
- You are strongly encouraged to go through the manual while trying the software on your own as early as possible.
- The course covers basic elements of supervised and unsupervised learning methods with applications in several signal processing problems. Review of key concepts of probability theory and numerical optimization techniques will be provided. A precise list of course contents will be made available at a later point of time.

References

- [1] V. N. Vapnik, *Statistical Learning Theory*, Wiley, 1998.
- [2] V. N. Vapnik, *The Nature of Statistical Learning Theory*, 2nd ed., Springer, 2000.
- [3] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*, 2nd ed., Springer, 2009.
- [4] H. Stark and J. W. Woods, *Probability and Random Processes with Applications to Signal Processing*, 3rd ed., Prentice Hall, 2002.
- [5] A. Antoniou and W.-S. Lu, *Practical Optimization—Algorithms and Engineering Applications*, Springer 2007.
- [6] S. Boyd and L. Vandenberghe, *Convex Optimization*, Cambridge University Press, 2004.
- [7] C. M. Bishop, *Pattern Recognition and Machine Learning*, Springer, 2006.
- [8] Y. S. Abu-Mostafa, M. Magdon-Ismail, and H.-T. Lin, *Learning from Data – A Short Course*, AMLbook.com, 2012.
- [9] F. Rosenblatt, “The perception: A probabilistic model for information storage and organization in the brain,” *Psychological Review*, vol. 65, no. 6, pp. 386-408, 1958.
- [10] G.-X. Yuan, C.-H. Ho, and C.-J. Lin, “Recent advances of large-scale linear classification,” *Proc. IEEE*, vol. 100, no. 9, pp. 2584-2603, Sept. 2012.
- [11] T. Poggio and S. Smale, “The mathematics of learning: Dealing with data,” *Notices of the AMS*, vol. 50, no. 5, pp. 537-544, May 2003.
- [12] G. H. Golub and C. F. Van Loan, *Matrix Computations*, 3rd edition, The Johns Hopkins University Press, 1989.

Chapter 1 Preliminaries

1.1 Probability Theory

We begin by presenting a brief review of the **probability theory**. The very reason of doing so is that contemporary machine learning methodologies are developed in a probabilistic framework [1] - [3]. In other words, most part of machine learning makes sense only when you look at it from a **statistical perspective**.

1.1.1 Definition of Probability

A. Classical Probability

There are different kinds of probability: probability as intuition, probability as the ratio of favorable to total outcomes (also known as classical theory of probability), probability as a measure of frequency of occurrence, and probability based on axiomatic theory [4].

Definition 1.1 **Random experiment and basic event**

An experiment E is said to be a *random experiment* if it can be performed repeatedly and its outcomes are *not* deterministic but probabilistic. An outcome of a random experiment, denoted by ω , is called a *basic event*. ■

Definition 1.2 **Sample space**

The set of all basic events, denoted by Ω , is called *sample space*. ■

Example 1.1

Let E be the experiment of picking a ball at random from a box containing N identical balls, where the balls have been marked as #1, #2, ..., to # N . Obviously E is random experience because its outcome has N possibilities. If the number on the ball that has been picked is i , then the outcome is a basic event and is denoted by ω_i . Obviously, the **sample space is**

$$\Omega = \{\omega_1, \dots, \omega_N\}. \quad \blacksquare$$

Example 1.2

Let E be that of observing the average daytime temperature in a city in October. Obviously E is a random experiment. If we denote the outcome of experiment E by ω_a , then the sample space in

this case may be described as $\Omega = \{\omega_a, -\infty < \omega_a < \infty\}$. Note that Ω contains infinite number of basic events. ■

Definition 1.3 **General event**

A general event (or simply called it event), say A , consists of several basic events, hence a general event is a subset of sample space Ω . ■

Example 1.3

Continue from Example 1.1. Suppose one randomly picks a ball from the box, one can state that

- Event A of getting a ball whose number is no greater than 3 is a general event because it consists of three basic events: $A = \{\omega_1, \omega_2, \omega_3\}$ which is a proper subset of $\Omega = \{\omega_1, \dots, \omega_N\}$.
- Event B of getting a ball with number even numbers is a general event because $B =$

$$\{\omega_2, \omega_4, \dots, \omega_{100}\} \subset \Omega. \quad \blacksquare$$

The classical probability theory applies to the cases where the random experiment E obeys two basic assumptions: (i) the total number of outcomes, N , is finite; and (ii) all outcomes (i.e. the results of basic events) are equally likely. Consequently, the probability of a (general) *event* A , denoted by $P(A)$, is obtained *a priori* (i.e., relating to reasoning that proceeds from theoretical deduction rather than from observation or experience) by counting the number of ways N_e that event A can occur, then computing the ratio N_e/N as the probability. This is,

$$P(A) = \frac{N_e}{N} \quad (1.1)$$

Example 1.4

Continue from Example 1.3.

Let A be the event of obtaining a ball whose number is no greater than 3. The probability $P(A) = 3/100 = 0.03$.

Let B be the event of obtaining a ball that is even-numbered. The probability $P(B) = 50/100 = 0.5$.

Let C be the event of obtaining a ball whose number is a multiple of 3. The probability $P(C) = 33/100 = 0.33$.

■

The major problems with the classical theory of probability are that it cannot handle outcomes that are not equally likely; and it cannot deal with infinite number of outcomes.

B. Probability Based On Axiomatic Theory

Developed by A. N. Kolmogorov in 1930's, the probability based on axiomatic theory is the one that is followed by most modern texts on the subject. Kolmogorov's probability theory is set-theoretic in that a random event is considered as a subset of a sample space in an abstract setting where the two basic assumptions made in the classical probability theory are no longer necessary.

In order to present Kolmogorov's axiomatic definition of probability, basic elements of *set algebra* and the sigma algebra (σ – algebra) are sketched first.

Set Algebra

- A *set* is a collection of objects (or elements).
- A set B is said to be a *subset* of set A , and this relation is denoted by $B \subseteq A$, if B is contained within set A . Let for example set A denote all Victoria residents and set B be all Victoria residents whose height is between 5.5 and 6.5 feet. Obviously set B is a (proper) subset of set A .
- Let A and B be two sets in space Ω . The *union* (sum) of A and B , denoted by $A \cup B$ or $A + B$, is the set of elements that are in at least one of the sets A or B .
- The *intersection* (or set product) of A and B , denoted by $A \cap B$ or AB , is the set of elements that are in both A and B .
- The *empty set*, denoted by ϕ , is a set that contains no objects.
- Let A be a set in space Ω , the complement of A , denoted by A^c , is a set of all elements that are not in A . Obviously, $A \cup A^c = \Omega$ and $A \cap A^c = \phi$.

- Two sets A and B are said to be equal if both $B \subseteq A$ and $A \subseteq B$ hold.
- The difference of A and B , denoted by $A - B$, is the set of elements that are in A but not in B . It follows that

$$A - B = AB^c \quad \text{and} \quad B - A = BA^c$$

- The *exclusive-or* of sets A and B , denoted by $A \oplus B$, is the set of elements that are in A or B , but not in both. Obviously, we can write

$$A \oplus B = (A - B) \cup (B - A)$$

- Sets A and B are said to be disjoint if they have no elements in common, thus $AB = \emptyset$.
- It can readily be verified that

$$(A \cup B)^c = A^c \cap B^c \quad \text{and} \quad (A \cap B)^c = A^c \cup B^c$$

which can be extended by induction to the general case of n sets as

$$\left(\bigcup_{i=1}^n A_i \right)^c = \bigcap_{i=1}^n A_i^c \quad \text{and} \quad \left(\bigcap_{i=1}^n A_i \right)^c = \bigcup_{i=1}^n A_i^c$$

Definition 1.4 σ -Algebra

Let Ω be a set and we consider a collection of subsets of Ω , denoted by \mathcal{F} . \mathcal{F} is said to be a σ -Algebra if

1. $\emptyset \in \mathcal{F}$ and $\Omega \in \mathcal{F}$,
2. If $A \in \mathcal{F}$ and $B \in \mathcal{F}$, then $A \cup B \in \mathcal{F}$ and $A \cap B \in \mathcal{F}$.
3. If $A \in \mathcal{F}$ then $A^c \in \mathcal{F}$.
4. \mathcal{F} is closed under countable set of unions, intersections, and combinations. Hence if A_1, \dots, A_n, \dots belong to \mathcal{F} , then

$$\bigcup_{i=1}^{\infty} A_i \in \mathcal{F} \quad \text{and} \quad \bigcap_{i=1}^{\infty} A_i \in \mathcal{F} \quad \blacksquare$$

Why is σ -Algebra relevant?

Recall in the classical theory of probability the number of basic events must be finite and all outcomes are equally likely, i.e. the probability of every basic event is the same. The axiomatic probability theory built on σ -algebra allows us to remove these two fundamental limitations. Consider a random experiment and the associated sample space Ω . In the axiomatic probability theory we only consider the collection \mathcal{F} of those subsets of Ω that form a σ -algebra. These subsets are called events. It turns out that \mathcal{F} usually includes all subsets of engineering and science interest.

Definition 1.5 *Axiomatic probability*

Given a sample space Ω and a σ -algebra \mathcal{F} of Ω , probability $P[\cdot]$ is a set function that

assigns to every event $A \in \mathcal{F}$ that obeys the three axioms:

- (1) $P[A] \geq 0$.
- (2) $P[\Omega] = 1$.
- (3) $P[A \cup B] = P[A] + P[B]$ if $A \cap B = \emptyset$.

From these axioms, one can establish the following important properties:

- (4) $P[\emptyset] = 0$.
- (5) $P[A \cap B^c] = P[A] - P[A \cap B]$.
- (6) $P[A] = 1 - P[A^c]$.
- (7) $P[A \cup B] = P[A] + P[B] - P[A \cap B]$.
- (8) $P\left[\bigcup_{i=1}^n A_i\right] = \sum_{i=1}^n P[A_i]$ if $A_i \cap A_j = \emptyset$ for all $i \neq j$. ■

Definition 1.6 *Probability space*

The triple of a sample space Ω , a collection of (general) events \mathcal{F} that form a σ -algebra, and a probability measure P , namely (Ω, \mathcal{F}, P) , is called a probability space. ■

1.1.2 Joint Probability, Conditional Probability, and Independence

A. Joint Probability

The *joint probability* of two events A and B is defined as the probability of “events A and B both occur”, namely $P[A \cap B]$ or $P[AB]$.

B. Conditional Probability

Often times people want to know the probability of “event A occurs” given that event B has occurred. This is called *conditional probability* and is denoted by $P[A|B]$. It is intuitively clear that because of the presence of the condition that “event B has occurred”, $P[A|B]$ is in general different from $P[A]$.

Example 1.5

Suppose there are four identical balls in a box, which are marked as #1, #2, #3, and #4. The random experiment is to pick up a ball from the box. Now consider two events: Event A is that of obtaining ball #4; and event B is that of obtaining an even-numbered ball. Obviously, we have $P[A] = 1/4$ and $P[A|B] = 1/2$. ■

A question that naturally arises is how to define and compute conditional probability in general circumstances? Let A and B are two events in a random experiment. Each outcome of the random experiment must fall into one of the four cases: (1) A occurs and B does not; (2) B occurs and A does not; (3) A and B both occur; and (4) A and B both do not occur.

Suppose one applies above analysis to Example 1.5, repeats the experiment n times, and denotes the times that each of the four cases occurs by n_1, n_2, n_3 , and n_4 , respectively. We can state that

- $n_1 + n_2 + n_3 + n_4 = n$
- the frequency of event $B = F_n(B) = \frac{n_2 + n_3}{n}$.
- the frequency of event $AB = F_n(AB) = \frac{n_3}{n}$.
- given that event B has occurred, the frequency of event $A = F_n(A|B) = \frac{n_3}{n_2 + n_3}$.

It follows that

$$F_n(A|B) = \frac{F_n(AB)}{F_n(B)} \text{ provided that } F_n(B) > 0 \quad (1.2)$$

Based on (1.2), conditional probability is defined as follows.

Definition 1.7 *Conditional probability*

Let (Ω, \mathcal{F}, P) be a probability space, $A \in \mathcal{F}, B \in \mathcal{F}$ with $P[B] > 0$. The conditional probability of event A given that event B has occurred is defined by

$$P[A|B] = \frac{P[AB]}{P[B]} \quad (1.3)$$

Similarly, the conditional probability of event B given that event A has occurred is defined by

$$P[B|A] = \frac{P[AB]}{P[A]} \quad (1.4)$$

provided that $P[A] > 0$. From (1.3) and (1.4) it follows that

$$P[AB] = P[A]P[B|A] = P[B]P[A|B] \quad \blacksquare \quad (1.5)$$

Example 1.6

Consider a binary communication system with a two-symbol alphabet, i.e., 0 and 1. Let X and Y be the transmitted and received symbols, respectively. Here the sample space is given by $\Omega = \{(X, Y) : X = 0 \text{ or } 1, Y = 0 \text{ or } 1\} = \{(0, 0), (0, 1), (1, 0), (1, 1)\}$. Suppose that the communication system is slightly corrupted by noise such that

$$\begin{aligned} P[Y = 1 | X = 1] &= 0.92, & P[Y = 0 | X = 1] &= 0.08 \\ P[Y = 0 | X = 0] &= 0.92, & P[Y = 1 | X = 0] &= 0.08 \end{aligned}$$

And by design, $P[X = 0] = P[X = 1] = 0.5$. Under these circumstances, we have

$$\begin{aligned}
P[X = 0, Y = 0] &= P[X = 0] \cdot P[Y = 0 | X = 0] = 0.5 \times 0.92 = 0.46 \\
P[X = 0, Y = 1] &= P[X = 0] \cdot P[Y = 1 | X = 0] = 0.5 \times 0.08 = 0.04 \\
P[X = 1, Y = 0] &= P[X = 1] \cdot P[Y = 0 | X = 1] = 0.5 \times 0.08 = 0.04 \\
P[X = 1, Y = 1] &= P[X = 1] \cdot P[Y = 1 | X = 1] = 0.5 \times 0.92 = 0.46
\end{aligned}$$

■

Properties of Conditional Probability

Property 1 (Multiplication formula)

Let A_1, A_2, \dots, A_n be n events with $n \geq 2$ and $P[A_1 A_2 \cdots A_{n-1}] > 0$, then

$$P[A_1 A_2 \cdots A_n] = P[A_1] \cdot P[A_2 | A_1] \cdot P[A_3 | A_1 A_2] \cdots P[A_n | A_1 A_2 \cdots A_{n-1}] \quad (1.6)$$

Proof:

Using (1.3), we see that the right-hand side of (1.6) is equal to

$$P[A_1] \cdot \frac{P[A_1 A_2]}{P[A_1]} \cdot \frac{P[A_1 A_2 A_3]}{P[A_1 A_2]} \cdots \frac{P[A_1 A_2 \cdots A_n]}{P[A_1 A_2 \cdots A_{n-1}]} = P[A_1 A_2 \cdots A_n] \quad \blacksquare$$

Property 2 (Unconditional probability)

Let A_1, A_2, \dots, A_n be mutually exclusive events such that $\bigcup_{i=1}^n A_i = \Omega$ with $P[A_i] > 0$ for all i .

