

# Part IV

## The Appendices



# Appendix A

## Mathematical and Statistical Topics

### A.1 Trigonometry

Trigonometry is a branch of mathematics that studies triangles. The study of trigonometry deals with the relationships between the lengths of the sides and measures of the angles of a right triangle. As we will see, in two-dimensions, trigonometry can also be used to describe the relationships dealing with angles and the motion and shape of waves. This tutorial begins with the most basic of definitions. Although basic, the quantities defined below are fundamental to trigonometry.

**Definition A.1.1** RIGHT TRIANGLE.

*If the measures of one of the angles of a triangle is  $\pi/2$  ( $90^\circ$ ), then the triangle is a **right triangle**. The side of the triangle that is opposite the right angle is the **hypotenuse** of the right triangle. The length of the hypotenuse will be denoted by  $r$ .*

The *triangle postulate* is the statement that the sum of the measures of all the angles of a triangle is equal to  $\pi$  ( $180^\circ$ ). So for a right triangle, the two non-right angles must have measures that are less than  $\pi/2$ .

For any right triangle, the relationships between the lengths of the sides and the measures of the angles have specific definitions. These functions are the basic building blocks of trigonometry.

**Definition A.1.2** BASIC TRIGONOMETRIC FUNCTIONS.

*Let  $\theta$  denote the measure of an angle of a right triangle. Let  $\theta$  denote an angle of a right triangle,  $a$  denote the length of the side opposite  $\theta$ , and  $b$  denote the length of the side adjacent to  $\theta$ . Then for a right triangle, the trigonometric functions **sine**, **cosine**, and **tangent** are defined in terms of the ratios of the lengths of sides adjacent and opposite  $\theta$  and the length of the hypotenuse.*

$$\sin \theta = \frac{\text{length of opposite side}}{\text{length of hypotenuse}} = \frac{a}{r} \quad (\text{A.1})$$

$$\cos \theta = \frac{\text{length of adjacent side}}{\text{length of hypotenuse}} = \frac{b}{r} \quad (\text{A.2})$$

$$\tan \theta = \frac{\text{length of opposite side}}{\text{length of adjacent side}} = \frac{a}{b} = \frac{\sin \theta}{\cos \theta}. \quad (\text{A.3})$$

For an illustration of these definitions, see Figure A.1.

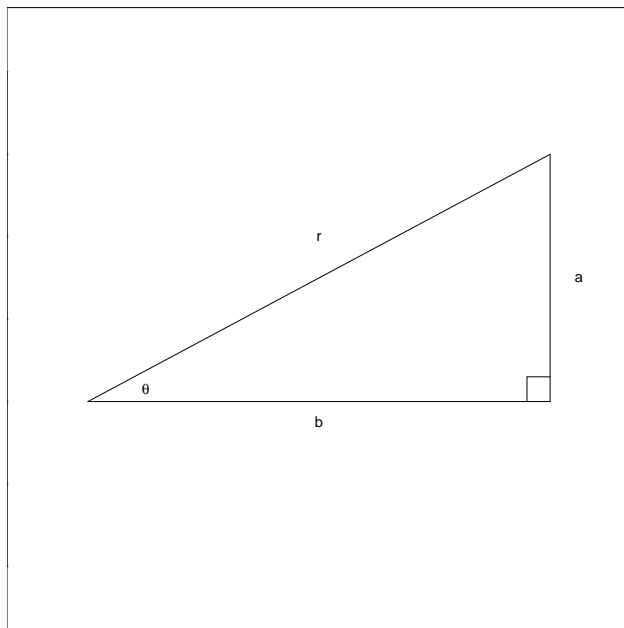


Figure A.1: Illustration of relationships of sides of triangle to angle  $A$  as described in Definition A.1.1.

The famous Pythagorean Theorem says that for a right triangle,

$$a^2 + b^2 = r^2.$$

**Definition A.1.3** RECIPROCAL OF BASIC TRIGONOMETRIC FUNCTIONS.

The reciprocals of the sine, cosine, and tangent functions are called **secant**, **cosecant**, and **cotangent**, respectively; that is,

$$\sec \theta = \frac{1}{\cos \theta}, \quad \csc \theta = \frac{1}{\sin \theta}, \quad \text{and} \quad \cot \theta = \frac{1}{\tan \theta}. \quad (\text{A.4})$$

Figure A.2 displays sine and cosine as a function of  $\theta$  for  $0 \leq \theta \leq 6\pi$ . In this figure, the sine and cosine function go through three cycles. Examination of this graph illustrate that the sine and cosine function are periodic, having period  $2\pi$ , as is the tangent function, which is not included in the graph. Mathematically, this is written as

$$\sin \theta = \sin(\theta + 2\pi), \quad \cos \theta = \cos(\theta + 2\pi), \quad \text{and} \quad \tan \theta = \tan(\theta + 2\pi).$$

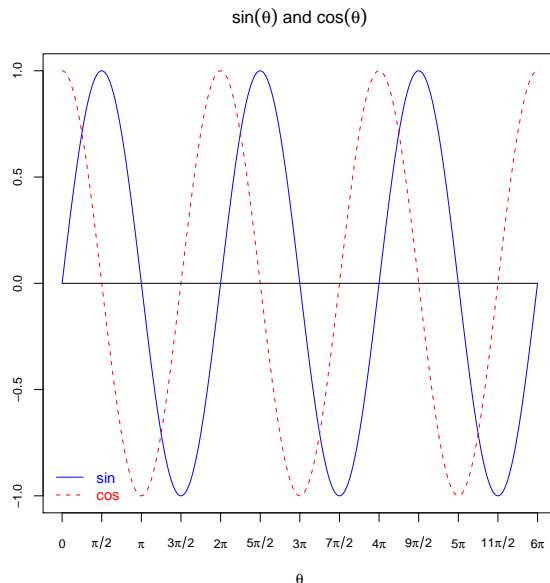


Figure A.2: Graphs of cosine and sine as functions of  $\theta$ .

Figure A.2 also readily illustrate how, in two-dimensions, trigonometry can also be used to describe the relationships dealing with angles and the motion and shape of waves.

For  $0 \leq \theta \leq \pi/2$ , it is also the case that

$$\begin{aligned}\sin(\theta) &= -\sin(\theta + \pi) = -\sin(\theta - \pi) = \sin(2\pi - \theta), \\ \cos(\theta) &= -\cos(\theta + \pi) = -\cos(\theta - \pi) = \cos(2\pi - \theta)\end{aligned}\tag{A.5}$$

but

$$\tan(\theta) = \tan(\theta + \pi) = \tan(\theta - \pi).$$

Figure A.3 illustrates this relationship for the sine function. Because the sine, cosine, and tangent functions are periodic, then, strictly speaking, they do not have a mathematical inverse. To define their inverse, it is necessary to restrict the values of  $\theta$  so that the functions are one-to-one and onto; that is, are bijective. Additionally, for the secant, cosecant, and cotangent, the value of  $\theta$  must not be so that the denominator is zero.

**Definition A.1.4** INVERSE OF BASIC TRIGONOMETRIC FUNCTIONS.

*For restricted domain of the trigonometric functions, we may define their inverses in the following manner.*

- (a) For  $-\pi/2 \leq \theta \leq \pi/2$ , the inverse of the sine function is called the **arcsine** and is defined by

$$\theta = \arcsin \phi \quad \text{if} \quad \phi = \sin \theta.$$

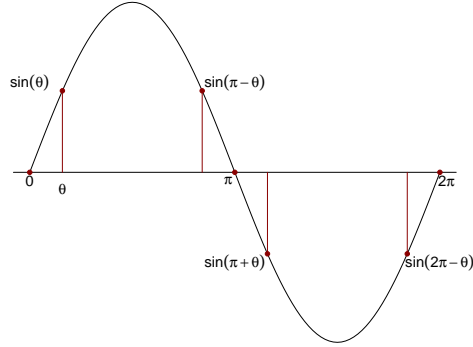


Figure A.3: Sine function for  $0 \leq \theta \leq 2\pi$  illustrating equation (A.5).