Let B be any event defined over the probability space of A_i 's. Then

$$P[B] = \sum_{i=1}^n P[A_i] P[B | A_i] \quad (1.7)$$

Proof:

From

$$B = B \Omega = B \bigcup_{i=1}^n A_i = \bigcup_{i=1}^n B A_i$$

where $\{B A_i\}$ are mutually exclusive, it follows that

$$P[B] = P\left[\bigcup_{i=1}^n B A_i\right] = \sum_{i=1}^n P[B A_i] = \sum_{i=1}^n P[A_i] \cdot P[B | A_i] \quad \blacksquare$$

Property 3 (Bayes formula)

Let A_1, A_2, \dots, A_n be mutually exclusive events such that $\bigcup_{i=1}^n A_i = \Omega$ with $P[A_i] > 0$ for all i .

Let B be any event with $P[B] > 0$. Then

$$P[A_j | B] = \frac{P[B | A_j] \cdot P[A_j]}{\sum_{i=1}^n P[B | A_i] \cdot P[A_i]} \quad (1.8)$$

Proof:

By following the definition of conditional probability and the formula of unconditional probability in (1.7), we obtain

$$P[A_j | B] = \frac{P[BA_j]}{P[B]} = \frac{P[B | A_j] \cdot P[A_j]}{\sum_{i=1}^n P[B | A_i] \cdot P[A_i]} \quad \blacksquare$$

The Bayes formula finds many applications in science and engineering. The terms in (1.8) bear various names: $P[A_j|B]$ is often called *a posteriori* probability of A_j given B ; $P[B|A_j]$ is called the *a priori* probability of B given A_j ; and $P[A_j]$ is called the *causal* or *a priori* probability of A_j . Typically *a priori* probabilities are estimated from past measurements or presupposed by experience, while *a posteriori* probabilities are measured or computed from observations [4].

Example 1.7

Suppose there are three boxes that look identical, each contains certain number of red and blue balls that are identical except the color. The number of red and blue balls in box i are r_i and b_i , respectively, for $i = 1, 2, 3$. The experiment in question is to randomly pick a box, then randomly pick a ball. The outcome was a red ball. Given that it is a red ball, compute the probability of the ball belongs to box 1.

Solution

Let A_i be the event of the ball in question belongs to box i for $i = 1, 2, 3$. Obviously, $\{A_i \text{ for } i = 1, 2, 3\}$ are mutually exclusive, $\bigcup_{i=1}^3 A_i = \Omega$, and $P[A_1] = P[A_2] = P[A_3] = 1/3$. We also define event B as “it is a red ball”. With these definitions, the problem we need to address is to compute the conditional probability $P[A_1|B]$.

Clearly, formula (1.8) is applicable if we are able to compute conditional probabilities $P[B|A_i]$ for $i = 1, 2, 3$. These conditional probabilities are found to be

$$P[B | A_i] = \frac{r_i}{r_i + b_i} \quad \text{for } i = 1, 2, 3.$$

Hence (1.8) yields

$$P[A_1 | B] = \frac{\frac{1}{3} \frac{r_1}{r_1 + b_1}}{\frac{1}{3} \frac{r_1}{r_1 + b_1} + \frac{1}{3} \frac{r_2}{r_2 + b_2} + \frac{1}{3} \frac{r_3}{r_3 + b_3}} = \frac{1}{1 + \frac{r_2(r_1 + b_1)}{r_1(r_2 + b_2)} + \frac{r_3(r_1 + b_1)}{r_1(r_3 + b_3)}} \quad \blacksquare$$

C. Independence

Consider two events A and B . We have seen that in general probability $P[A]$ differs from the conditional probability $P[A|B]$ as long as the occurrence of event B has an impact on occurrence of event A . In other words, if $P[A] = P[A|B]$ happens, then the occurrence of B has no impact on the occurrence of A , and we will say events A and B are mutually independent. Note that $P[A] = P[A|B]$ in conjunction with (1.3) leads to

$$P[AB] = P[A] \cdot P[B] \quad (1.9)$$

This motivates the following definition.

Definition 1.8 *Independence*

- Events A and B are said to be independent from each other if (1.9) holds.
- n events A_1, A_2, \dots, A_n are said to be mutually independent if

$$P[A_1 A_2 \cdots A_n] = P[A_1] \cdot P[A_2] \cdots P[A_n] \quad \blacksquare \quad (1.10)$$

1.1.3 Random Variables, Distribution Function, and Probability Density

A. Random Variables

Let E be a random experiment associated with sample space Ω . Corresponding to each outcome ω , suppose there is a real-valued function $\xi(\omega)$. In probability theory, one is concerned not only with the value of $\xi(\omega)$, but also with the probability of $\xi(\omega)$ taking certain values.

Example 1.8

In a junior-high school a student (that is an ω) is randomly selected, whose height is recorded as $\xi(\omega)$. One of the valid questions in this case would be “what is the probability of the height no greater than x cm?”, namely, what is $P[\xi(\omega) \leq x]$? \blacksquare

Evidently, to address the question in the associated probability space (Ω, \mathcal{F}, P) , it is necessary to assure that $\xi(\omega) \leq x$ belongs to σ -algebra \mathcal{F} so that $P[\xi(\omega) \leq x]$ is well defined. In probability theory functions of this kind are called *random variables*.

Definition 1.9 *Random variables*

Let $\xi(\omega)$ be a real-valued function defined over sample space Ω that is associated with a probability space (Ω, \mathcal{F}, P) . If, for any real x , $(\xi(\omega) \leq x)$ is an event, i.e., $(\xi(\omega) \leq x) \in \mathcal{F}$, then $\xi(\omega)$ is called a random variable. \blacksquare

The introduction of random variables was a major event in the development of modern probability theory as the concept made it possible to extend the probability theory to include studies of random variables and a variety of related issues. One of the issues is *distribution function* of a random variable.

B. Distribution Function and Probability Density

Definition 1.10 *Distribution function of a random variable*

Let $\xi(\omega)$ be a random variable, the distribution function of ξ , denoted by $\Phi_\xi(x)$, is defined as the probability of event $(\xi(\omega) \leq x)$. Namely,

$$\Phi_\xi(x) = P[\xi(\omega) \leq x] \quad \blacksquare \quad (1.11)$$

Several basic properties of distribution functions follow:

- $\Phi_\xi(x)$ is monotonically non-decreasing, i.e.,

$$x_2 \geq x_1 \text{ implies that } \Phi_\xi(x_2) \geq \Phi_\xi(x_1) \quad (1.12)$$

- $$\Phi_\xi(-\infty) = 0, \Phi_\xi(\infty) = 1 \quad (1.13)$$

- $$P[x_1 < \xi(\omega) \leq x_2] = \Phi_\xi(x_2) - \Phi_\xi(x_1) \quad (1.14)$$

- $$P[\xi(\omega) > x] = 1 - \Phi_\xi(x) \quad (1.15)$$

- $$P[\xi(\omega) = x] = \Phi_\xi(x) - \Phi_\xi(x-0) \quad (1.16)$$

where $\Phi_\xi(x-0) = \lim_{a \uparrow x} \Phi_\xi(a)$.

Types of distribution functions

There are several types of distribution functions, with the most important being *discrete type* and *continuous type*.

Discrete type distribution

A discrete random variable ξ takes on discrete values x_0, x_1, x_2, \dots (finite or countable infinite), with probabilities p_0, p_1, p_2, \dots , respectively. Thus a discrete distribution is characterized by a two-row matrix, known as *density matrix*, with finite or countable infinite columns as

$$\xi \Leftrightarrow \begin{bmatrix} x_0 & x_1 & x_2 & \cdots \\ p_0 & p_1 & p_2 & \cdots \end{bmatrix} \quad (1.17)$$

Obviously we have $p_i \geq 0$ and $\sum_i p_i = 1$. The distribution function in this case is given by

$$F(x) = \sum_{i: x_i \leq x} p_i \quad (1.18)$$

which is a staircase function that is continuous from right, see Fig. 1.1.

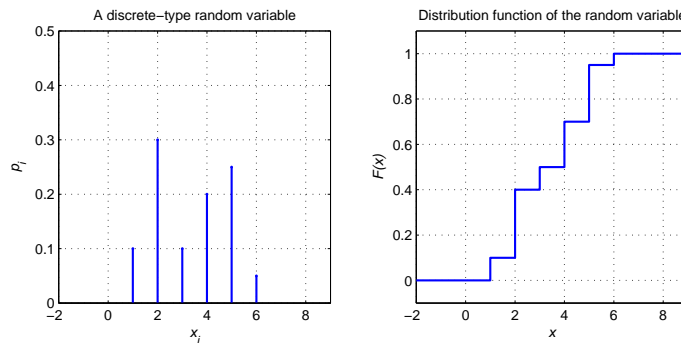


Fig. 1.1 Left: A discrete-type random variable with $x_0 = 1, x_1 = 2, \dots, x_5 = 6$ and $p_0 = 0.1, p_1 = 0.3, p_2 = 0.1, p_3 = 0.2, p_4 = 0.25$, and $p_5 = 0.05$; Right: The distribution function of the random variable.

Example 1.9 *Bernoulli distribution*

Also known as *binomial distribution*, the Bernoulli distribution is discrete and assumes the form of

$$\begin{bmatrix} 0 & 1 & \cdots & n \\ p_0 & p_1 & \cdots & p_n \end{bmatrix}$$

where

$$p_k = \binom{n}{k} p^k q^{n-k} \quad p \geq 0, q \geq 0, p + q = 1 \quad (1.19)$$

with

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

being the number of combinations of taking k elements out of n elements. Note that p_k in (1.19) are exactly the coefficient of polynomial $(p + q)^n$, hence the name of binomial distribution. The reason it also bears the name of Bernoulli distribution has to do with the well-known *Bernoulli trial* which is a random experiment with exactly two possible outcomes: "success" and "failure". Denote the probability of success and failure by p and q , respectively, and suppose the experiment is repeated n time independently. It is obvious that there are 2^n possible outcomes; the number of outcomes with k successes and $(n - k)$ failures is equal to $\binom{n}{k}$; and the probability

of "one such event occurs" is $p^k q^{n-k}$. Therefore, the probability of having k successes and $(n - k)$ failures is precisely given by p_k as seen in (1.19). ■

Continuous type distribution

A distribution of random variable ξ is said to be of continuous type if its distribution function assumes the form

$$\Phi_\xi(x) = \int_{-\infty}^x \varphi_\xi(y) dy \quad (1.20)$$

where $\varphi_\xi(x) \geq 0$ is called *density function* or *density*. Since $\Phi_\xi(\infty) = 1$, we have

$$\int_{-\infty}^{\infty} \varphi_\xi(y) dy = 1$$

Example 1.10 *Uniform distribution*

Let $a < b$. Consider the random variable ξ that takes a value over the interval $[a, b]$ uniformly randomly. The associated distribution, called uniform distribution, is of continuous type with the density function given by

$$\varphi_{\xi}(x) = \begin{cases} \frac{1}{b-a} & \text{if } x \in [a, b] \\ 0 & \text{if } x \notin [a, b] \end{cases} \quad (1.21)$$

Given the density in (1.21), the distribution function is found to be

$$\Phi_{\xi}(x) = \int_{-\infty}^x \varphi_{\xi}(y) dy = \begin{cases} 0 & \text{if } x < a \\ \frac{x-a}{b-a} & \text{if } a \leq x \leq b \\ 1 & \text{if } x > b \end{cases} \quad (1.22)$$

See Figs. 1.2 and 1.3 for plots of the density $\varphi_{\xi}(x)$ and distribution function $\Phi_{\xi}(x)$. ■

C. Expectation, Variance and Covariance

An important operation involving random variables is that of finding weighted average of a random variable ξ or a (measurable) function of ξ .

Density function of uniform distribution

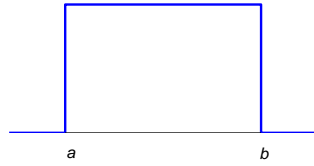


Figure 1.2

Distribution function of uniform random variable

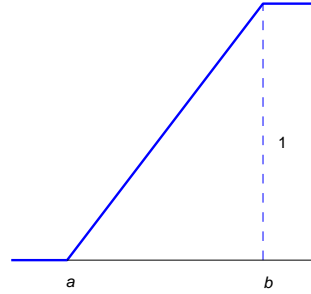


Figure 1.3

Definition 1.11 *Expectation*

The average of a random variable ξ under its probability distribution is called the expectation of

ξ , denoted by $E[\xi]$.

For a discrete type random variable ξ that takes values x_0, x_1, x_2, \dots with probability p_0, p_1, p_2, \dots (recall the density matrix given by (1.17), namely

$$\begin{bmatrix} x_0 & x_1 & x_2 & \cdots \\ p_0 & p_1 & p_2 & \cdots \end{bmatrix}$$

), the expectation of ξ , also known as mean value of ξ , is defined as

$$E[\xi] = \sum_i x_i p_i$$

Since the random variable ξ here takes values x_i , for the sake of notation convenience we denote the random variable by x , and write the above expression as

$$E[x] = \sum_i x_i p_i \quad (1.23)$$

For a continuous type random variable ξ with density $\varphi_\xi(x)$, its expectation is defined as

$$E[x] = \int_{-\infty}^{\infty} x \varphi_\xi(x) dx \quad (1.24)$$

The concept of expectation can be extended to expectation of a (measurable) function, say $g(\xi)$, of random variable ξ as

$$E[g(x)] = \sum_i g(x_i) p_i \quad (1.25)$$

for discrete-type case, or

$$E[g(x)] = \int_{-\infty}^{\infty} g(x) \varphi_\xi(x) dx \quad (1.26)$$

for continuous-type case. ■

Definition 1.12 *Variance*

The variance of a function g of random variable ξ provides a measure of how much there is in g around its mean value $E[g(x)]$:

$$\text{var}[g(x)] = E\left[\left(g(x) - E[g(x)]\right)^2\right] = E[g(x)^2] - \left(E[g(x)]\right)^2 \quad (1.27)$$

In particular, with $g(x) = x$ (1.27) yields the variance of a random variable ξ itself as

$$\text{var}[x] = E[x^2] - \left(E[x]\right)^2 \quad \blacksquare \quad (1.28)$$

Definition 1.13 *Covariance*

Covariance is concerned with two random variables x and y and quantifies the extent to which x and y vary together, namely,

$$\text{cov}[x, y] = E_{x,y}\left[(x - E[x])(y - E[y])\right] = E_{x,y}[xy] - E[x]E[y] \quad (1.29)$$

where $E_{x,y}$ denotes expectation with respect to joint density $\varphi(x, y)$, namely,

$$E_{x,y}[xy] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy \varphi(x, y) dx dy \quad \blacksquare$$

D. Independence of Random Variables and Conditional Distribution

Definition 1.14 *Independence of random variables*

Random variables $\xi_1(\omega), \dots, \xi_n(\omega)$ with distribution functions $\Phi_1(x_1), \dots, \Phi_n(x_n)$ are said to be *mutually independent* if

$$\Phi(x_1, \dots, x_n) = \Phi_1(x_1) \cdots \Phi_n(x_n) \quad (1.30)$$

holds, where $\Phi(x_1, \dots, x_n)$ is the joint distribution of $\xi_1(\omega), \dots, \xi_n(\omega)$, namely,

$$\Phi(x_1, \dots, x_n) = P[\xi_1 \leq x_1, \dots, \xi_n \leq x_n] \quad \blacksquare$$

From (1.30), it follows that

$$P[\xi_1 \leq x_1, \dots, \xi_n \leq x_n] = P[\xi_1 \leq x_1] \cdots P[\xi_n \leq x_n] = \prod_{i=1}^n P[\xi_i \leq x_i] \quad (1.31)$$

In words, if random variables ξ_1, \dots, ξ_n are mutually independent, the joint probability of the event $(\xi_1 \leq x_1, \dots, \xi_n \leq x_n)$ is equal to product of the probability of individual event $(\xi_i \leq x_i)$.

For continuous type of random variables, (1.30) implies that $\xi_1(\omega), \dots, \xi_n(\omega)$ are mutually independent if the joint probability density equals the product of the individual probability density almost everywhere (a.e.), i.e.,

$$\varphi(x_1, \dots, x_n) = \varphi_1(x_1) \cdots \varphi_n(x_n) \quad (\text{a.e.}) \quad (1.32)$$

Definition 1.15 *Conditional distribution*

For two discrete-type random variables ξ and η with

$$\xi \Leftrightarrow \begin{bmatrix} x_0 & x_1 & x_2 & \cdots \\ p_0 & p_1 & p_2 & \cdots \end{bmatrix} \text{ and } \eta \Leftrightarrow \begin{bmatrix} y_0 & y_1 & y_2 & \cdots \\ q_0 & q_1 & q_2 & \cdots \end{bmatrix}$$

where $p_i > 0$ and $q_i > 0$ are assumed, denote the joint probability

$$P[\xi = x_i, \eta = y_j] = p_{i,j}$$

hence we can write

$$p_i = \sum_j p_{i,j} \text{ and } q_j = \sum_i p_{i,j}$$

which in turn implies that

$$P[\xi = x_i | \eta = y_j] = \frac{P[\xi = x_i, \eta = y_j]}{P[\eta = y_j]} = \frac{p_{i,j}}{q_j} = \frac{p_{i,j}}{\sum_i p_{i,j}} \quad (1.33)$$

and

$$P[\eta = y_j | \xi = x_i] = \frac{p_{i,j}}{p_i} = \frac{p_{i,j}}{\sum_j p_{i,j}} \quad (1.34)$$

The conditional distribution refers to the probability of event $(\xi \leq x \text{ given } \eta = y_j)$ or that of $(\eta \leq y \text{ given } \xi = x_i)$. By using (1.33) and (1.34), we obtain

$$P[(\xi \leq x | \eta = y_j)] = \frac{\sum_{i: x_i \leq x} p_{i,j}}{\sum_i p_{i,j}} \quad (1.35)$$

and

$$P[\eta \leq y | \xi = x_i] = \frac{\sum_{j: y_j \leq y} p_{i,j}}{\sum_j p_{i,j}} \quad (1.36)$$

From (1.35) and (1.36), we see that conditional distributions involving two discrete-type random variables can be expressed (and evaluated) using their joint probability distribution $p_{i,j}$.

For two continuous-type random variables ξ and η with joint probability density $\varphi(x, y)$, the conditional distributions are given by

$$P[\xi \leq x | \eta = y] = \frac{\int_{-\infty}^x \varphi(z, y) dz}{\int_{-\infty}^{\infty} \varphi(z, y) dz} \quad (1.37)$$

and

$$P[\eta \leq y | \xi = x] = \frac{\int_{-\infty}^y \varphi(x, z) dz}{\int_{-\infty}^{\infty} \varphi(x, z) dz} \quad (1.38)$$

By writing (1.37) as

$$P[\xi \leq x | \eta = y] = \int_{-\infty}^x \left(\frac{\varphi(z, y)}{\int_{-\infty}^{\infty} \varphi(z, y) dz} \right) dz$$

one may interpolate the integrand in the above expression as the *conditional distribution density* and denote it as $\varphi(x | y)$:

$$\varphi(x | y) = \frac{\varphi(x, y)}{\int_{-\infty}^{\infty} \varphi(z, y) dz} \quad (1.39)$$

In this way, (1.37) becomes

$$P[\xi \leq x | \eta = y] = \int_{-\infty}^x \varphi(z | y) dz$$

Similarly, by defining the conditional distribution density

$$\varphi(y | x) = \frac{\varphi(x, y)}{\int_{-\infty}^{\infty} \varphi(x, z) dz} \quad (1.40)$$

(1.38) can be expressed as

$$P[\eta \leq y | \xi = x] = \int_{-\infty}^y \varphi(z | x) dz \quad \blacksquare$$

1.1.4 Gaussian Distribution

Also known as the *normal distribution*, the Gaussian distribution played an extremely important role in the development of probability theory, and is arguably the most useful among all probabilistic distributions. The life span of bulbs produced under practically same manufacturing conditions, for example, obeys a Gaussian distribution. This is also true for several other measures of products that are manufactured in quantity under the same conditions. Gaussian distribution is also encountered in many natural, biological, and social events/phenomena: velocities of gas molecules; errors in measuring a physical object; heights/weights of biological species, yearly precipitations of a certain city, etc. A common thread of these random variables is that they are cumulative synthesis of many small (i.e. of minor importance), independent random components.

A. One-Dimensional Gaussian Distribution

Definition 1.16 Gaussian distribution

A one-dimensional Gaussian distribution is a continuous-type distribution whose density is given by

$$\varphi_{\mu, \sigma}(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2 / 2\sigma^2} \quad (1.41)$$

where μ and σ are two real-valued parameters. Hence the Gaussian distribution function is given by

$$\Phi_{\mu, \sigma}(x) = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^x e^{-(y-\mu)^2 / 2\sigma^2} dy \quad (1.42)$$

A common notation for the one-dimensional normal distribution is $\mathcal{N}(x, \mu, \sigma^2)$. ■

By (1.24) and (1.28), the expectation and variance of a Gaussian random variable are found to be

$$E[x] = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\infty} x e^{-(x-\mu)^2 / 2\sigma^2} dx = \mu \quad (1.43)$$

and

$$\text{var}[x] = E[x^2] - E[x]^2 = \sigma^2 \quad (1.44)$$

respectively. In other words, we see that the Gaussian distribution is characterized by its expectation and variance. Fig. 1.4 illustrates the meaning of μ as the mean value of x using two Gaussian density functions with different μ 's, while Fig. 1.5 illustrates the meaning of σ^2 as the variance of x by several Gaussian density functions with different σ 's.

Example 1.11

Let x be a Gaussian random variable with $\mu = 2$ and $\sigma = 1$. Compute the probability of x being in between $\mu - 3$ and $\mu + 3$ (i.e. x falls into the interval $[-1, 5]$).

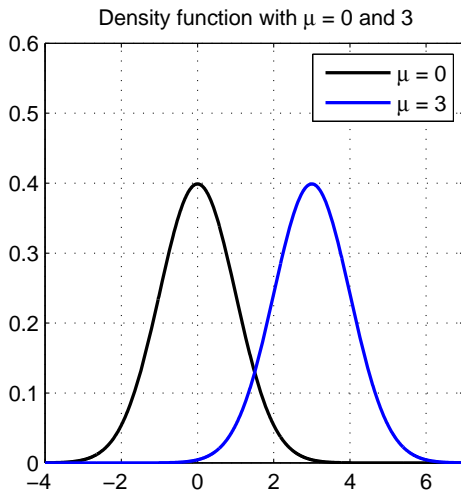


Figure 1.4

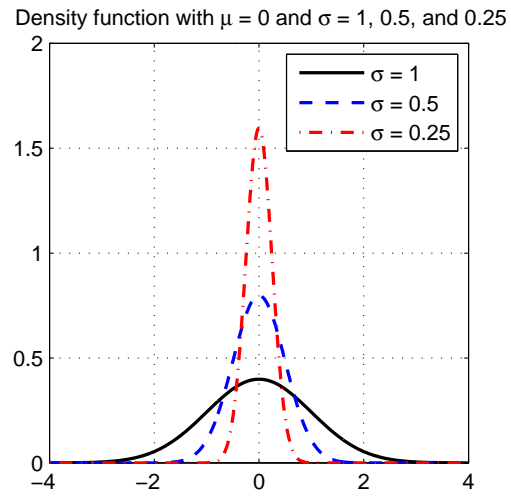


Figure 1.5

Solution

From (1.14) and (1.42) it follows that

$$P[-1 \leq \xi \leq 5] = \frac{1}{\sqrt{2\pi}} \int_{-1}^5 e^{-(y-2)^2/2} dy = \frac{1}{\sqrt{2\pi}} \int_{-3}^3 e^{-z^2/2} dz \approx 99.7\% \quad \blacksquare$$

B. Multidimensional Gaussian Distribution

Let us consider n random variables $\xi_1(\omega), \dots, \xi_n(\omega)$. The joint distribution of these random variables is defined by

$$\Phi(x_1, x_2, \dots, x_n) = P[\xi_1(\omega) \leq x_1, \xi_2(\omega) \leq x_2, \dots, \xi_n(\omega) \leq x_n] \quad (1.45)$$

If $\xi_1(\omega), \dots, \xi_n(\omega)$ are of continuous type, then $\Phi(x_1, x_2, \dots, x_n)$ can be expressed as

$$\Phi(\mathbf{x}) = \int_{-\infty}^{x_1} \dots \int_{-\infty}^{x_n} \varphi(y_1, \dots, y_n) dy_n \dots dy_1 \quad (1.46)$$

where $\mathbf{x} = [x_1 \ x_2 \ \cdots \ x_n]^T$ and $\varphi(x_1, \dots, x_n)$ is the density function.

As expected, the most important continuous-type multidimensional distribution is Gaussian distribution whose density in matrix notation is given by

$$\varphi(\mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{n/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\} \quad (1.47)$$

where $\boldsymbol{\Sigma}$ is a symmetric and positive definite matrix of size n by n , called *covariance matrix*, and $|\boldsymbol{\Sigma}|$ denotes the determinant of $\boldsymbol{\Sigma}$. Oftentimes the density of normal distribution is denoted by

$\mathcal{N}(\mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$. It can be verified that the expectation of Gaussian $\boldsymbol{\xi}$ is $\boldsymbol{\mu}$, and the covariance of between the individual components (as random variables themselves) of $\boldsymbol{\xi}$ is given by $\boldsymbol{\Sigma}$.