- (b) For  $0 \leq \theta \leq \pi$ , the inverse of the cosine function is called the **arccosine** and is defined by

$$\theta = \arccos \phi \quad \text{if} \quad \phi = \cos \theta.$$

- (c) For  $-\pi/2 \leq \theta \leq \pi/2$ , the inverse of the tangent function is called the **arctangent** and is defined by

$$\theta = \arctan \phi \quad \text{if} \quad \phi = \tan \theta. \quad (\text{A.6})$$

- (d) For  $-\pi/2 \leq \theta \leq \pi/2$  and  $\theta \neq 0$ , the inverse of the cosecant is called the **arccosecant** and is defined by

$$\theta = \operatorname{arccsc} \phi \quad \text{if} \quad \phi = \csc \theta.$$

- (e) For  $0 \leq \theta \leq \pi$  and  $\theta \neq \pi/2$ , the inverse of the secant is called the **arcsecant** and is defined by

$$\theta = \operatorname{arcsec} \phi \quad \text{if} \quad \phi = \sec \theta.$$

- (f) For  $0 \leq \theta \leq \pi$ , the inverse of the cotangent is called the **arccotangent** and is defined by

$$\theta = \operatorname{arccot} \phi \quad \text{if} \quad \phi = \cot \theta.$$

Equation (1.1) in Theoretical Exercise T1.2 in Chapter 1, provides an alternative definition that is used throughout the text.

Table A.1: Some useful trigonometric identities.

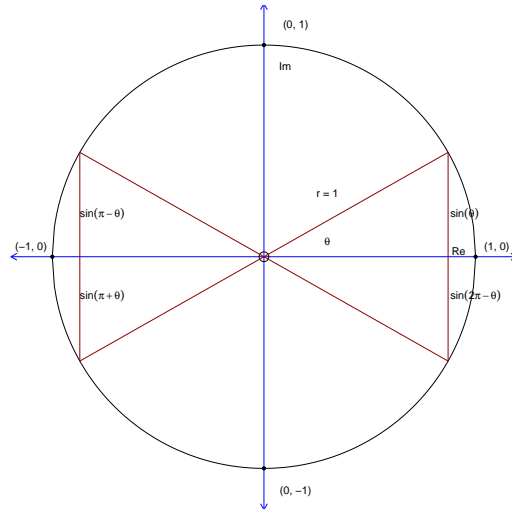
Double-angle Identities		
$\sin(2\theta) = 2 \sin \theta \cos \theta = 2 \tan \theta / (1 + \tan^2 \theta)$	$\cos(2\theta) = \cos^2 \theta - \sin^2 \theta = 1 - 2 \sin^2 \theta$	
$\tan(2\theta) = 2 \tan \theta / (1 - \tan^2 \theta)$	$\cot(2\theta) = (\cot^2 \theta - 1) / (2 \cot \theta)$	
Half-angle Identities		
$\sin(\theta/2) = \pm \sqrt{(1 - \cos \theta)/2}$	$\cos(\theta/2) = \pm \sqrt{(1 + \cos \theta)/2}$	
$\tan(\theta/2) = \pm \sqrt{(1 - \cos \theta)/(1 + \cos \theta)} = \sin \theta / (1 + \cos \theta) = (1 - \cos \theta) / \sin \theta$		
Product-to-Sum Identities		
$2 \cos \alpha \cos \beta = \cos(\alpha + \beta) + \cos(\alpha - \beta)$	$2 \sin \alpha \sin \beta = \cos(\alpha - \beta) - \cos(\alpha + \beta)$	
$2 \sin \alpha \cos \beta = \sin(\alpha + \beta) + \sin(\alpha - \beta)$	$2 \cos \alpha \sin \beta = \sin(\alpha + \beta) - \sin(\alpha - \beta)$	
Sum-to-Product Identities		
$\sin \alpha + \sin \beta = 2 \sin[(\alpha + \beta)/2] \cos[(\alpha - \beta)/2]$	$\cos \alpha + \cos \beta = 2 \cos[(\alpha + \beta)/2] \cos[(\alpha - \beta)/2]$	
$\sin \alpha - \sin \beta = 2 \sin[(\alpha - \beta)/2] \cos[(\alpha + \beta)/2]$	$\cos \alpha - \cos \beta = -2 \sin[(\alpha + \beta)/2] \sin[(\alpha - \beta)/2]$	
$\tan(\alpha \pm \beta) = (\tan \alpha \pm \tan \beta) / (1 \mp \tan \alpha \tan \beta)$		
Other Useful Identities		
$2 \sin^2 \theta = 1 - \cos(2\theta)$	$2 \cos^2 \theta = 1 + \cos(2\theta)$	$\sin^2 \theta + \cos^2 \theta = 1$ (Pythagorean Identity)

When the measure of  $\theta$  is restricted to be between 0 and  $\pi/2$ , some interesting relationships between the previously defined trigonometric functions are the following.

$$\begin{aligned}\sin \theta &= \cos\left(\frac{\pi}{2} - \theta\right) \\ \cos \theta &= \sin\left(\frac{\pi}{2} - \theta\right) \\ \tan \theta &= \cot\left(\frac{\pi}{2} - \theta\right)\end{aligned}$$

Identities such as these, that interrelate trigonometric functions are referred to as **trigonometric identities**. Table A.1 contains a set of trigonometric identities that are useful in the some of the exercises and proofs in this text, especially those that relate to the spectral analysis of time series. There are many other identities. Common identities that are omitted from this table can, for the most part, be easily derived using those identities that do appear in the table. In the last row of the table is the “Pythagorean Identity,” which is easily proven using the Pythagorean Theorem.

Figure A.1 illustrate that the trigonometric functions can be translated onto the complex mathematical plane, where the horizontal axis represents the space of real numbers and the vertical axis the space of imaginary numbers. A famous result, known as Euler's formula,



state that for any real number  $\theta$ ,

$$e^{i\theta} = \cos \theta + i \sin \theta, \quad (\text{A.7})$$

where  $i = \sqrt{-1}$  is the **imaginary number**. For an imaginary number  $z = a + bi$ , let  $\Re(z) = a$  be the real part of  $z$  and  $\Im(z) = b$  be the imaginary part of  $z$ . Then, for any real number  $\theta$ , a little algebra will show that the sine and cosine functions may be written

$$\cos \theta = \Re(e^{i\theta}) = \frac{e^{i\theta} + e^{-i\theta}}{2} \quad (\text{A.8})$$

$$\sin \theta = \Im(e^{i\theta}) = \frac{e^{i\theta} - e^{-i\theta}}{2i}. \quad (\text{A.9})$$

## A.2 Fourier Series and Spectral Representations

The ideas of frequency domain analysis of data are an important part of time series analysis. We have attempted to motivate these ideas intuitively in the main part of this book. In this appendix, we describe in more detail the mathematics behind the ideas. For more details, the reader is encouraged to consult Chapter 4 of Priestley (1981) or Chapter 7 of Anderson (1994). We recall that the basic idea of the frequency domain analysis of time series is to express various quantities as the “sum” of orthogonal sinusoids.



### A.2.1 Fourier Series for Periodic Functions

Suppose that  $f$  is a function defined on a finite interval  $[a, b]$ . By this we also include functions that are periodic or period  $b - a$ . Throughout this book we have taken our interval to be  $[0, 1]$ . If it is not, we can construct the function  $g(y) = f[(x - a)/(b - a)]$  and apply the results to the function  $g$ . The interval  $[0, 1]$  is natural for spectral densities as then frequency can be thought of as cycles per unit of time. Thus for monthly data, frequency  $1/12$  refers to possible cycles of length 12 months. The interval traditionally considered in mathematics is  $[-\pi, \pi]$ , which means for time series that one has to perform mental arithmetic involving  $\pi$  in order to convert frequency to a physically meaningful quantity. Using the interval  $[0, 1]$  also simplifies many of the formulas in the subject.

Our aim is to decompose  $f$  into the sum of sinusoids that are independent in some sense. Thus we begin by trying to approximate  $f$  by a sum of sinusoids of periods  $1, 1/2, \dots, 1/n$  for an arbitrary positive integer  $n$ .