An important special case of multidimensional Gaussian distribution is when the individual random variables are mutually independent, the density function in this case becomes separable as

$$\varphi(x_1, \dots, x_n) = \varphi_1(x_1) \cdots \varphi_n(x_n)$$

with each $\varphi_i(x_i)$ being a one-dimensional Gaussian density, i.e.,

$$\varphi_i(x_i) = \frac{1}{\sigma_i \sqrt{2\pi}} e^{-(x_i - \mu_i)^2 / 2\sigma_i^2}$$

Therefore the probability density of n mutually independent Gaussian distribution is given by

$$\varphi(x_1, \dots, x_n) = \prod_{i=1}^n \frac{1}{\sigma_i \sqrt{2\pi}} e^{-(x_i - \mu_i)^2 / 2\sigma_i^2} = \frac{1}{(2\pi)^{n/2}} \frac{1}{\sigma_1 \cdots \sigma_n} \exp\left(-\frac{1}{2} \sum_{i=1}^n \frac{(x_i - \mu_i)^2}{\sigma_i^2}\right) \quad (1.48)$$

It can be readily verified that (1.47) coincides with (1.48) by assigning $\boldsymbol{\mu} = [\mu_1 \ \cdots \ \mu_n]$ and

$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_n^2 \end{bmatrix}$$

In other words, (1.47) becomes joint density of n independent Gaussian distributions if and only if $\boldsymbol{\Sigma}$ is *diagonal*.

1.1.5 Likelihood Function and Log-Likelihood

Consider a continuous-type distribution with density $\varphi(x, \mathbf{w})$ where vector \mathbf{w} collects the parameters involved. Let \mathcal{D} be a data set $\mathcal{D} = \{x_1, x_2, \dots, x_N\}$ that are drawn independently from the same distribution. In literature, a data set of this type is said to be *independently and identically distributed* (i.i.d.). Because of its probabilistic independence, given parameter \mathbf{w} the probability

density of an i.i.d. data set \mathcal{D} assumes the form

$$p(\mathcal{D} | \mathbf{w}) = \prod_{i=1}^N \varphi(x_i, \mathbf{w}) \quad (1.49)$$

which is called *likelihood function* of the probability distribution. A popular approach for determining the parameters of a probability distribution using an observed data is to maximize the likelihood function with respect to the parameters. Because logarithm is a monotonically increasing function and logarithm of the likelihood function simplifies subsequent mathematical analysis, one maximizes the log-likelihood instead:

$$\log p(\mathcal{D} | \mathbf{w}) = \sum_{i=1}^n \log \varphi(x_i, \mathbf{w}) \quad (1.50)$$

Example 1.12

Consider the 1-D Gaussian distribution in (1.41), i.e.

$$\varphi_{\mu, \sigma}(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2}$$

Following (1.50), its log-likelihood function for an observed i.i.d. data set, denoted by $L(\mu, \sigma^2, \mathcal{D})$, is given by

$$L(\mu, \sigma^2, \mathcal{D}) = -\frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2 - \frac{N}{2} \log \sigma^2 - \frac{N}{2} \log(2\pi)$$

Maximizing $L(\mu, \sigma^2, \mathcal{D})$ with respect to $\{\mu, \sigma^2\}$ yields *maximum likelihood* (ML) estimates

$$\mu_{ML} = \frac{1}{N} \sum_{i=1}^N x_i \quad (1.51)$$

and

$$\sigma_{ML}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu_{ML})^2 \quad \blacksquare \quad (1.52)$$

1.2 Optimization Methods

1.2.1 Unconstrained Optimization [5]

Unconstrained optimization studies problems of the form

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{x}) \quad (1.53)$$

where $f(\mathbf{x})$ is a real-valued twice continuously differentiable function.

- Gradient of $f(\mathbf{x})$ is defined by

$$\nabla f(\mathbf{x}) = \mathbf{g}(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$$

- Hessian of $f(\mathbf{x})$ is defined by

$$\nabla^2 f(\mathbf{x}) = \mathbf{H}(\mathbf{x}) = \left\{ \frac{\partial^2 f}{\partial x_i \partial x_j} \right\}_{i,j=1,2,\dots,n}$$

Note that Hessian is always square and symmetric.

- A point \mathbf{x} is called a *stationary point* if $\nabla f(\mathbf{x}) = \mathbf{0}$. A geometric interpretation of the concept in the case of $n = 1$ or 2 is that \mathbf{x} is a stationary point if the tangent (plane) of $f(\mathbf{x})$ at \mathbf{x} is in parallel with the \mathbf{x} -axis (plane).

A. First-Order Necessary Condition

If \mathbf{x}^* is a local minimum point (minimizer), then

$$\nabla f(\mathbf{x}^*) = \mathbf{0} \quad (1.54)$$

In other words, if \mathbf{x}^* is a minimum point, then it must be a stationary point.

◇ A geometric interpretation of the 1st-order necessary condition when $n = 1$ or 2 is that the tangent plane of $f(\mathbf{x})$ at a minimizer must be in parallel to the \mathbf{x} -axis (plane).

◇ $\nabla f(\mathbf{x}^*) = \mathbf{0}$ is not sufficient for \mathbf{x}^* to be a local minimizer.

B. Second-Order Necessary Conditions

If \mathbf{x}^* is a local minimum point (minimizer), then

$$\nabla f(\mathbf{x}^*) = \mathbf{0} \quad \text{and} \quad \nabla^2 f(\mathbf{x}^*) \succeq \mathbf{0} \quad (1.55)$$

◇ $\nabla^2 f(\mathbf{x}^*) \succeq \mathbf{0}$ (i.e. the Hessian $\nabla^2 f(\mathbf{x}^*)$ is positive semidefinite if and only if the eigenvalues of $\nabla^2 f(\mathbf{x}^*)$ are nonnegative.

◇ $\nabla f(\mathbf{x}^*) = \mathbf{0}$ and $\nabla^2 f(\mathbf{x}^*) \succeq \mathbf{0}$ are not sufficient for \mathbf{x}^* to be a local minimizer.

C. Second-Order Sufficient Conditions

Point \mathbf{x}^* is a local minimizer, if

$$\nabla f(\mathbf{x}^*) = \mathbf{0} \quad \text{and} \quad \nabla^2 f(\mathbf{x}^*) \succ \mathbf{0} \quad (1.56)$$

D. Convex Sets and Convex Functions

- A set $S \subseteq \mathbb{R}^n$ is said to be convex if

$$\mathbf{x}_1, \mathbf{x}_2 \in S \Rightarrow \alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2 \in S \quad \text{for } 0 \leq \alpha \leq 1 \quad (1.57)$$

- A function $f(\mathbf{x})$ is said to be convex over convex set S if

$$f(\alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2) \leq \alpha f(\mathbf{x}_1) + (1 - \alpha) f(\mathbf{x}_2) \quad (1.58)$$

for any $\mathbf{x}_1, \mathbf{x}_2 \in S$ and $0 \leq \alpha \leq 1$.

- A function $f(\mathbf{x}) \in C^1$ (i.e. continuously differentiable) is convex if and only if

$$f(\mathbf{x}_1) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{x}_1 - \mathbf{x}) \quad (1.59)$$

for any $\mathbf{x}, \mathbf{x}_1 \in S$. Property (1.59) is an important one as it actually characterizes convex functions and has a geometrical interpretation – all tangents of $f(\mathbf{x})$ lie below the function graph $(\mathbf{x}, f(\mathbf{x}))$ as illustrated in Fig. 1.6.

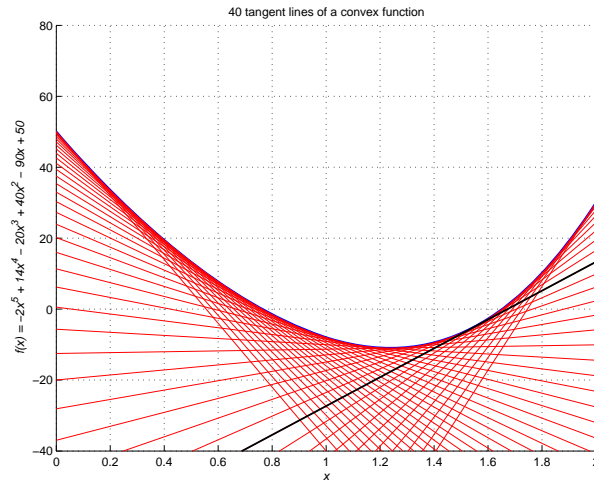


Figure 1.6 All supporting tangents of a convex function lie underneath it.

- A function $f(\mathbf{x}) \in C^2$ (i.e. twice continuously differentiable) is convex over a convex set S (which contains more than one point) if and only if $\nabla^2 f(\mathbf{x}) \succeq \mathbf{0}$ for $\mathbf{x} \in S$.

- If function $f(\mathbf{x}) \in C^1$ is convex over a convex set S and \mathbf{x}^* is a point in S that satisfies $\nabla f(\mathbf{x}^*) = \mathbf{0}$, then \mathbf{x}^* is a global minimizer in S . This is because for any $\mathbf{x} \in S$ we have

$$f(\mathbf{x}) \geq f(\mathbf{x}^*) + \nabla f(\mathbf{x}^*)^T (\mathbf{x} - \mathbf{x}^*) = f(\mathbf{x}^*)$$

◊ In other words, for smooth convex functions, the first-order necessary condition $\nabla f(\mathbf{x}^*) = \mathbf{0}$ is sufficient for \mathbf{x}^* to be a minimizer.

E. A General Structure of Unconstrained Optimization Algorithms

Step 1: Choose an initial point \mathbf{x}_0 and a convergence tolerance ε . Set a counter $k = 0$ for the number of iterations.

Step 2: Determine a search direction \mathbf{d}_k for reducing $f(\mathbf{x})$ from $f(\mathbf{x}_k)$.

Step 3: Determine a step size $\alpha_k \geq 0$ such that $f(\mathbf{x}_k + \alpha \mathbf{d}_k)$ is minimized, and construct

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k.$$

Step 4: If $\|\alpha_k \mathbf{d}_k\| < \varepsilon$, stop and output a solution \mathbf{x}_{k+1} , otherwise set $k := k + 1$ and repeat from Step 2.

F. A Backtracking Algorithm for Line Search (BLS) [6]

Assume \mathbf{d}_k is a descent direction of $f(\mathbf{x})$ at \mathbf{x}_k .

Step 1: Select constants $\rho \in (0, 0.5)$ and $\gamma \in (0, 1)$ (e.g. $\rho = 0.1$, $\gamma = 0.5$). Set $\alpha = 1$.

Step 2: While

$$f(\mathbf{x}_k + \alpha \mathbf{d}_k) > f(\mathbf{x}_k) + \rho \alpha \mathbf{g}_k^T \mathbf{d}_k \tag{1.60}$$

Set $\alpha = \gamma \alpha$.

Step 3: Output $\alpha_k = \alpha$.

◊ How BLS works: the right-hand side of (1.60) is a line in α passing through point $(0, f(\mathbf{x}_k))$ with a negative slope because \mathbf{d}_k is a descent direction:

$$\rho \mathbf{g}_k^T \mathbf{d}_k = \rho \left. \frac{df(\mathbf{x}_k + \alpha \mathbf{d}_k)}{d\alpha} \right|_{\alpha=0} < 0$$

See Fig. 1.7 for an illustration. Note from the figure that a value of α less than α_0 may be

considered acceptable because such an α satisfies

$$f(\mathbf{x}_k + \alpha \mathbf{d}_k) < f(\mathbf{x}_k) + \rho \alpha \mathbf{g}_k^T \mathbf{d}_k < f(\mathbf{x}_k)$$

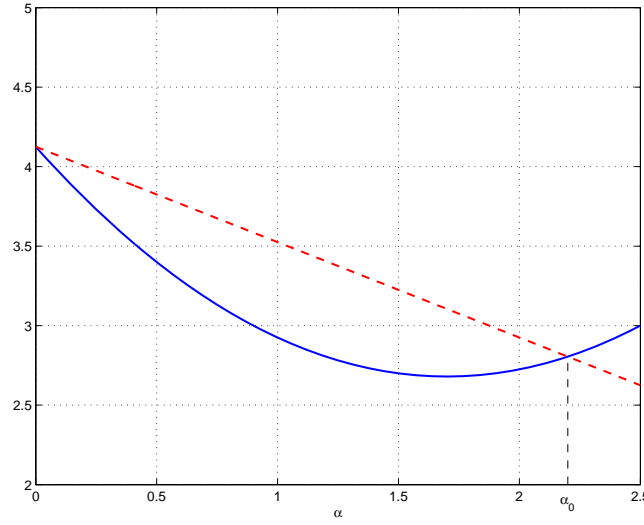


Figure 1.7 Both $f(\mathbf{x}_k + \alpha \mathbf{d}_k)$ (blue curve) and $f(\mathbf{x}_k) + \rho \alpha \mathbf{g}_k^T \mathbf{d}_k$ (red dashed line) are shown as functions of α . An acceptable α is found by gradually reducing α to a value less than α_0 (which is signified by $f(\mathbf{x}_k + \alpha \mathbf{d}_k) < f(\mathbf{x}_k) + \rho \alpha \mathbf{g}_k^T \mathbf{d}_k$).

- Two MATLAB functions, namely `bt_lsearch.m` and `inex_lsearch.m`, for line search are available from the link below:

<http://www.ece.uvic.ca/~wslu/Talk.html>

G. The Steepest Descent Method (SDM)

- The SDM was originated by A. Cauchy (1848). It remains to be one of the most popular minimization techniques because its simplicity although the algorithm is usually slower relative to another popular algorithm –Newton’s algorithm.
- The SMD follows the general algorithmic structure described above, where the search direction is computed as

$$\mathbf{d}_k = -\nabla f(\mathbf{x}_k) \quad (1.61)$$

This choice of \mathbf{d}_k may be understood by the Taylor expansion of $f(\mathbf{x})$ at \mathbf{x}_k up to its 1st-order:

$$f(\mathbf{x}_k + \mathbf{d}) \approx f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{d} \quad (\text{assuming } \|\mathbf{d}\| \text{ is small})$$

From which we see the value of $f(\mathbf{x})$ is reduced at highest rate at $\mathbf{x}_k + \mathbf{d}_k$ if $\mathbf{d}_k = -\nabla f(\mathbf{x}_k)$ is taken.

H. Why Is SDM Slow?

The slow convergence of the SDM has to do with the “condition” of the Hessian at \mathbf{x}_k : the progress made by the k th SDM iteration relative to the preceding progress can be estimated as

$$f(x_{k+1}) - f(x^*) \leq \left(\frac{1-r}{1+r} \right)^2 (f(x_k) - f(x^*)) \quad (1.62)$$

where r is the ratio of the smallest over largest eigenvalues of the Hessian at \mathbf{x}_k (which is denoted by \mathbf{H}_k). From (1.62) the SDM is expected to be slow when ratio r is small, or stated in another way, when the *condition number* of \mathbf{H}_k is large.

◇ The condition number of a square nonsingular matrix is the ratio of its largest over smallest eigenvalues. A square nonsingular matrix is said to be *ill-conditioned* if its condition number is very large.

I. Scaling and Normalization

The eigenvalues of Hessian of an objective function and, in turn, the performance of the steepest-descent method tend to depend to a large extent on the choice of variables. For example, in one and the same two-dimensional problem, whether the contours look nearly circular or extremely elliptical has to do with how the units of design components $\{x_i, i = 1, \dots, n\}$ are chosen. Consequently, the rate of convergence can often be improved by scaling the variables properly through variable transformation. A possible approach to scaling is to let

$$\mathbf{x} = \mathbf{T}\mathbf{y}$$

where \mathbf{T} is an $n \times n$ diagonal matrix, and then solve the problem

$$\underset{\mathbf{y}}{\text{minimize}} \quad h(\mathbf{y}) = f(\mathbf{x})|_{\mathbf{x}=\mathbf{T}\mathbf{y}}$$

The gradient and Hessian of the new problem are calculated as

$$\mathbf{g}_h = \mathbf{T}\mathbf{g}_x \quad \text{and} \quad \mathbf{H}_h = \mathbf{T}^T \mathbf{H} \mathbf{T} \quad (1.63)$$

respectively. We note that both the steepest-descent direction and the eigenvalues associated with the transformed problem are altered, hence the performance of an SDM is expected to be improved by adequate scaling.

The choice of \mathbf{T} tends to depend heavily on the problem at hand and, as a result, no general rules can be stated. For many problems in machine learning, however, the variables are always related to real-world quantities (e.g. number of rooms in a typical house, age of a house, total living area of a house in feet², etc.) and are known to vary in certain ranges. In cases like these, it is often beneficial to normalize each variable x , assuming to be in the range $[a, b]$, to a standard range say $[-0.5, 0.5]$, by linear transform

$$\bar{x} = \frac{x - \mu}{b - a} \quad \text{with} \quad \mu = \frac{a + b}{2} \quad (1.64)$$

J. Newton's Method

Newton's method is well known for its fast convergence at a cost of increased complexity.

- Newton's method also follows the general algorithmic structure as outlined in Sec. 1.2.1.E, where the search direction is computed as

$$\mathbf{d}_k = -(\nabla^2 f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k), \quad \text{i.e.,} \quad \mathbf{d}_k = -\mathbf{H}_k^{-1} \mathbf{g}_k \quad (1.65)$$

◇ This choice of \mathbf{d}_k may be understood as the minimizer of the 2nd-order approximate of $f(\mathbf{x})$ in a small vicinity of \mathbf{x}_k :

$$f(\mathbf{x}_k + \mathbf{d}) \approx f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T \nabla^2 f(\mathbf{x}_k) \mathbf{d}$$

Figure 1.8 illustrates the first two Newton steps for the convex function

$$f(x) = x^2 + 0.05e^{-x} + 8$$

With $x_0 = -5$, the first two Newton iterations yield $x_1 = -3.15$ and $x_2 = -0.7929$ while the global minimizer is at $x^* = 0.0242$.

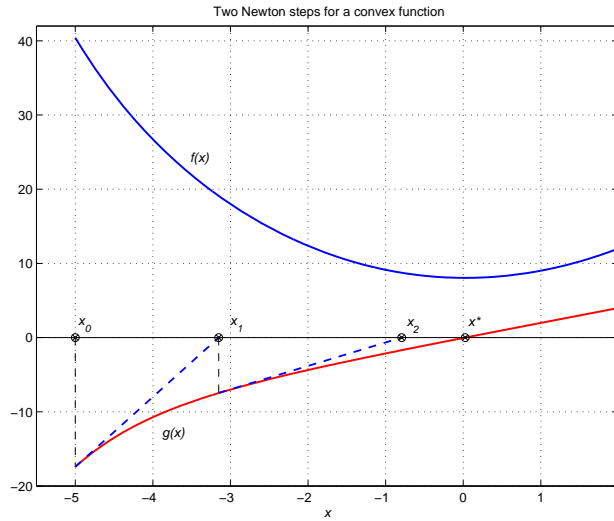


Figure 1.8 Two Newton steps for a convex function.

- It is important to note that Newton's direction given by (1.65) is effective only when \mathbf{H}_k is *positive definite*. In case \mathbf{H}_k is not positive definite, the search direction in Newton's method is modified to

$$\mathbf{d}_k = -\hat{\mathbf{H}}_k^{-1} \mathbf{g}_k$$

where $\hat{\mathbf{H}}_k = \frac{\mathbf{H}_k + \beta \mathbf{I}}{1 + \beta}$ with a sufficiently large β so that $\mathbf{H}_k + \beta \mathbf{I} \succ \mathbf{0}$.

K. Quasi-Newton Methods

One may take a unifying look at the SDM and Newton method by writing their search direction as

$$\mathbf{d}_k = -\mathbf{S}_k \mathbf{g}_k \quad (1.66)$$

where

$$S_k = \begin{cases} \mathbf{I} & \text{for steepest descent} \\ \mathbf{H}_k^{-1} & \text{for Newton} \end{cases}$$

Newton's method is known for its fast convergence and solution accuracy. However, this fast convergence is achieved at the cost of increased computational complexity which has largely to do with the computation of the inverse Hessian in every Newton iteration. Quasi-Newton algorithms are developed to provide convergent rates comparable with that of Newton algorithm with reduced complexity. Actually quasi-Newton algorithms are based only on gradient information and do not require explicit evaluation of the Hessian and its inverse.

Quasi-Newton algorithms follow the general algorithmic structure described earlier, where search direction \mathbf{d}_k is evaluated using (1.66) where matrix S_k is updated in each iteration using gradient information. Two most well-known quasi-Newton algorithms are those developed by Davidon, Fletcher, and Powell (DFP), and by Broyden, Fletcher, Goldfarb, and Shanno (BFGS), respectively. The BFGS and DFP formulas for updating matrix S_k are given below.

◇ Quantities required: $\boldsymbol{\delta}_k = \mathbf{x}_{k+1} - \mathbf{x}_k$, $\boldsymbol{\gamma}_k = \mathbf{g}_{k+1} - \mathbf{g}_k$

◇ Davidon-Fletcher-Powell (DFP) updating formula: $S_0 = \mathbf{I}$, and

$$S_{k+1} = S_k + \frac{\boldsymbol{\delta}_k \boldsymbol{\delta}_k^T}{\boldsymbol{\delta}_k^T \boldsymbol{\gamma}_k} - \frac{S_k \boldsymbol{\gamma}_k \boldsymbol{\gamma}_k^T S_k}{\boldsymbol{\gamma}_k^T S_k \boldsymbol{\gamma}_k} \quad (1.67)$$

◇ Broyden-Fletcher-Goldfarb-Shanno (BFGS) updating formula: $S_0 = \mathbf{I}$, and

$$S_{k+1} = S_k + \left(1 + \frac{\boldsymbol{\gamma}_k^T S_k \boldsymbol{\gamma}_k}{\boldsymbol{\gamma}_k^T \boldsymbol{\delta}_k} \right) \frac{\boldsymbol{\delta}_k \boldsymbol{\delta}_k^T}{\boldsymbol{\gamma}_k^T \boldsymbol{\delta}_k} - \frac{\boldsymbol{\delta}_k \boldsymbol{\gamma}_k^T S_k + S_k \boldsymbol{\gamma}_k \boldsymbol{\delta}_k^T}{\boldsymbol{\gamma}_k^T \boldsymbol{\delta}_k} \quad (1.68)$$

1.2.2 Constrained Optimization

A general constrained optimization problem assumes the form

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to :} && a_i(\mathbf{x}) = 0 \quad \text{for } i = 1, 2, \dots, p \\ & && c_j(\mathbf{x}) \leq 0 \quad \text{for } j = 1, 2, \dots, q \end{aligned} \quad (1.69\text{a-c})$$

◇ Feasible region: $\{\mathbf{x} : a_i(\mathbf{x}) = 0 \text{ for } i = 1, 2, \dots, p, \text{ and } c_j(\mathbf{x}) \leq 0 \text{ for } j = 1, 2, \dots, q\}$

◇ A point \mathbf{x} is said to be *feasible* if it is in the feasible region (inside or on the boundary).

◇ An inequality constraint $c_j(\mathbf{x}) \leq 0$ is said to be *active* at a feasible point \mathbf{x} , if $c_j(\mathbf{x}) = 0$. An inequality constraint $c_j(\mathbf{x}) \leq 0$ is said to be *inactive* at a feasible point \mathbf{x} if $c_j(\mathbf{x}) < 0$.

A. Karush-Kuhn-Tucker (KKT) Conditions

The Karush-Kuhn-Tucker (KKT) conditions are a set of *first-order necessary conditions* for point \mathbf{x}^* to be a local minimizer of problem (1.69). These conditions are fundamentally important in the theory and practice of constrained optimization. The KKT conditions can be stated as follows.

KKT Conditions

If \mathbf{x}^* is a local minimizer of the constrained problem (1.69), then the following conditions, known as *KKT conditions*, are satisfied:

$$(a) \quad a_i(\mathbf{x}^*) = 0 \quad \text{for } i = 1, 2, \dots, p \quad (1.70a)$$

$$(b) \quad c_j(\mathbf{x}^*) \leq 0 \quad \text{for } j = 1, 2, \dots, q \quad (1.70b)$$

(c) there exist Lagrange multipliers λ_i^* for $1 \leq i \leq p$ and μ_j^* for $1 \leq j \leq q$ such that

$$\nabla f(\mathbf{x}^*) + \sum_{i=1}^p \lambda_i^* \nabla a_i(\mathbf{x}^*) + \sum_{j=1}^q \mu_j^* \nabla c_j(\mathbf{x}^*) = \mathbf{0} \quad (1.70c)$$

$$(d) \quad \mu_j^* c_j(\mathbf{x}^*) = 0 \quad \text{for } j = 1, 2, \dots, q \quad (1.70d)$$

$$(e) \quad \mu_j^* \geq 0 \quad \text{for } j = 1, 2, \dots, q \quad (1.70e)$$

◇ Conditions (1.70a) and (1.70b) are merely the constraints imposed in the problem, see (1.69b) and (1.69c).

◇ Condition (1.70c) states that at a local minimizer the gradient of the objective function is a linear combination of the gradients of the constraint functions, and the coefficients in the combination are (up to a sign) Lagrange multipliers.

◇ If we define the *Lagrangian* of problem (1.69) as

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{i=1}^p \lambda_i a_i(\mathbf{x}) + \sum_{j=1}^q \mu_j c_j(\mathbf{x}) \quad (1.71)$$

then (1.70c) can be expressed as

$$\nabla_{\mathbf{x}} L(\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) = \mathbf{0} \quad (1.72)$$

◇ Conditions in (1.70d) are called *complementarity conditions*. It follows that if the k th inequality constraint is inactive at \mathbf{x}^* , then $\mu_k^* = 0$. In other words, the equation in (1.70c) only includes those terms $\nabla c_j(\mathbf{x}^*)$ with $c_j(\mathbf{x})$ active at \mathbf{x}^* .

◇ The total number of equations in the KKT conditions, namely p (from (1.70a)) + n (from (1.70c)) + q (from (1.70d)), are equal to the total number of components in unknown vectors $\{\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*\}$.

B. Convex Programming Problems

Definition 1.17 Convex programming

A constrained minimization problem (1.69) is said to be a *convex programming* (CP) problem if the objective function is convex and the feasible region defined by the constraints, i.e., $\{\mathbf{x}: a_i(\mathbf{x}) = 0 \text{ for } i = 1, 2, \dots, p, \text{ and } c_j(\mathbf{x}) \leq 0 \text{ for } j = 1, 2, \dots, q\}$ is a convex set. In words, a CP problem is a problem of minimizing a convex function over a convex region. ■

- It can be shown that the feasible region $\{\mathbf{x}: a_i(\mathbf{x}) = 0 \text{ for } i = 1, 2, \dots, p, \text{ and } c_j(\mathbf{x}) \leq 0 \text{ for } j = 1,$

$2, \dots, q\}$ is convex if (i) the equality constraint functions $a_i(\mathbf{x})$ are *affine*, and (ii) the inequality constraint functions $c_j(\mathbf{x})$ are *convex*.