**Theorem A.2.1** FOURIER SERIES APPROXIMATION.

Let  $f$  be a square integrable function defined on  $[0, 1]$ ; that is  $\int_0^1 f^2(x) dx < \infty$ . Let

$$g_j(x) = a_j \cos(2\pi jx) + b_j \sin(2\pi jx) = C_j \cos(2\pi jx - \phi_j),$$

be a sinusoid of period  $1/j$ . Let

$$f_n(x) = \frac{a_0}{2} + \sum_{j=1}^n g_j(x),$$

where  $a_0$  is some constant. Then the values of  $a_0, a_1, \dots, a_n$  and  $b_1, b_2, \dots, b_n$  that minimize

$$\int_0^1 \left| f(x) - \frac{a_0}{2} - \sum_{j=1}^n g_j(x) \right|^2 dx$$

are given by

$$a_j = 2 \int_0^1 \cos(2\pi jx) f(x) dx, \quad j = 0, 1, \dots, n \quad (\text{A.10})$$

$$b_j = 2 \int_0^1 \sin(2\pi jx) f(x) dx, \quad j = 1, \dots, n. \quad (\text{A.11})$$

The theorem yields that the sum of sinusoids that is closest to  $f$  in the integrated square error sense (which is analogous to the sum of square errors in regression analysis) has the coefficients given in equations (A.10) and (A.11). Note that the function  $f$  need not be square integrable in order to calculate the coefficients. In fact,  $f$  needs only to be absolutely integrable; that is  $\int_0^1 |f(x)| dx < \infty$ , since (for example)

$$\int_0^1 |f(x) \cos(2\pi jx)| dx \leq \int_0^1 |f(x)| |\cos(2\pi jx)| dx \leq \int_0^1 |f(x)| dx < \infty,$$

since  $|\cos(2\pi jx)| \leq 1$ . Note that a square integrable function is absolutely integrable, but not necessarily vice versa. Consider, for example,  $f(x) = 1/\sqrt{x}$ .

**Definition A.2.1** FOURIER COEFFICIENTS & PARTIAL SUM FOURIER APPROXIMANT.

Let  $f$  be an absolutely integrable function defined on  $[0, 1]$ , and let

$$\begin{aligned} a_j &= 2 \int_0^1 \cos(2\pi jx) f(x) dx, \quad j = 0, 1, \dots, \\ b_j &= 2 \int_0^1 \sin(2\pi jx) f(x) dx, \quad j = 1, 2, \dots, \end{aligned}$$

and

$$f_n(x) = \frac{a_0}{2} + \sum_{j=1}^n [a_j \cos(2\pi jx) + b_j \sin(2\pi jx)].$$

Then the  $a_j$  and  $b_j$  are called the **cosine** and **sine Fourier coefficients** of  $f$ , while  $f_n$  is called the  $n^{\text{th}}$  **partial sum Fourier approximant** to  $f$ .

In general, we will have that  $f \neq f_n$ . To accomplish harmonic analysis of  $f$  then, we hope that  $f_n \rightarrow f$  in some sense, and also that there is some sense in which the sinusoids  $g_1, g_2, \dots$  are orthogonal. In the next theorem, we present results showing that in many cases, we do have the convergence and orthogonality.

**Theorem A.2.2** PROPERTIES OF FOURIER SERIES.

Let  $f$  be an absolutely integrable function defined on the interval  $[0, 1]$  and let  $f_1, f_2, \dots$  be it approximating Fourier sums. Then

(a) The sinusoids  $f_j(x) = a_j \cos(2\pi jx) + b_j \sin(2\pi jx)$  are orthogonal in the sense that

$$\int_0^1 g_j(x) g_k(x) dx = \begin{cases} 0, & j \neq k \\ a_j^2 + b_j^2, & j = k. \end{cases}$$

(b) We have the following results about pointwise convergence of  $f_n(x)$ .

(1) At a point  $x$  where  $f$  is continuous on the left and right derivatives of  $f$  exist,

$$\lim_{n \rightarrow \infty} f_n(x) = f(x).$$

(2) At a point  $x$  where  $f$  is discontinuous but the left and right derivatives of  $f$  exist,

$$\lim_{n \rightarrow \infty} f_n(x) = \frac{f(x^+) + f(x^-)}{2},$$

where  $f(x^+)$  and  $f(x^-)$  denote the right and left limits of  $f$  at the point  $x$ ; that is, the approximating Fourier sums converge to the average of the two values of  $f$  at the discontinuity.

(3) If  $f$  is continuous and has square integrable derivative, then  $f_n \rightarrow f$  absolutely and uniformly.

(c) If  $f$  is square integrable, then

$$\lim_{n \rightarrow \infty} \int_0^1 |f(x) - f_n(x)|^2 dx = 0,$$

while

$$2 \int_0^1 f^2(x) dx = \frac{a_0^2}{2} + \sum_{k=1}^{\infty} (a_k^2 + b_k^2).$$

**Implications:** This theorem gives us our harmonic analysis for two types of functions – absolutely integrable (part (b)) and square integrable (part (c)). For square integrable functions  $f_n \rightarrow f$  in the integrated squared error sense, which need not mean that  $f_n(x)$  converges pointwise to  $f(x)$ . Note that the second equation in part (c) (which is called Parseval’s identity) is analogous to a decomposition of “total sum of squares” in analysis of variance, where the integral plays the role of the sum, into a sum of squares of amplitudes of the sinusoids for periods  $1, 1/2, 1/3, \dots$ . For absolutely integrable functions, part (b1) shows that under general conditions,  $f_n(x)$  does converge to  $f(x)$  pointwise, while part (b2) explains the behavior of the Fourier approximations for the ideal bandpass filter in Section 4.3. Thus at the two points of discontinuity in the transfer function,  $f_n(x)$  converges to the average of 0 and 1; that is to 0.5. Finally, part (a) shows that as long as the Fourier coefficients exist, then the sinusoids are orthogonal where now the inner product is defined as an integral rather than as a sum as it was in Section 2.2.3.

## The Complex Form of Fourier Series

Recall Euler’s equation given in equation (A.7)

$$e^{i\theta} = \cos \theta + i \sin \theta,$$

along with equations (A.8) and (A.9). Applying these results to the previous discussion provides a complex form of Fourier series. Specifically, write

$$\cos(2\pi jx) = \frac{e^{2\pi i jx} + e^{-2\pi i jx}}{2} \quad \text{and} \quad \sin(2\pi jx) = \frac{(-i)(e^{2\pi i jx} - e^{-2\pi i jx})}{2},$$

which gives

$$f_n(x) = \frac{a_0}{2} + \sum_{j=1}^n [a_j \cos(2\pi jx) + b_j \sin(2\pi jx)] = \sum_{j=-n}^n r_j e^{-2\pi i jx},$$

where

$$r_j = \frac{a_j + ib_j}{2} \quad \text{and} \quad r_{-j} = \bar{r}_j = \frac{a_j - ib_j}{2}, \quad j = 0, 1, \dots$$

Note that

$$r_j = \int_0^1 f(x) e^{2\pi i jx} dx.$$

Thus we have that  $f$  and the  $r$ s are Fourier pairs and

$$f(x) = \sum_{j=-\infty}^{\infty} r_j e^{-2\pi i j x} \quad \text{and} \quad r_j = \int_0^1 f(x) e^{2\pi i j x} dx.$$

The  $r_j$  are called the complex Fourier coefficients of  $f$ .

## A.2.2 Spectral Representation of $R$

In Appendix A.2.1, we were concerned with whether a function could be represented as a sum of sinusoids. In time series, we have the converse problem; that is, we begin with the autocovariance sequence and ask whether there is a function having the  $R$ s as its Fourier coefficients. The simplest, although not the most general, answer to this question is given in the following theorem.