Hence a general CP problem assumes the form

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to:} && \mathbf{a}_i^T \mathbf{x} = b_i \quad \text{for } i = 1, 2, \dots, p \\ & && c_j(\mathbf{x}) \leq 0 \quad \text{for } j = 1, 2, \dots, q \end{aligned} \quad (1.73\text{a-c})$$

where $f(\mathbf{x})$ and $c_j(\mathbf{x})$ for $j = 1, 2, \dots, q$ are convex.

- An important property of CP problems is that the KKT conditions become *both necessary and sufficient*. Therefore, for CP problems $\{\mathbf{x}^*, \boldsymbol{\lambda}^*, \boldsymbol{\mu}^*\}$ satisfying the KKT conditions gives a local solution (minimizer) \mathbf{x}^* .
- Another important property of CP problems is that any local minimizer of a CP problem is a global minimizer. By combining the above two properties, one concludes that \mathbf{x}^* is a global minimizer of CP problem (1.73) if and only if the KKT conditions are satisfied.

C. Important Classes of CP Problems

• **Linear programming (LP)**

Linear programming problems are of great importance both theoretically and practically. Many problems in science and engineering can be accurately or approximately modeled as LP problems.

The *standard-form* LP problem is given by

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) = \mathbf{c}^T \mathbf{x} \\ & \text{subject to:} && \mathbf{Ax} = \mathbf{b} \\ & && \mathbf{x} \geq \mathbf{0} \end{aligned} \quad (1.74\text{a-c})$$

Another class of popular LP problem known as *alternative-form* LP problems, assumes the form

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) = \mathbf{c}^T \mathbf{x} \\ & \text{subject to:} && \mathbf{Ax} \leq \mathbf{b} \end{aligned} \quad (1.75\text{a-b})$$

• **Quadratic programming (QP)**

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \mathbf{x}^T \mathbf{H} \mathbf{x} + \mathbf{p}^T \mathbf{x} \\ & \text{subject to:} && \mathbf{Ax} \leq \mathbf{b} \end{aligned} \quad (1.76\text{a-b})$$

where \mathbf{H} is symmetric and positive semidefinite ($\mathbf{H} \succeq \mathbf{0}$).

• **QP with quadratic constraints**

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \mathbf{x}^T \mathbf{H} \mathbf{x} + \mathbf{p}^T \mathbf{x} \\ & \text{subject to:} && \mathbf{x}^T \mathbf{Q}_i \mathbf{x} + \mathbf{q}_i^T \mathbf{x} + r_i \leq 0 \quad \text{for } 1 \leq i \leq q \end{aligned} \quad (1.77\text{a-b})$$

where \mathbf{H} and \mathbf{Q}_i are symmetric and positive semidefinite.

- **Semidefinite programming (SDP)**

- ◇ **Primal SDP**

$$\begin{aligned} & \text{minimize} && \mathbf{C} \cdot \mathbf{X} \\ & \text{subject to:} && \mathbf{A}_i \cdot \mathbf{X} = b_i \quad \text{for } 1 \leq i \leq p \\ & && \mathbf{X} \succeq \mathbf{0} \end{aligned} \tag{1.78a-c}$$

where variable \mathbf{X} and data \mathbf{A}_i and \mathbf{C} are symmetric matrices and $\mathbf{C} \cdot \mathbf{X}$ denotes the inner product in the space of real symmetric matrices defined by

$$\mathbf{C} \cdot \mathbf{X} = \sum_{i=1}^n \sum_{j=1}^n c_{i,j} x_{i,j}$$

- ◇ **Dual SDP**

$$\begin{aligned} & \text{minimize} && \mathbf{c}^T \mathbf{x} \\ & \text{subject to:} && \mathbf{F}(\mathbf{x}) \succeq \mathbf{0} \end{aligned} \tag{1.79a-b}$$

where $\mathbf{F}(\mathbf{x}) = \mathbf{F}_0 + \sum_{i=1}^p x_i \mathbf{F}_i$.

- **Second-order cone programming (SOCP)**

$$\begin{aligned} & \text{minimize} && \mathbf{b}^T \mathbf{x} \\ & \text{subject to:} && \|\mathbf{A}_j^T \mathbf{x} + \mathbf{c}_j\| \leq \mathbf{b}_j^T \mathbf{x} + d_j \quad \text{for } 1 \leq j \leq q \end{aligned} \tag{1.80a-b}$$

D. Lagrange Dual [6]

Duality is an important concept associated with constrained convex problems. To study the Lagrange dual, we need a concept called *Lagrange dual function*. Consider the general CP problem

$$\begin{aligned} & \text{minimize} && f(\mathbf{x}) \\ & \text{subject to:} && \mathbf{a}_i^T \mathbf{x} = b_i \quad \text{for } 1 \leq i \leq p \\ & && c_j(\mathbf{x}) \leq 0 \quad \text{for } 1 \leq j \leq q \end{aligned} \tag{1.81a-c}$$

where $f(\mathbf{x})$ and $c_j(\mathbf{x})$ are convex and recall its Lagrangian

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{i=1}^p \lambda_i (\mathbf{a}_i^T \mathbf{x} - b_i) + \sum_{j=1}^q \mu_j c_j(\mathbf{x})$$

Definition 1.18 *Lagrange dual function*

The *Lagrange dual function* of problem (1.81) is defined as

$$q(\boldsymbol{\lambda}, \boldsymbol{\mu}) = \inf_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \tag{1.82}$$

for $\boldsymbol{\lambda} \in R^p$ and $\boldsymbol{\mu} \in R^q$ with $\boldsymbol{\mu} \geq \mathbf{0}$. Note that the Lagrangian is *convex* with respect to \mathbf{x} . If

$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ is unbounded from below for some $\{\boldsymbol{\lambda}, \boldsymbol{\mu}\}$, then the value of $q(\boldsymbol{\lambda}, \boldsymbol{\mu})$ is assigned to $-\infty$.

Property 1 $q(\boldsymbol{\lambda}, \boldsymbol{\mu})$ is a concave function with respect to $\{\boldsymbol{\lambda}, \boldsymbol{\mu}\}$. The property follows from

that fact that for $\boldsymbol{\lambda}_1, \boldsymbol{\lambda}_2 \in R^p$ and $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2 \in R^q$ with $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2 \geq \mathbf{0}$ and for $t \in (0, 1)$, we have

$$\begin{aligned} q(t\boldsymbol{\lambda}_1 + (1-t)\boldsymbol{\lambda}_2, t\boldsymbol{\mu}_1 + (1-t)\boldsymbol{\mu}_2) &= \inf_{\mathbf{x}} L(\mathbf{x}, t\boldsymbol{\lambda}_1 + (1-t)\boldsymbol{\lambda}_2, t\boldsymbol{\mu}_1 + (1-t)\boldsymbol{\mu}_2) \\ &= \inf_{\mathbf{x}} \left[(t+1-t)f(\mathbf{x}) + \sum_{i=1}^p (t\lambda_{1,i} + (1-t)\lambda_{2,i})(a_i^T \mathbf{x} - b_i) + \sum_{j=1}^q (t\mu_{1,j} + (1-t)\mu_{2,j})c_j(\mathbf{x}) \right] \\ &\geq t \cdot \inf_{\mathbf{x}} \left[f(\mathbf{x}) + \sum_{i=1}^p \lambda_{1,i}(a_i^T \mathbf{x} - b_i) + \sum_{j=1}^q \mu_{1,j}c_j(\mathbf{x}) \right] + (1-t) \cdot \inf_{\mathbf{x}} \left[f(\mathbf{x}) + \sum_{i=1}^p \lambda_{2,i}(a_i^T \mathbf{x} - b_i) + \sum_{j=1}^q \mu_{2,j}c_j(\mathbf{x}) \right] \\ &= t \cdot q(\boldsymbol{\lambda}_1, \boldsymbol{\mu}_1) + (1-t) \cdot q(\boldsymbol{\lambda}_2, \boldsymbol{\mu}_2) \end{aligned}$$

Definition 1.19 *Lagrange dual problem*

The *Lagrange dual problem* with respect to problem (1.81) is defined as

$$\begin{aligned} &\underset{\boldsymbol{\lambda}, \boldsymbol{\mu}}{\text{maximize}} && q(\boldsymbol{\lambda}, \boldsymbol{\mu}) \\ &\text{subject to:} && \boldsymbol{\mu} \geq \mathbf{0} \end{aligned} \tag{1.83a-b}$$

Property 2 For any \mathbf{x} feasible for problem (1.81) and $\{\boldsymbol{\lambda}, \boldsymbol{\mu}\}$ feasible for problem (1.83), we have

$$f(\mathbf{x}) \geq q(\boldsymbol{\lambda}, \boldsymbol{\mu}) \tag{1.84}$$

This is because

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{i=1}^p \lambda_i (a_i^T \mathbf{x} - b_i) + \sum_{j=1}^q \mu_j c_j(\mathbf{x}) = f(\mathbf{x}) + \sum_{j=1}^q \mu_j c_j(\mathbf{x}) \leq f(\mathbf{x})$$

thus

$$q(\boldsymbol{\lambda}, \boldsymbol{\mu}) = \inf_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \leq L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \leq f(\mathbf{x})$$

We call the convex minimization problem in (1.81) the *primal* problem and the concave maximization problem in (1.83) the *dual* problem. The property (1.84) leads naturally to the concept of *duality gap* between the primal and dual objectives, which is defined by

$$\delta(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) - q(\boldsymbol{\lambda}, \boldsymbol{\mu}) \tag{1.85}$$

It follows that for feasible $\{\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}\}$ the duality gap is always nonnegative.

Property 3 Let \mathbf{x}^* be a solution of the primal problem in (1.81). Then the dual function at any feasible $\{\boldsymbol{\lambda}, \boldsymbol{\mu}\}$ serves as a lower bound of the optimal value of the primal objective, $f(\mathbf{x}^*)$, namely,

$$f(\mathbf{x}^*) \geq q(\boldsymbol{\lambda}, \boldsymbol{\mu}) \quad (1.86)$$

This property follows immediately from (1.84) by taking the minimum of $f(\mathbf{x})$ on its left-hand side.

- A question that naturally arises is what is the tightest lower bound the dual function $q(\boldsymbol{\lambda}, \boldsymbol{\mu})$ can offer? Obviously the answer is found by maximizing the dual function $q(\boldsymbol{\lambda}, \boldsymbol{\mu})$ on the right-hand side of (1.86) subject to $\boldsymbol{\mu} \geq \mathbf{0}$, which is exactly the Lagrange dual problem as formulated in (1.83). Therefore, if we denote the solution of (1.83) by $(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$, then we have

$$f(\mathbf{x}^*) \geq q(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*) \quad (1.87)$$

Based on (1.79), we now introduce the concept of *strong* and *weak* duality as follows.

Definition 1.20 *Strong and weak duality*

Let \mathbf{x}^* and $(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ be solutions of primal problem (1.81) and dual problem (1.83), respectively.

We say strong duality holds if $f(\mathbf{x}^*) = q(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$, i.e., the optimal duality gap is zero; and a weak duality holds if $f(\mathbf{x}^*) > q(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$. ■

- It can be shown that if the primal problem is strictly feasible, i.e., there exists \mathbf{x} satisfying

$$\begin{aligned} \mathbf{a}_i^T \mathbf{x} &= b_i \quad \text{for } 1 \leq i \leq p \\ c_j(\mathbf{x}) &< 0 \quad \text{for } 1 \leq j \leq q \end{aligned}$$

(this is to say that the interior of the feasible region of problem (1.81) is nonempty), then strong duality holds, i.e., the optimal duality gap is zero.

Example 1.13 Find the Lagrange dual of the LP problem

$$\begin{aligned} &\text{minimize} \quad \mathbf{c}^T \mathbf{x} \\ &\text{subject to:} \quad \mathbf{A}\mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0} \end{aligned}$$

Solution We write $\mathbf{x} \geq \mathbf{0}$ as $-\mathbf{x} \leq \mathbf{0}$ hence the Lagrangian of the LP problem is given by

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \mathbf{c}^T \mathbf{x} + (\mathbf{A}\mathbf{x} - \mathbf{b})^T \boldsymbol{\lambda} - \mathbf{x}^T \boldsymbol{\mu}$$

thus

$$q(\boldsymbol{\lambda}, \boldsymbol{\mu}) = \inf_{\mathbf{x}} \{ \mathbf{c}^T \mathbf{x} + (\mathbf{A}\mathbf{x} - \mathbf{b})^T \boldsymbol{\lambda} - \mathbf{x}^T \boldsymbol{\mu} \} = \inf_{\mathbf{x}} \{ (\mathbf{c} + \mathbf{A}^T \boldsymbol{\lambda} - \boldsymbol{\mu})^T \mathbf{x} - \mathbf{b}^T \boldsymbol{\lambda} \} \quad (1.88)$$

For given $\{\boldsymbol{\lambda}, \boldsymbol{\mu}\}$ such that $\mathbf{c} + \mathbf{A}^T \boldsymbol{\lambda} - \boldsymbol{\mu} \neq \mathbf{0}$, by (1.88) we have $q(\boldsymbol{\lambda}, \boldsymbol{\mu}) = -\infty$. Therefore to

deal with a well-defined dual function $q(\boldsymbol{\lambda}, \boldsymbol{\mu})$ we assume $\mathbf{c} + \mathbf{A}^T \boldsymbol{\lambda} - \boldsymbol{\mu} = \mathbf{0}$ which leads to

$$q(\lambda, \mu) = \inf_x (-b^T \lambda) = -b^T \lambda$$

and the Lagrange dual of the LP problem in question is given by

$$\begin{aligned} & \underset{\lambda, \mu}{\text{maximize}} && -b^T \lambda \\ & \text{subject to:} && \mu \geq 0 \end{aligned}$$

Since $c + A^T \lambda - \mu = 0$, $\mu = c + A^T \lambda$ so the above problem becomes

$$\begin{aligned} & \underset{\lambda}{\text{maximize}} && -b^T \lambda \\ & \text{subject to:} && -c - A^T \lambda \leq 0 \end{aligned}$$

which is equivalent to

$$\begin{aligned} & \underset{\lambda}{\text{minimize}} && b^T \lambda \\ & \text{subject to:} && (-A^T) \lambda \leq c \end{aligned}$$

Example 1.14 Find the Lagrange dual of the QP problem

$$\begin{aligned} & \underset{x}{\text{minimize}} && \frac{1}{2} x^T H x + p^T x \\ & \text{subject to:} && A x \leq b \end{aligned} \tag{1.89a-b}$$

where H is positive definite.