### **Theorem A.2.3** SPECTRAL REPRESENTATION OF $R$ .

*If  $\{R(v) : v \in \mathcal{Z}\}$  is a positive definite sequence of numbers such that  $R(v) = R(-v)$  and  $\sum_{v=-\infty}^{\infty} |R(v)| < \infty$ , then there exists a continuous function  $f$  such that*

$$(a) \quad f(\omega) = \sum_{v=-\infty}^{\infty} R(v) e^{-2\pi i v \omega} \quad \text{for } \omega \in [0, 1].$$

$$(b) \quad R(v) = \int_0^1 f(\omega) e^{2\pi i v \omega} d\omega \quad \text{for } v \in \mathcal{Z}.$$

$$(c) \quad f(\omega) = f(1 - \omega).$$

This theorem gives a harmonic analysis of an absolutely summable autocovariance function in terms of the spectral density function  $f$ . Unfortunately, the covariance function is not always absolutely summable, as illustrated by a harmonic process in Section 5.2. We can still obtain a harmonic analysis of a non-absolutely summable  $R$  as an integral, but the integral is of a special type called a Lebesgue-Stieltjes integral. In the spectral representation of  $R$  in Theorem 4.2 we wrote

$$R_v = \int_0^1 \cos(2\pi v \omega) dF(\omega),$$

where the integral is of the Lebesgue-Stieltjes type. Instead of defining this type of integral, which is rather complicated, we describe in the next theorem a way to calculate and interpret it in the special cases with which we are concerned.

### **Theorem A.2.4** LEBESGUE-STIELTJES INTEGRAL.

*Let  $g$  and  $F$  be bounded functions defined on  $[0, 1]$  and for a positive integer  $m$ , consider partitioning  $[0, 1]$  into  $2^m$  intervals*

$$\left[0, \frac{1}{2^m}\right], \quad \left[\frac{1}{2^m}, \frac{2}{2^m}\right], \quad \dots, \quad \left[\frac{2^m - 1}{2^m}, 1\right].$$

Thus if we let  $\omega_{m,j} = j/2^m$ , for  $j = 0, 1, \dots, 2^m$ , the  $k^{\text{th}}$  interval is from  $\omega_{m,k-1}$  to  $\omega_{m,k}$ , for  $k = 1, \dots, 2^m$ . Now let

$$I_m = \sum_{k=1}^{2^m} g(\omega'_{m,k}) [F(\omega_{m,k}) - F(\omega_{m,k-1})],$$

where  $\omega'_{m,k} = (\omega_{m,k-1} + \omega_{m,k})/2$ . Then if  $g$  is Lebesgue-Stieltjes integrable with respect to  $F$ , and  $g$  is continuous (and therefore also Riemann-Stieltjes integrable), we have that  $\lim_{m \rightarrow \infty} I_m$  exists, is finite, and

$$\int_0^1 g(x) dF(x) = \lim_{m \rightarrow \infty} I_m.$$

From this result we have that

$$R_v = \int_0^1 \cos(2\pi v\omega) dF(\omega) = \lim_{m \rightarrow \infty} \sum_{k=1}^{2^m} \cos(2\pi v\omega'_{m,k}) [F(\omega_{m,k}) - F(\omega_{m,k-1})],$$

which means that the autocovariance sequence  $R$  can be regarded as the sum of a large number of cosines with amplitudes given by “increments” of the spectral distribution function  $F$ . Thus, if there is a large increase in  $F$  between frequencies  $\omega_{m,k-1}$  and  $\omega_{m,k}$ , the sinusoid having frequency  $\omega'_{m,k}$  will be prominent in the decomposition of  $R$ . In fact, if the only large increase in  $F$  is in the interval containing  $\omega'_{m,k}$ , then  $R$  will appear almost exactly sinusoidal. The logical limit of this idea is the autocovariance of the harmonic process.

## The Harmonic Process

Consider the harmonic process having a single frequency component of  $\omega_0 \in (0, 0.5)$ ; that is,  $R_v = \sigma^2 \cos(2\pi v\omega_0)$ . Define

$$F(\omega) = \begin{cases} 0, & 0 \leq \omega < \omega_0 \\ \frac{\sigma^2}{2}, & \omega_0 \leq \omega < 1 - \omega_0 \\ \sigma^2 & 1 - \omega_0 \leq \omega \leq 1, \end{cases}$$

and let  $\omega_m$  be the center of the interval in the partition of  $[0, 1]$  at the  $m^{\text{th}}$  step of the limiting process. If  $\omega_0$  is on the boundary between two intervals, we choose  $\omega_m$  to be the center of the one containing the jump in  $F$ . Thus not that  $1 - \omega_m$  is in the center of the interval containing  $1 - \omega_0$  and we have

$$\begin{aligned} & \lim_{m \rightarrow \infty} \sum_{k=1}^{2^m} \cos(w\pi v\omega'_{m,k}) [F(\omega_{m,k}) - F(\omega_{m,k-1})] \\ &= \lim_{m \rightarrow \infty} \left[ \frac{\sigma^2}{2} \cos\{2\pi v\omega_m\} + \frac{\sigma^2}{2} \cos\{2\pi v(1 - \omega_m)\} \right] \\ &= \sigma^2 \cos(2\pi v\omega_0), \end{aligned}$$

since  $\omega_m \rightarrow \omega_0$  and  $\cos(2\pi v\omega) = \cos\{2\pi v(1 - \omega)\}$ . Thus this definition of the spectral distribution function does indeed lead to the correct autocovariance function. We note that in the general theory of representing sequences, the function  $F$  is unique, and thus our definition of  $F$  is the only one that represents  $R$ .

In general, if  $F$  is absolutely continuous, we can see that the Lebesgue-Stieltjes integral reduces to the usual Riemann integral since then

$$\begin{aligned} \int_0^1 \cos(w\pi v\omega) dF(\omega) &= \lim_{m \rightarrow \infty} \frac{1}{2^m} \sum_{k=1}^{2^m} \cos(2\pi v\omega'_{m,k}) \left[ \frac{F(\omega_{m,k}) - F(\omega_{m,k-1})}{1/2^m} \right] \\ &\doteq \lim_{m \rightarrow \infty} \frac{1}{2^m} \sum_{k=1}^{2^m} \cos(2\pi v\omega'_{m,k}) \frac{dF(\omega'_{m,k})}{d\omega} \\ &= \int_0^1 \cos(2\pi v\omega) f(\omega) d\omega \end{aligned}$$

by the definition of the Riemann integral.

### A.2.3 Spectral Representation of $X$

It is possible to obtain a harmonic analysis of the time series  $X$  itself (see part (4) of Theorem 4.2). The representation is different than that for  $R$  in three ways. First, the representation of  $R$  is in terms of cosines only. This is because  $R$  is an even function; that is  $R(v) = R(-v)$ . Second, the increments  $F(\omega_k) - F(\omega_j)$  are replaced by increments of two uncorrelated stochastic processes (more on this below). Finally, since the representation will be as a limit of sums of linear combinations of random variables, the usual limit cannot be used, rather the limit is a limit in mean square (see Definition e in Appendix A.6).

A stochastic process  $\{C(\omega), \omega \in [0, 1]\}$  is an indexed collection of random variables. The random variable  $C(\omega_2) - C(\omega_1)$  for  $0 \leq \omega_1 < \omega_2 \leq 1$  is called an **increment** of  $C$ . If increments defined on non-overlapping intervals are uncorrelated, we say that  $C$  has **uncorrelated increments**.

Thus the spectral representation of  $X$  means essentially that

$$\begin{aligned} X_t = \text{l.i.m.} \sum_{k=1}^{2^m} & \left[ \cos(2\pi v\omega'_{m,k}) \{C(\omega_{m,k}) - C(\omega_{m,k-1})\} \right. \\ & \left. + \sin(2\pi v\omega'_{m,k}) \{S(\omega_{m,k}) - S(\omega_{m,k-1})\} \right] \end{aligned}$$

where l.i.m. denotes limit in mean square. This again allows us to think of a particular infinitely long realization from  $X$  as the sum of sinusoids whose amplitudes have been chosen according to some random mechanism, determined by the processes  $C$  and  $S$ , whose properties (such as variance) are described by the spectral distribution function  $F$ .

## A.3 The Fast Fourier Transform

Many of the calculations in time series analysis require the calculation of the discrete Fourier transform of a set of numbers (see Section 2.2.3 for example). In this section we describe the fast Fourier transform (FFT). The basic ideas were developed and popularized by Cooley & Tukey (1965) and Gentleman & Sande (1966). The FFT is a classic example of the use of recursion to greatly reduce the number of operations required to solve a computational problem. Although there are a wide variety of algorithms known as FFTs, almost all of them have two essential elements: (1) a recursion that allows one to find the DFT of  $n$  points as a simple combination of DFTs for subsets of the points, and (2) a reordering of the original data so as to minimize the number of complex exponentials that need to be calculated and/or stored while the recursion is being performed.

### A.3.1 The FFT Recursion

Let  $W_n = e^{2\pi i/n}$ . Then we can write the DFT of  $X_1, \dots, X_n$  as

$$Z_k = \sum_{t=1}^n X_t W_n^{(t-1)(k-1)}, \quad k = 1, \dots, n.$$

**Theorem A.3.1** THE FFT RECURSION.

*Suppose  $n$  can be factored into  $n = rs$ . Define the  $r \times s$  matrix  $\mathbf{V}$  to have  $j^{\text{th}}$  column*

$$X_j, X_{s+j}, \dots, X_{(r-1)s+j}, \quad j = 1, \dots, s.$$

*Then*

$$Z_k = \sum_{\ell=1}^s W_n^{(k-1)(\ell-1)} Z_{\ell, \text{mod}(k-1, r)+1},$$

*where  $Z_{\ell,1}, \dots, Z_{\ell,r}$  is the DFT of the  $\ell^{\text{th}}$  column of  $\mathbf{V}$ .*

**Proof:** We can get  $Z_k$  by adding over the elements of  $\mathbf{V}$ :

$$\begin{aligned} Z_k &= \sum_{\ell=1}^s \sum_{m=1}^r V_{m\ell} W_n^{(k-1)[(m-1)s+(\ell-1)]} \\ &= \sum_{\ell=1}^s W_n^{(k-1)(\ell-1)} \sum_{m=1}^r V_{m\ell} W_n^{(k-1)(m-1)s} \end{aligned}$$

since  $V_{m\ell} = X_{(m-1)s+\ell}$ . But  $s/n = 1/r$  which gives  $W_n^s = W_r$ . Furthermore

$$k-1 = \left\lfloor \frac{k-1}{r} \right\rfloor r + (k-1) - \left\lfloor \frac{k-1}{r} \right\rfloor r = \left\lfloor \frac{k-1}{r} \right\rfloor r + \text{mod}(k-1, r),$$

and thus

$$W_n^{(k-1)(m-1)s} = W_r^{(m-1)\{\lfloor \frac{k-1}{r} \rfloor r + \text{mod}(k-1, r)\}} = W_r^{(m-1)\text{mod}(k-1, r)},$$

since  $W_r^{\left[\frac{k-1}{r}\right]r} = 1$  and  $e^{2\pi i} = 1$ . Thus

$$\begin{aligned} Z_k &= \sum_{\ell=1}^s W_n^{(k-1)(\ell-1)} \sum_{m=1}^r V_{m\ell} W_r^{(m-1)\text{mod}(k-1,r)} \\ &= \sum_{\ell=1}^s W_n^{(k-1)(\ell-1)} Z_{\ell, \text{mod}(k-1,r)+1}. \end{aligned}$$

Thus it takes  $sr^2 = nr$  operations to calculate DFTs of all  $s$  columns of  $\mathbf{V}$  and then  $ns$  operations to calculate all  $n$   $Z$ s, making a total of  $nr + ns = n(r + s)$  operations instead of  $n^2$ . Further, if  $r$  is itself factorable, we can use the same idea to more rapidly find the DFTs of the columns of  $\mathbf{V}$ . This process can be continued until all of the prime factors of  $n$  are exhausted, giving a total of  $n(p_1 + \dots + p_K)$  operations, where  $p_1, \dots, p_K$  are the prime factors of  $n$ .

### A.3.2 Reordering Data

The recursions given in Theorem A.3.1 can be further improved if we can order the calculations in such a way that all of the one using a certain complex exponential can be done at the same time. To illustrate this, consider the case of  $n = 8$  and  $s = 2$ :

$$\begin{aligned} Z_1 &= Z_{1,1} + W_8^0 Z_{2,1} \\ Z_2 &= Z_{1,2} + W_8^1 Z_{2,2} \\ Z_3 &= Z_{1,3} + W_8^2 Z_{2,3} \\ Z_4 &= Z_{1,4} + W_8^3 Z_{2,4} \\ Z_5 &= Z_{1,1} + W_8^4 Z_{2,1} \\ Z_6 &= Z_{1,2} + W_8^5 Z_{2,2} \\ Z_7 &= Z_{1,3} + W_8^6 Z_{2,3} \\ Z_8 &= Z_{1,4} + W_8^7 Z_{2,4}, \end{aligned}$$

where  $Z_{1,\cdot}$  is the DFT of the odd numbered  $X$ s and  $Z_{2,\cdot}$  is the DFT of the even numbered  $X$ s. Thus

$$\begin{aligned} Z_{1,1} &= Z_{1,1}^{(1)} + W_4^0 Z_{1,1}^{(2)} & Z_{2,1} &= Z_{2,1}^{(1)} + W_4^0 Z_{2,1}^{(2)} \\ Z_{1,2} &= Z_{1,2}^{(2)} + W_4^1 Z_{1,2}^{(2)} & Z_{2,2} &= Z_{2,2}^{(1)} + W_4^1 Z_{2,2}^{(2)} \\ Z_{1,3} &= Z_{1,1}^{(1)} + W_4^2 Z_{1,1}^{(2)} & Z_{2,3} &= Z_{2,1}^{(1)} + W_4^2 Z_{2,1}^{(2)} \\ Z_{1,4} &= Z_{1,2}^{(2)} + W_4^3 Z_{1,2}^{(2)} & Z_{2,4} &= Z_{2,2}^{(1)} + W_4^3 Z_{2,2}^{(2)} \end{aligned}$$

where

$$\begin{aligned} Z_{1,1}^{(1)} &= X_1 + W_2^0 X_5 \\ Z_{1,2}^{(1)} &= X_1 + W_2^1 X_5 \\ Z_{1,1}^{(2)} &= X_3 + W_2^0 X_7 \end{aligned}$$



$$\begin{aligned}
Z_{1,2}^{(2)} &= X_3 + W_2^1 X_7 \\
Z_{2,1}^{(1)} &= X_2 + W_2^0 X_6 \\
Z_{2,2}^{(1)} &= X_2 + W_2^1 X_6 \\
Z_{2,1}^{(2)} &= X_4 + W_2^0 X_8 \\
Z_{2,2}^{(2)} &= X_4 + W_2^1 X_8.
\end{aligned}$$

Thus if we start by permuting the  $X$ s, we can summarize the calculations by

$$\begin{bmatrix} X_1 \\ X_5 \\ X_3 \\ X_7 \\ X_2 \\ X_6 \\ X_4 \\ X_8 \end{bmatrix} \longrightarrow \begin{bmatrix} Z_{1,1}^{(1)} \\ Z_{1,1}^{(2)} \\ Z_{2,1}^{(1)} \\ Z_{2,1}^{(2)} \\ Z_{1,2}^{(1)} \\ Z_{1,2}^{(2)} \\ Z_{2,2}^{(1)} \\ Z_{2,2}^{(2)} \end{bmatrix} \longrightarrow \begin{bmatrix} Z_{1,1} \\ Z_{2,1} \\ Z_{1,2} \\ Z_{2,2} \\ Z_{1,3} \\ Z_{2,3} \\ Z_{1,4} \\ Z_{2,4} \end{bmatrix} \longrightarrow \begin{bmatrix} Z_1 \\ Z_2 \\ Z_3 \\ Z_4 \\ Z_5 \\ Z_6 \\ Z_7 \\ Z_8 \end{bmatrix},$$

where the first arrow means first combine all four consecutive pairs using  $W_2^0$  and then again using  $W_2^1$ . The second arrow means combine the first two pairs within  $W_4^0$ , the next two using  $W_4^1$ , the next two using  $W_4^2$ , and the last two using  $W_4^3$ . In the last step we combine the first pair using  $W_8^0$ , the next pair using  $W_8^1$ , and so on. Thus at each step we are fixing the complex exponential and then doing all combinations using that value.