Solution The Lagrangian of the QP problem is given by

$$L(x, \mu) = \frac{1}{2} x^T H x + p^T x + \mu^T (A x - b)$$

Hence

$$q(\mu) = \inf_x \left\{ \frac{1}{2} x^T H x + p^T x + \mu^T (A x - b) \right\} \tag{1.90}$$

where the infimum (defined as the largest lower bound of the objective function) is attained at $x = -H^{-1}(p + A^T \mu)$. By substituting this solution into (1.90), we obtain

$$q(\mu) = -\frac{1}{2} \mu^T A H^{-1} A^T \mu - \mu^T (A H^{-1} p + b) - \frac{1}{2} p^T H^{-1} p \tag{1.91}$$

If we let $P = A H^{-1} A^T$, $t = A H^{-1} p + b$ and neglect the constant term in (1.91), the dual function

becomes $q(\mu) = -\frac{1}{2} \mu^T P \mu - \mu^T t$, hence the Lagrange dual of (1.89) is given by

$$\begin{aligned} & \underset{\mu}{\text{minimize}} && \frac{1}{2} \mu^T P \mu + \mu^T t \\ & \text{subject to:} && \mu \geq 0 \end{aligned} \tag{1.92}$$

Note that by definition matrix P in (1.92) is positive definite, therefore the dual problem is also a convex QP problem, but with simpler constraints in comparison with the primal problem in (1.89). In addition, if the number of constraints involved in the primal problem is smaller than n , so is the size of the dual problem. ■

E. CVX

• *The material of this section is largely based on the User Guide version 1.22 of the software written by M. Grant and S. Boyd.*

- **What is CVX?**

Designed by Michael Grant and Stephen Boyd, with input from Yinyu Ye, `cvx` is a modeling system for disciplined convex programming that are convex optimization problems described by a limited set of construction rules. `cvx` solves standard problems such as LP, QP, SOCP, and SDP problems. Compared to directly using a solver such as SeDuMi for these problems, `cvx` greatly simplifies the task of specifying the problem. `cvx` also solves more complex convex optimization problems, including many involving nonsmooth functions, such as l_1 norm.

- **Examples**

Example 1.15 We first consider the most basic convex optimization problem, least-squares. In a least-squares problem, we seek $\mathbf{x} \in \mathbb{R}^n$ that minimizes $\|\mathbf{Ax} - \mathbf{b}\|_2$, where $\mathbf{A} \in \mathbb{R}^{m \times n}$ is skinny and full rank (i.e., $m \geq n$ and $\text{rank}(\mathbf{A}) = n$). Let us create some test problem data for m , n , \mathbf{A} , and \mathbf{b} in MATLAB:

```
m = 16; n = 8;
A = randn(m,n);
b = randn(m,1);
```

(We chose small values of m and n to keep the output readable.) Then the least squares solution $\mathbf{x} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}$ is easily computed using the backslash operator:

```
x_ls = A\b;
```

Using `cvx`, the same problem can be solved as follows:

```
cvx_begin
    variable x(n);
    minimize( norm(A*x-b) );
cvx_end
```

When MATLAB reaches command `cvx_end`, the least-squares problem is solved, and the Matlab variable \mathbf{x} is overwritten with the solution of the least-squares problem, i.e., $(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}$. Now \mathbf{x} is an ordinary length- n numerical vector, identical to what would be obtained in the traditional approach, at least to within the accuracy of the solver. In addition, two additional Matlab variables are created:

- `cvx_optval`, which contains the value of the objective function; i.e., $\|\mathbf{Ax} - \mathbf{b}\|_2$;
- `cvx_status`, which contains a string describing the status of the calculation. In this case, `cvx_status` would contain the string `Solved`. See Appendix C of the user guide for a list of the possible values of `cvx_status` and their meaning.
- `cvx_slvtol`: the tolerance level achieved by the solver.
- `cvx_slvitr`: the number of iterations taken by the solver.

Example 1.16 Suppose we wish to add some simple upper and lower bounds to the

least-squares problem above. That is, we wish to solve the CP problem

$$\begin{aligned} & \text{minimize} && \| \mathbf{Ax} - \mathbf{b} \|_2 \\ & \text{subject to:} && \mathbf{l} \leq \mathbf{x} \leq \mathbf{u} \end{aligned}$$

where \mathbf{l} and \mathbf{u} are given data, vectors with the same dimension as the variable \mathbf{x} . We can no longer use the simple backslash notation to solve this problem, but it can be transformed into a QP, which can be solved without difficulty if you have some form of QP software available.

Let us provide some numeric values for \mathbf{l} and \mathbf{u} :

```
bnds = randn(n,2);
l = min(bnds,[],2);
u = max(bnds,[],2);
```

Then if you have the MATLAB Optimization Toolbox, you can use the `quadprog` function to solve the problem as follows:

```
x_qp = quadprog(2*A'*A,-2*A'*b,[],[],[],[],l,u);
```

This actually minimizes the square of the norm, which is the same as minimizing the norm itself. In contrast, the `cvx` specification is given by

```
cvx_begin
    variable x(n);
    minimize(norm(A*x-b));
    subject to
        l <= x <= u;
cvx_end
```

Example 1.17 `cvx` supports constraints more complex than simple bounds as seen in Example 1.16. For example, let us define new matrices \mathbf{C} and \mathbf{d} in MATLAB as follows:

```
p = 4;
C = randn(p,n);
d = randn(p,1);
```

Now let us add an equality constraint and a nonlinear inequality constraint to the original least-squares problem:

```
cvx_begin
    variable x(n);
    minimize(norm(A*x-b));
    subject to
        C*x == d;
        norm(x,Inf) <= 1;
cvx_end
```

Example 1.18 For our final example in this section, let us show how traditional MATLAB code and `cvx` specifications can be mixed to form and solve multiple optimization problems. The following code solves the problem of minimizing $\| \mathbf{Ax} - \mathbf{b} \|_2 + \gamma \| \mathbf{x} \|_1$, for a logarithmically spaced vector of (positive) values of γ . This gives us points on the optimal trade-off curve

between $\|Ax - b\|_2$ and $\|x\|_1$. An example of this curve is given in Fig. 1.9.

```
gamma = logspace(-2,2,20);
l2norm = zeros(size(gamma));
l1norm = zeros(size(gamma));
fprintf(1, ' gamma norm(x,1) norm(A*x-b)\n');
fprintf(1, '-----\n');
for k = 1:length(gamma),
    fprintf(1,'%8.4e', gamma(k));
    cvx_begin quiet
        variable x(n);
        minimize(norm(A*x-b)+gamma(k)*norm(x,1));
    cvx_end
    l1norm(k) = norm(x,1);
    l2norm(k) = norm(A*x-b);
    fprintf(1,' %8.4e %8.4e\n',l1norm(k),l2norm(k));
end
plot(l1norm, l2norm);
xlabel('norm (x,1)');
ylabel('norm (A*x-b)');
grid
```

- **Basic rules**

- (1) All `cvx` models must be preceded by command `cvx_begin` and terminated with command `cvx_end`.

Command `cvx_begin` accepts several variants. For example, `cvx_begin quiet` prevents the model from producing screen output while it is been solved.

- (2) **Data Types for Variables**

All variables must be declared using command `variable` (or `variables`). Variables can be real or complex; and scalar, vector, matrix, or n-dimensional array. Variables can be matrices with structure,

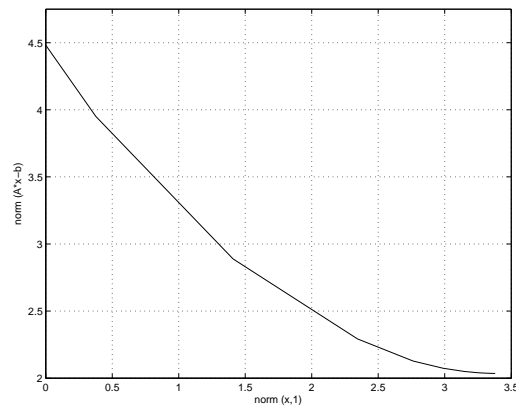


Fig. 1.9

like symmetry or bandedness. The variable structure is given by including descriptive keywords

following the name and size of the variable. Here are some examples of `cvx` variable:

```
variable w(50) complex;
variable X(20,10);
variable Y(50,50) symmetric;
variable Z(100,100) hermitian toeplitz;
```

- **Structure Keywords**

Version 1.22 of `cvx` supports the following structure keywords:

```
banded(lb,ub) complex diagonal hankel hermitian lower_bidiagonal
lower_hessenberg lower_triangular scaled_identity skew_symmetric
symmetric toeplitz tridiagonal upper_bidiagonal upper_hankel
upper_hessenberg upper_triangular
```

(3) Objective functions

Declaring an objective function requires the use of the `minimize` or `maximize` function, as appropriate. The objective function in a call to `minimize` must be convex; the objective function in a call to `maximize` must be concave. At most one objective function may be declared in a given `cvx` specification, and the objective function must have a scalar value.

If no objective function is specified, the problem is interpreted as a *feasibility* problem, which is the same as performing a minimization with the objective function set to zero. In this case, `cvx_optval` is either 0, if a feasible point is found, or `+Inf`, if the constraints are not feasible.

(4) Constraints

Version 1.22 of `cvx` supports the following constraint types:

- Equality `==` constraints, where both the left- and right-hand sides are affine functions of the optimization variables.
- Less-than `<=` inequality constraints, where the left-hand expression is convex, and the right-hand expression is concave.
- Greater-than `>=` constraints, where the left-hand expression is concave, and the right-hand expression is convex.

These equality and inequality operators work for arrays. When both sides of the constraint are arrays of the same size, the constraint is imposed elementwise. For example, if a and b are $m \times n$ matrices, then $a \leq b$ is interpreted by `cvx` as mn (scalar) inequalities, i.e., each entry of a must be less than or equal to the corresponding entry of b . `cvx` also handles cases where one side is a scalar and the other is an array. This is interpreted as a constraint for each element of the array, with the (same) scalar appearing on the other side. As an example, if a is an $m \times n$ matrix, then $a \geq 0$ is interpreted as mn inequalities: each element of the matrix must be nonnegative.

(5) Functions

The base `cvx` function library includes a variety of convex, concave, and affine functions which accept `cvx` variables or expressions as arguments. Many are common MATLAB functions such as `sum`, `trace`, `diag`, `sqrt`, `max`, and `min`, re-implemented as needed to support `cvx`; others are new functions not found in MATLAB. A complete list of the functions in the base library can be found in §B of the user guide. It's also possible to add your own new functions; see §5 of the user guide.

(6) Sets

◇ `cvx` supports the definition and use of convex sets. The base library includes the cone of positive semidefinite $n \times n$ matrices, the second-order or Lorentz cone, and various norm balls. A complete list of sets supplied in the base library is given in §B of the user guide.

◇ Since MATLAB does not have a set membership operator, `cvx` adopts a slightly different syntax to require that an expression is in a set. To represent a set, a function is used to return an unnamed variable that is required to be in the set. Consider, for example, S_+^n , the cone of symmetric positive semidefinite $n \times n$ matrices. In `cvx`, this is represented by function `semidefinite(n)`, which returns an unnamed new variable, that is constrained to be positive semidefinite. To require that the matrix expression `X` be symmetric positive semidefinite, the syntax `X == semidefinite(n)` is used. The literal meaning of this is that `x` is constrained to be equal to some unnamed variable, which is required to be an $n \times n$ symmetric positive semidefinite matrix. This is equivalent to saying that `X` must be symmetric positive semidefinite.

As an example, consider the constraint that a (matrix) variable `x` is a correlation matrix, i.e., it is symmetric, has unit diagonal elements, and is positive semidefinite.

◇ In `cvx` we can declare such a variable and impose such constraints using

```
variable X(n,n) symmetric;
X == semidefinite(n);
diag(X) == ones(n,1);
```

The second line here imposes the constraint that `X` be positive semidefinite (you can read `'=='` here as `'is'`, so the second line can be read as `'X is positive semidefinite'`). The left-hand side of the third line is a vector containing the diagonal elements of `x`, whose elements we require to be equal to one. Incidentally, `cvx` allows us to simplify the third line to

```
diag(X) == 1;
```

because `cvx` follows the MATLAB convention of handling array/scalar comparisons by comparing each element of the array independently with the scalar.

◇ Sets can be combined in affine expressions, and we can constrain an affine expression to be in a convex set. For example, we can impose constraints of the form

$$A^*X^*A' - X == B^*\text{semidefinite}(n)^*B';$$

where `x` is an $n \times n$ symmetric variable matrix, and **A** and **B** are $n \times n$ constant matrices. This constraint requires that $AXA^T - X = BYB^T$, for some $Y \in S_+^n$.