To determine the index  $k$  of  $X_j$  in the permuted array, one need only write down the binary representation of  $j-1$ , reverse the bits, and then add one. For example,  $X_5$  is mapped to position 2 since  $5-1=4$ , and the binary representation of 4 is  $100 \rightarrow 001 \rightarrow 010 = 2$ . Note that if  $X_j$  gets mapped to position  $k$ , then  $X_k$  gets mapped to position  $j$ .

There are many variations on this theme and further simplifications that can be made. For some of the functions in *Timeslab in R*, we use the algorithm described by Singleton (1969) which is in some ways still the standard by which newer methods are judged. We chose this algorithm because it performs the transform in place and will work for a wide variety of different prime factors, although as we saw in Section 2.2.4, it does require that  $n$  not have any large prime factors. Other functions in *Timeslab in R* use an FFT that is the default in *R*. This version is a modification of the Singleton (1969) FFT that was presented in Singleton (1979).

## A.4 Matrix Operations

In this section we consider three operations on matrices: (1) the modified Cholesky decomposition of a positive definite matrix, (2) the Gram-Schmidt decomposition of a nonsingular  $n \times m$  matrix, and (3) the matrix sweep operator. Each of these is useful as both a computational and theoretical tool.

### A.4.1 The Modified Cholesky Decomposition

An  $n \times n$  matrix  $\mathbf{V}$  is said to be positive definite or positive semidefinite if for any nonzero  $n$ -dimensional vector  $\mathbf{h}$ , we have  $\mathbf{h}^T \mathbf{V} \mathbf{h} > 0$  or  $\mathbf{h}^T \mathbf{V} \mathbf{h} \geq 0$  respectively. In statistics, it is often necessary to find what is called the square root of a positive definite matrix.

**Definition A.4.1** SQUARE ROOT OF A MATRIX.

The square root of a symmetric, positive definite  $n \times n$  matrix  $\mathbf{V}$  is an  $n \times n$  matrix  $\mathbf{A}$  satisfying  $\mathbf{V} = \mathbf{A} \mathbf{A}^T$ . We denote such a matrix by  $\mathbf{V}^{1/2}$ . The inverse square root of  $\mathbf{V}$  is the inverse of  $\mathbf{V}^{1/2}$  and is denoted  $\mathbf{V}^{-1/2}$ . The transposes of  $\mathbf{V}^{-1}$  and  $\mathbf{V}^{-1/2}$  are denoted  $\mathbf{V}^{-T}$  and  $\mathbf{V}^{-T/2}$  respectively.

The next theorem shows that there exists a unique matrix square root for any positive definite matrix and also shows how to find it. Note that a triangular matrix is called **unit** if it has ones on the main diagonal.

**Theorem A.4.1** MODIFIED CHOLESKY DECOMPOSITION.

Let  $\mathbf{V}$  be a symmetric  $n \times n$  matrix. Then

- (a)  $\mathbf{V}$  is positive definite if and only if there exists a unique unit lower triangular  $n \times n$  matrix  $\mathbf{L}$  and a unique diagonal  $n \times n$  matrix  $\mathbf{D}$  having positive diagonal elements, such that

$$\mathbf{V} = \mathbf{L} \mathbf{D} \mathbf{L}^T.$$

- (b) This factorization is called the modified Cholesky decomposition (MCD) of  $\mathbf{V}$ , and if the decomposition exists we can calculate the elements  $\mathbf{L}$  and  $\mathbf{D}$  one row at a time by  $D_{11} = V_{11}$  and for  $i = 2, \dots, n$

$$L_{ij} = \frac{V_{ij} - \sum_{\ell=1}^{j-1} L_{i\ell} D_{\ell\ell} L_{j\ell}}{D_{jj}}, \quad j = 1, \dots, i-1$$

$$D_{ii} = V_{ii} - \sum_{\ell=1}^{i-1} D_{\ell\ell} L_{i\ell}^2.$$

- (c) The MCD of  $\mathbf{V}$  is nested; that is, if  $\mathbf{V}_k, \mathbf{L}_k$ , and  $\mathbf{D}_k$  represent the upper left-hand  $k \times k$  parts of  $\mathbf{V}, \mathbf{L}$  and  $\mathbf{D}$ , respectively, then

$$\mathbf{V}_k = \mathbf{L}_k \mathbf{D}_k \mathbf{L}_k^T, \quad k = 1, \dots, n,$$

and thus for  $k$  greater than or equal to  $i$  and  $j$ , the  $(i, j)^{th}$  elements of  $\mathbf{V}_k, \mathbf{L}_k$  and  $\mathbf{D}_k$  can be denoted by  $V_{ij}, L_{ij}$ , and  $D_{ij}$ , respectively.

- (d) The unique square root of a positive definite matrix  $\mathbf{V}$  is given by  $\mathbf{V}^{1/2} = \mathbf{L} \mathbf{D}^{1/2}$  where  $\mathbf{D}^{1/2} = \text{Diag}(D_{11}^{1/2}, \dots, D_{nn}^{1/2})$ .
- (e) The inverse square root of a positive definite matrix  $\mathbf{V}$  is given by  $\mathbf{V}^{-1/2} = \mathbf{D}^{-1/2} \mathbf{L}^{-1}$  where  $\mathbf{D}^{-1/2} = \text{Diag}(D_{11}^{-1/2}, \dots, D_{nn}^{-1/2})$ .

## The mchol Function

If  $\mathbf{A}$  is a symmetric  $n \times n$  matrix ( $\mathbf{A}$ ), then the function

`mchol(A)`

will attempt to find the modified Cholesky factors  $\mathbf{L}$  and  $\mathbf{D}$  of  $\mathbf{A}$ . It returns a list containing three objects: `error`,  $\mathbf{L}$ , and  $\mathbf{D}$ . If the algorithm judges that  $\mathbf{A}$  is not positive definite, that is, if a diagonal element  $\mathbf{D}$  is found that is less than  $10^{-25}$ , the function will return an error value of `error = 1` with null values for  $\mathbf{L}$  and  $\mathbf{D}$ . If  $\mathbf{A}$  is judged to be positive definite, then `mchol` will return `error = 0`, and will return the factors  $\mathbf{L}$  and  $\mathbf{D}$  in the matrices  $\mathbf{L}$  and  $\mathbf{D}$ . This in addition to finding the decomposition in possible, `mchol` provides an easy check for positive definiteness of a matrix.

## A.4.2 The Gram-Schmidt Decomposition

An important operation in many areas of mathematics and statistics is to find what is called an orthogonal basis for a matrix; that is, from one set of vectors, find a new set by some linear transformation such that the inner product of any two different new vectors is zero. We will use the Gram-Schmidt decomposition (see Clayton (1971)) to accomplish this aim.

**Theorem A.4.2** GRAM-SCHMIDT DECOMPOSITION.

*If  $\mathbf{X}$  is an  $n \times p$  matrix having full rank  $p$ , then there exists an  $n \times p$  orthogonal matrix  $\mathbf{Q}$  (that is,  $\mathbf{Q}^T \mathbf{Q}$  is diagonal) and a unit upper triangular  $p \times p$  matrix  $\mathbf{R}$  such that*

$$\mathbf{X} = \mathbf{Q}\mathbf{R}.$$

*This decomposition of  $\mathbf{X}$  is called its **Gram-Schmidt decomposition** (GSD).*

Note that  $\mathbf{X}^T \mathbf{X} = \mathbf{R}^T \mathbf{Q}^T \mathbf{Q} \mathbf{R} = \mathbf{R}^T \mathbf{V} \mathbf{R}$  where  $\mathbf{R}^T$  is unit lower triangular and  $\mathbf{V}$  is diagonal with positive diagonal elements. Since the modified Cholesky decomposition  $\mathbf{X}^T \mathbf{X} = \mathbf{L} \mathbf{D} \mathbf{L}^T$  of  $\mathbf{X}^T \mathbf{X}$  is unique, we have that  $\mathbf{R}^T = \mathbf{L}$  and  $\mathbf{V} = \mathbf{D}$ .

## The gram.schmidt Function

If the matrix  $\mathbf{A}$  is  $n \times p$ , then the command

`gram.schmidt(A)`

will attempt to find the factors returns a list containing three objects: `error`,  $\mathbf{Q}$ , and  $\mathbf{R}$ . If there is an error, then `error = 1`, and the values of  $\mathbf{Q}$  and  $\mathbf{R}$  will be null. If there is no error, then `error = 0`, and  $\mathbf{Q}$  will contain the  $\mathbf{Q}$  factor, and  $\mathbf{R}$  the  $\mathbf{R}$  factor of the GSD.

The function uses what is known as the “Gram-Schmidt decomposition algorithm”. In other words, if we let  $\mathbf{Q}_0 = \mathbf{X}$ , and  $\mathbf{Q}_1, \dots, \mathbf{Q}_p = \mathbf{Q}$  be a sequence of  $n \times p$  matrices, and

$\mathbf{q}_{ij}$  denote the  $j^{\text{th}}$  column of  $\mathbf{Q}_i$ , then `gram.schmidt` calculates for  $i = 1, \dots, p-1$

$$\begin{aligned} d_i &= \mathbf{q}_{i-1,i}^T \mathbf{q}_{i-1,i} \\ R_{ij} &= \frac{\mathbf{q}_{i-1,j}^T \mathbf{q}_{i-1,i}}{d_i}, \quad j = i+1, \dots, p \\ \mathbf{q}_{ij} &= \mathbf{q}_{i-1,j} - R_{ij} \mathbf{q}_{i-1,i}, \quad j = i+1, \dots, p. \end{aligned}$$

If  $|d_i| < 10^{-25}$  for any  $i$ , then `gram.schmidt` concludes that  $\mathbf{X}$  is singular, and an error will occur (`error = 1`).

### A.4.3 The Sweep Operator

Many of the computational and theoretical results in regression, multivariate analysis, and time series can be expressed succinctly using what is called the sweep operator.

**Definition A.4.2** SWEEP OPERATOR.

Let  $\mathbf{A}$  be an  $n \times n$  matrix. The process of sweeping  $\mathbf{A}$  on its  $k^{\text{th}}$  diagonal element, denoted by  $\mathbf{B} = \text{SWEEP}(k) \mathbf{A}$ , is the process of forming the matrix  $\mathbf{B}$  by (assuming  $A_{kk} \neq 0$ )

$$\begin{aligned} B_{kk} &= \frac{1}{A_{kk}} \\ B_{ik} &= -\frac{A_{ik}}{A_{kk}}, \quad i \neq k \quad (k^{\text{th}} \text{ column}) \\ B_{kj} &= \frac{A_{kj}}{A_{kk}}, \quad j \neq k \quad (k^{\text{th}} \text{ row}) \\ B_{ij} &= A_{ij} - \frac{A_{ik}A_{kj}}{A_{kk}}, \quad i \neq k, j \neq k. \end{aligned}$$

If  $\mathbf{B}_1 = \text{SWEEP}(k_1) \mathbf{A}$ ,  $\mathbf{B}_2 = \text{SWEEP}(k_2) \mathbf{B}_1, \dots, \mathbf{B} = \text{SWEEP}(k_r) \mathbf{B}_{r-1}$ , we write  $\mathbf{B} = \text{SWEEP}(k_1, \dots, k_r) \mathbf{A}$ , and say that  $\mathbf{B}$  is the result of sweeping  $\mathbf{A}$  on its diagonals  $k_1, \dots, k_r$ .

The sweep operator has many useful properties, as we see in the next theorem.

**Theorem A.4.3** PROPERTIES OF THE SWEEP OPERATOR.

Let

$$\mathbf{A} = \begin{bmatrix} \mathbf{B} & \mathbf{C} \\ \mathbf{D} & \mathbf{E} \end{bmatrix},$$

where  $\mathbf{B}, \mathbf{C}, \mathbf{D}$ , and  $\mathbf{E}$  are  $r \times r, r \times s, s \times r$ , and  $s \times s$ , respectively. Then

(a)  $\text{SWEEP}(k_1, \dots, k_p) \mathbf{A}$  can be found by sweeping the diagonals in any order, for example,

$$\text{SWEEP}(k_1, k_2) \mathbf{A} = \text{SWEEP}(k_2, k_1) \mathbf{A}.$$

(b)  $\text{SWEEP}(k, k) \mathbf{A} = \mathbf{A}$ ; that is, sweeping a second time on a given diagonal undoes the effect of a previous sweeping on that diagonal.

(c)  $\text{SWEEP}(1, \dots, r) \mathbf{A} = \begin{bmatrix} \mathbf{B}^{-1} & \mathbf{B}^{-1}\mathbf{C} \\ -\mathbf{DB}^{-1} & \mathbf{E} - \mathbf{DB}^{-1}\mathbf{C} \end{bmatrix}$  if  $\mathbf{B}$  is nonsingular.

(d)  $\text{SWEEP}(r+1, \dots, r+s) \mathbf{A} = \begin{bmatrix} \mathbf{B} - \mathbf{CE}^{-1}\mathbf{D} & -\mathbf{CE}^{-1} \\ \mathbf{E}^{-1}\mathbf{D} & \mathbf{E}^{-1} \end{bmatrix}$  if  $\mathbf{E}$  is nonsingular.

(e)  $\text{SWEEP}(1, \dots, r+s) \mathbf{A} = \mathbf{A}^{-1}$  and

$$\mathbf{A}^{-1} = \begin{bmatrix} \mathbf{B}^{-1} + \mathbf{B}^{-1}\mathbf{C}(\mathbf{E} - \mathbf{DB}^{-1}\mathbf{C})^{-1}\mathbf{DB}^{-1} & -\mathbf{B}^{-1}\mathbf{C}(\mathbf{E} - \mathbf{DB}^{-1}\mathbf{C})^{-1} \\ -(\mathbf{E} - \mathbf{DB}^{-1}\mathbf{C})^{-1}\mathbf{DB}^{-1} & (\mathbf{E} - \mathbf{DB}^{-1}\mathbf{C})^{-1} \end{bmatrix}$$

if  $\mathbf{B}$ ,  $\mathbf{E}$ , and  $(\mathbf{E} - \mathbf{DB}^{-1}\mathbf{C})^{-1}$  are nonsingular.

The sweep operator can be used to solve many theoretical and computational problems. We consider two examples of its use. First, the basic quantities for the regression of  $\mathbf{y}$  on  $\mathbf{X}$  (where  $\mathbf{X}$  is  $n \times p$ ) can be found by applying property (c) (see Appendix A.5).

$$\text{SWEEP}(1, \dots, p) \begin{bmatrix} \mathbf{X}^T\mathbf{X} & \mathbf{X}^T\mathbf{y} \\ \mathbf{y}^T\mathbf{X} & \mathbf{y}^T\mathbf{y} \end{bmatrix} = \begin{bmatrix} \frac{1}{\sigma^2} \text{Var}(\hat{\beta}) & \hat{\beta} \\ -\hat{\beta}^T & \text{RSS} \end{bmatrix}.$$

To illustrate the theoretical utility of sweep, notice for matrix  $\mathbf{A}$  in Theorem A.4.3 that sweeping on all of its diagonals is the same as sweeping the matrix  $\text{SWEEP}(r+1, \dots, r+s) \mathbf{A}$  on diagonals  $1, \dots, r$ . Thus the upper left-hand corner in part (e) is the inverse of the upper left-hand corner in part (d) by doing part (c). Doing a similar calculation with  $\pm \mathbf{E}^{-1}$  replacing  $\mathbf{E}$  in  $\mathbf{A}$  gives the well-known matrix inversion lemma (see Rao (1973), p. 33, for example).

**Theorem A.4.4** MATRIX INVERSION LEMMA.

Let  $\mathbf{B}$  and  $\mathbf{E}$  be nonsingular  $r \times r$  and  $s \times s$  matrices and  $\mathbf{C}$  and  $\mathbf{D}$  be  $r \times s$  and  $s \times r$  matrices. Then

$$(\mathbf{B} \pm \mathbf{CED})^{-1} = \mathbf{B}^{-1} \mp \mathbf{B}^{-1}\mathbf{C}(\mathbf{E}^{-1} \pm \mathbf{DB}^{-1}\mathbf{C})^{-1}\mathbf{DB}^{-1}.$$

We will use this formula extensively when we consider recursive regression in Appendix A.5, that is, when we consider adjusting estimators for the addition or deletion of observations. It can also be used to motivate the Kalman Filter Algorithm discussed in Appendix A.6.