Example 1.19 Standard-form LP

```
function x = lps_cvx(c,A,b)
n = length(c);
cvx_begin quiet
    variable x(n);
    minimize(c'*x);
    subject to
        A*x == b;
```

```

    x >= 0;
cvx_end

```

Example 1.20 Alternative-form LP

```

function x = lpa_cvx(c,A,b)
n = length(c);
cvx_begin quiet
    variable x(n);
    minimize(c'*x);
    subject to
        A*x >= b;
cvx_end

```

Example 1.21 Convex QP

```

function x = qp_cvx(H,p,A,b,E,d)
n = length(p);
cvx_begin quiet
    variable x(n);
    minimize(0.5*x'*H*x+x'*p);
    subject to
        A*x == b;
        E*x >= d;
cvx_end

```

Example 1.22 SDP

```

% Program: sdp_cvx.m
% To solve semidefinite programming (SDP) problem
%           minimize    c'*x
%           subject to: F0 + x(1)*F1 + ... + x(p)*Fp >= 0
% using cvx, where ">= 0" means being positive semidefinite,
% and x = [x1 x2 ... xp]'.
% Input:
% c: constant vector of length p defining objective function.
% F = [F0 F1 ... Fp] with each Fi a symmetric square matrix.
% Output:
% x: solution of the SDP problem.
% Written by W.-S. Lu, University of Victoria.
% Last modified: Oct. 29, 2012.
% Example: Solve the SDP problem involved in Example 14.1 from PO, where
% c = [0 0 0 -1]';
% F = [F0 F1 F2 F3 F4]
% = [-2.0   0.5   0.6   0  -1   0   0   0  -1   0   0   0  -1   0   0;
%     0.5  -2.0  -0.4  -1   0   0   0   0   0   0   0  -1   0  -1   0;
%     0.6  -0.4  -3.0   0   0   0  -1   0   0   0  -1   0   0   0  -1];
% x = sdp_cvx(c,F);
function x = sdp_cvx(c,F)

```



```

p = length(c);
[n,m] = size(F);
Fx = F(:,1:n);
cvx_begin quiet
    variable x(p);
    minimize(c'*x);
    subject to
    for i = 1:p,
        Fx = Fx + x(i)*F(:,(i*n+1):((i+1)*n));
    end
    Fx == semidefinite(n);
cvx_end

```

Example 1.23 SOCP

```

% Program: socp_cvx.m
% To solve second-order cone programming (SOCP) problem
%     minimize b'*x
%     subject to: ||Ai'*x + ci|| <= bi'*x + di
%               for i = 1, ..., q
% using cvx, where x = [x1 x2 ... xm] and Ai is a matrix of size m by ni.
Thus there are a total of q
% 2nd-order (Lorentz) cones, each with size ni.
% Input:
% A = [A1 A2 ... Aq]
% B = [b1 b2 ... bq]
% b: constant vector of length m defining objective function
% c = [c1; c2; ...; cq]
% d = [d1 d2 ... dq]
% k = [n1 n2 ... nq]
% Output:
% x: solution of the SOCP problem.
% Written by W.-S. Lu, University of Victoria.
% Last modified: Oct. 30, 2012.
% Example: Solve the SOCP problem in Example 14.5 from PO, where
% b = [1 0 0 0 0]';
% A = [0 0 0 0 0 0;
%      -1 0 0.5 0 0 0;
%      0 1 0 1 0 0;
%      1 0 0 0 -0.7071 -0.3536;
%      0 -1 0 0 -0.7071 0.3536];
% z = zeros(5,1); B = [b z z];
% c = [0 0 -0.5 0 4.2426 -0.7071]';
% d = [0 1 1];
% k = [2 2 2];
% x = socp_cvx(A,B,b,c,d,k);
function x = socp_cvx(A,B,b,c,d,k)

```

```

m = length(b);
q = length(d);
t = 0;
cvx_begin quiet
    variable x(m);
    minimize(b'*x);
    subject to
        for i = 1:q,
            Ai = A(:,(t+1):(t+k(i)));
            bi = B(:,i);
            ci = c((t+1):(t+k(i)));
            di = d(i);
            norm(Ai'*x+ci) <= bi'*x + di;
            t = t + k(i);
        end
cvx_end

```

Problems

1.1 [7] Suppose we have three colored boxes r (red), b (blue), and g (green). Box r contains 3 apples, 4 oranges, and 3 limes; box b contains 1 apple, 1 orange, and 0 limes; and box g contains 3 apples, 3 oranges, and 4 limes. If a box is chosen at random with probabilities $p(r) = 0.2$, $p(b) = 0.2$, and $p(g) = 0.6$, and a piece of fruit is removed from the box (with equal probability of selecting any of the items in the box), then

- (i) what is the probability of selecting an apple?
- (ii) If we observe that the selected fruit is an orange, what is the probability that it came from the green box?

1.2 [7] There are 2 opaque bags. One bag has two blue balls and the other has one blue ball and one red ball. Suppose you pick a bag at random, and then pick a ball from that bag at random. You found the ball you picked was a blue ball. Now you pick the second ball from that same bag. Compute the probability that this ball is also blue?

1.3 Suppose we have a box that contains A red balls and B blue balls, one takes balls out of the box one by one, what is the probability that the balls that remain in the box are all red balls?

1.4 Prove the equality in Eq. (1.27).

1.5 Let x be a random variable with mean (expectation) $E[x] = \mu$ and variance $\text{Var}[x] = \sigma^2$.

Compute and express the mean and variance of random variable $ax + b$ in terms of μ and σ^2 , where a and b are constants.

1.6 Let x denote a set of temperature readings in Celsius (C). Treating x as a random variable, its average (expectation) and variance are known to be 22.5°C and $\sigma^2 = 0.4$. If we convert x to a set of temperature readings y in Fahrenheit (F) using $F = (9/5)C + 32$, what are the average and variance of the temperature data in Fahrenheit?

1.7 Let ξ_1 and ξ_2 be two independent random variables, find the distribution function of random variable $\xi = \xi_1 + \xi_2$ where

- (i) ξ_1 and ξ_2 are uniform distribution over intervals $(-5, 1)$ and $(1, 5)$, respectively.
- (ii) ξ_1 and ξ_2 are random variables with density functions $\varphi_1(x) = \varphi_2(x) = \frac{1}{2\alpha} e^{-|x|/\alpha}$ with some $\alpha > 0$.

1.8 Consider the convex QP problem

$$\underset{x}{\text{minimize}} \quad \frac{1}{2} \mathbf{x}^T \begin{bmatrix} 10 & -9 \\ -9 & 8.15 \end{bmatrix} \mathbf{x} + \mathbf{x}^T \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad (\text{P1.1})$$

- (i) Apply SDM to solve the problem with initial point $\mathbf{x}_0 = [1 \ 1]^T$. Exact line search is carried out using the closed-form formula

$$\alpha_k = \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_k^T \mathbf{H} \mathbf{g}_k}$$

where \mathbf{H} is the Hessian of the objective function and \mathbf{g}_k is the gradient of the objective function at \mathbf{x}_k . How many iterations the SDM has to take to achieve a solution \mathbf{x}_s with the accuracy of $\|\mathbf{x}_s - \mathbf{x}^*\|_2 \leq 10^{-5}$ where \mathbf{x}^* is the true minimizer of the problem which is known to be $\mathbf{x}^* = [-16.3 \quad -18]^T$?

(ii) Now solve the problem in part (i) by scaling the variable \mathbf{x} with $\mathbf{x} = \mathbf{T} \mathbf{y}$ first, where matrix \mathbf{T} is given by

$$\mathbf{T} = \begin{bmatrix} 0.25 & 4.5 \\ 0 & 5 \end{bmatrix}$$

Working with the scaled QP problem, how many SDM iterations are needed to achieve a solution with the same accuracy as in part (i)?

Explain why such a scaling matrix can improve SDM's performance? If this \mathbf{T} was not given, how do you find a similar \mathbf{T} to properly scale an ill-conditioned QP problem?

To earn credit, include the solution details, your MATLAB code, and all numerical results.

1.9 Consider the one-variable *polynomial curve fitting* problem where a polynomial model of the form

$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \cdots + w_m x^m = \sum_{i=0}^m w_i x^i \quad (\text{P1.1})$$

is used to fit a data set $\{(x_k, y_k), \text{ for } k = 1, \dots, n\}$, where $\mathbf{w} = [w_0 \quad w_1 \quad \cdots \quad w_m]^T$. The function $y(x, \mathbf{w})$ is a polynomial of order m in x , whose behavior for a fixed order m is determined by its coefficients $\{w_i, i = 0, 1, \dots, m\}$. Therefore, the *polynomial curve fitting* problem is about to find "optimal" $\{w_i, i = 0, 1, \dots, m\}$ such that the polynomial $y(x, \mathbf{w})$, when evaluated at $\{x_k, k = 1, 2, \dots, n\}$, produces $\{y(x_k, \mathbf{w}), k = 1, \dots, n\}$ that is the *closest* to the given data $\{y_k, k = 1, 2, \dots, n\}$ in a certain sense. The *least-squares* (LS) polynomial curve fitting is one of the most popular approach to deal with polynomial modeling problems, which in the present case amounts to finding an $(m + 1)$ -dimensional \mathbf{w} that solves the unconstrained problem

$$\underset{\mathbf{w}}{\text{minimize}} \quad e(\mathbf{w}) = \frac{1}{2} \sum_{k=1}^n |y(x_k, \mathbf{w}) - y_k|^2 \quad (\text{P1.2})$$

where $y(x, \mathbf{w})$ is given by (P1.1). Show that the optimal \mathbf{w} satisfies the linear system of equations

$A\mathbf{w} = \mathbf{t}$ where, for real-valued data set $\{(x_k, y_k), \text{ for } k = 1, \dots, n\}$, $A = (a_{i,j}) \in R^{(m+1) \times (m+1)}$ with $i, j = 0, 1, \dots, m$; and $\mathbf{t} = (t_i) \in R^{(m+1) \times 1}$ with $i = 0, 1, \dots, m$ are determined by the data set as

$$a_{i,j} = \sum_{k=1}^n (x_k)^{i+j} \quad \text{and} \quad t_i = \sum_{k=1}^n (x_k)^i y_k \quad (\text{P1.3})$$

To get credit, include all derivation details.

1.10 The purpose of this problem is to perform numerical evaluation of the LS data fitting method.

(i) Data preparation: the data to be used in what follows was generated from a 2nd order polynomial model that were contaminated by a small amount of Gaussian noise. Use the MATLAB code below to produce the data.

```
w_target = [9.5 -5.6 1]; % here the vector is associated with a poly with increasing order.
```

```
p_target = fliplr(w_target); % function polyval requires poly coefficients with decreasing order.
```

```
x = [0.5 2.5 3 3.95 4.05];
```

```
randn('state',16)
```

```
y = polyval(p_target,x) + 0.4*randn(1,5);
```

To see how the noisy data are related to the target polynomial, use MATLAB to plot the following:

```
t = 0:4.5/999:4.5;
```

```
z = polyval(p_target,t);
```

```
figure(1)
```

```
plot(t,z,'k')
```

```
hold on
```

```
plot(x,y,'ro','linewidth',1.5)
```

(ii) Apply the polynomial data fitting method from Prob. 1.7 to the data set $\{x, y\}$ generated in part (i) with a 2nd – order polynomial $y(x, \mathbf{w})$ (i.e. $m = 2$, see Eq. (P1.1)). To earn credit:

- Compute and report the numerical values of \mathbf{w} for the $y(x, \mathbf{w})$ that solves the minimization problem in (P1.2);
- Plot $y(t, \mathbf{w})$ versus $t = 0:4.5/999:4.5$ in Figure 1 generated in part (i) with the data set also displayed in the figure;
- Evaluate the approximation error, namely the value of the minimized objective function $e(\mathbf{w})$ (see (P1.2)) and comment the performance of the 2nd-order model obtained.

1.11 Repeat Problem 1.8 with a 4th-order polynomial $y(x, \mathbf{w})$.

1.12 Use CVX to solve the following LP and convex QP problems

(i)

$$\begin{aligned} &\text{minimize } f(\mathbf{x}) = 2x_1 + 9x_2 + 3x_3 \\ &\text{subject to: } -2x_1 + 2x_2 + x_3 - x_4 = 1 \\ &\quad \quad \quad x_1 + 4x_2 - x_3 - x_5 = 1 \\ &\quad \quad \quad x_i \geq 0 \quad \text{for } i = 1, 2, \dots, 5 \end{aligned}$$

(ii)

$$\begin{aligned} &\text{minimize } f(\mathbf{x}) = x_1 + x_2 \\ &\text{subject to: } \quad \quad \quad x_1 \leq 2 \\ &\quad \quad \quad \quad \quad \quad x_2 \leq 2 \\ &\quad \quad \quad \quad \quad \quad -x_2 \leq 5 \\ &\quad \quad \quad -2x_1 + x_2 \leq 2 \\ &\quad \quad \quad 2x_1 + x_2 \leq 4 \end{aligned}$$

(iii)

$$\begin{aligned} &\text{minimize } f(\mathbf{x}) = (x_1 - x_3)^2 + (x_2 - x_4)^2 \\ &\text{subject to: } \quad -x_1 \leq 0 \\ &\quad \quad \quad -x_2 \leq 0 \\ &\quad \quad \quad x_1 + 2x_2 \leq 2 \\ &\quad \quad \quad -x_4 \leq -2 \\ &\quad \quad \quad -x_3 - x_4 \leq -3 \\ &\quad \quad \quad x_3 + 2x_4 \leq 6 \end{aligned}$$

To earn credit, include CVX code and report numerical values of the solutions with an accuracy of no-less-than-eight decimal places.