## The swp Function

The **swp** function is called using

$$\text{swp}(\mathbf{A}, \mathbf{k1} = 1, \mathbf{k2} = \text{ncol}(\mathbf{A}))$$

where  $\mathbf{A}$  is a matrix,  $\mathbf{k1}$  is the index of the first diagonal element of  $\mathbf{A}$  on which to begin sweeping, and  $\mathbf{k2}$  is the index of the last diagonal. The default for  $\mathbf{k1} = 1$ , and for  $\mathbf{k2}$  is the last column of  $\mathbf{A}$ . The function returns a list containing the error indicator **error**, and a matrix  $\mathbf{A}$ . If **error** = 0, then  $\mathbf{A}$  is the result of sweeping the input matrix on diagonals  $\mathbf{k1}$  through  $\mathbf{k2}$  and **error**. If there is a diagonal element that is less than  $10^{-25}$ , then **error** = 1, and the matrix  $\mathbf{A}$  will be null.

## A.5 Least Squares Regression Analysis

### A.5.1 The General Linear Model

### A.5.2 Results for the General Linear Model

### A.5.3 The Case of Correlated Errors

### A.5.4 Special Types of Regression Analysis

#### Orthogonal Regression

#### Stepwise Regression

#### Recursive Regression

An important topic in statistics is recursive estimation. In regression analysis, this means given the results  $\hat{\beta}$ ,  $RSS$ , and  $\mathbf{V} = (\mathbf{X}^T \mathbf{X})^{-1}$  of regressing the  $m \times 1$  vector  $\mathbf{y}$  on the  $m \times p$  matrix  $\mathbf{X}$ , we seek the corresponding quantities after adding or deleting data  $\mathbf{y}_1$ , and  $\mathbf{X}_1$  of length  $n$ . Thus we seek

$$\begin{aligned}\hat{\beta}_N &= (\mathbf{X}^T \mathbf{X} \pm \mathbf{X}_1^T \mathbf{X}_1)^{-1} (\mathbf{X}^T \mathbf{y} \pm \mathbf{X}_1^T \mathbf{y}_1) \\ RSS_N &= (\mathbf{y}^T \mathbf{y} \pm \mathbf{y}_1^T \mathbf{y}_1) - \hat{\beta}_N^T (\mathbf{X}^T \mathbf{y} \pm \mathbf{X}_1^T \mathbf{y}_1) \\ \mathbf{V}_n &= (\mathbf{X}^T \mathbf{X} \pm \mathbf{X}_1^T \mathbf{X}_1)^{-1}\end{aligned}$$

where the  $\pm$  is  $+$  for adding observations and  $-$  for deleting. From the matrix inversion lemma (Theorem A.4.4) we can obtain recursive formulas for these quantities.

#### **Theorem A.5.1** RECURSIVE REGRESSION FORMULAS.

*The updated quantities  $\hat{\beta}_N$ ,  $RSS_N$ , and  $\mathbf{V}_N$  defined above are given by*

$$\begin{aligned}\hat{\beta}_N &= \hat{\beta} \pm \mathbf{A} \mathbf{X}_1^T (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} \mathbf{e}_1 \\ &= \hat{\beta} \pm (\mathbf{X}^T \mathbf{X} \pm \mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{e}_1 \\ RSS_N &= RSS \pm \mathbf{e}_1^T (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} \mathbf{e}_1 \\ &= RSS \pm \mathbf{e}_1^T \left[ \mathbf{I}_n \mp \mathbf{X}_1 (\mathbf{X}^T \mathbf{X} \pm \mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \right] \mathbf{e}_1 \\ \mathbf{V}_N &= \mathbf{A} \mp \mathbf{A} \mathbf{X}_1^T (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} \mathbf{X}_1 \mathbf{A} \\ &= (\mathbf{X}^T \mathbf{X} \pm \mathbf{X}_1^T \mathbf{X}_1)^{-1},\end{aligned}$$

where

$$\mathbf{H}_1 = \mathbf{X}_1 \mathbf{A} \mathbf{X}_1^T, \quad \mathbf{A} = (\mathbf{X}^T \mathbf{X})^{-1}, \quad \mathbf{e}_1 = \mathbf{y}_1 - \mathbf{X}_1 \hat{\beta}.$$

**Proof:** The basis of this theorem is the Matrix Inversion Lemma (Theorem A.4.4) and the facts below, easily verified by multiplication.

$$(\mathbf{I}_n \pm \mathbf{H}_1)^{-1} = \mathbf{I}_n \mp (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} \mathbf{H}_1 = \mathbf{I}_n \mp \mathbf{H}_1 (\mathbf{I}_n \pm \mathbf{H}_1)^{-1}.$$

Thus

$$\begin{aligned} \hat{\boldsymbol{\beta}}_N &= (\mathbf{X}^T \mathbf{X} \pm \mathbf{X}_1^T \mathbf{X}_1)^{-1} (\mathbf{X}^T \mathbf{y} \pm \mathbf{X}_1^T \mathbf{y}_1) \\ &= [\mathbf{A} \mp \mathbf{A} \mathbf{X}_1^T (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} \mathbf{X}_1 \mathbf{A}] (\mathbf{X}^T \mathbf{y} \pm \mathbf{X}_1^T \mathbf{y}_1) \\ &= \hat{\boldsymbol{\beta}} \mp \mathbf{A} \mathbf{X}_1^T (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} \hat{\mathbf{y}}_1 \pm \mathbf{A} \mathbf{X}_1^T \mathbf{y}_1 - \mathbf{Z} \mathbf{X}_1^T (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} \mathbf{H}_1 \mathbf{y}_1, \end{aligned}$$

where  $\hat{\mathbf{y}}_1 = \mathbf{X}_1 \mathbf{A} \mathbf{X}^T \mathbf{y} = \mathbf{X}_1 \hat{\boldsymbol{\beta}}$ , which gives

$$\begin{aligned} \hat{\boldsymbol{\beta}}_N - \hat{\boldsymbol{\beta}} &= \pm \mathbf{A} \mathbf{X}_1^T [\mathbf{I}_n \mp (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} \mathbf{H}_1] \mathbf{y}_1 \mp \mathbf{A} \mathbf{X}_1^T (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} \hat{\mathbf{y}}_1 \\ &= \pm \mathbf{A} \mathbf{X}_1^T (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} (\mathbf{y}_1 - \hat{\mathbf{y}}_1), \end{aligned}$$

which gives the first expression for  $\hat{\boldsymbol{\beta}}_N$ . To show the second, note that

$$\begin{aligned} (\mathbf{X}^T \mathbf{X} \pm \mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{e}_1 &= [\mathbf{A} \mp \mathbf{A} \mathbf{X}_1^T (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} \mathbf{X}_1 \mathbf{A}] \mathbf{X}_1^T \mathbf{e}_1 \\ &= \mathbf{A} \mathbf{X}_1^T [\mathbf{I}_n \mp (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} \mathbf{H}_1] \mathbf{e}_1 \\ &= \mathbf{A} \mathbf{X}_1^T (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} \mathbf{e}_1. \end{aligned}$$

To find  $RSS_N$ , we have that  $\pm \mathbf{H}_1 \mathbf{y}_1 = (\mathbf{I}_n \pm \mathbf{H}_1) \mathbf{y}_1 - \mathbf{y}_1$ , and so

$$\begin{aligned} RSS_N &= (\mathbf{y}^T \mathbf{y} \pm \mathbf{y}_1^T \mathbf{y}_1) - [\hat{\boldsymbol{\beta}}^T \pm \mathbf{e}_1^T (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} \mathbf{X}_1 \mathbf{A}] (\mathbf{X}^T \mathbf{y} \pm \mathbf{X}_1^T \mathbf{y}_1) \\ &= RSS \pm \mathbf{y}_1^T \mathbf{y}_1 \mp \hat{\boldsymbol{\beta}}^T \mathbf{X}_1^T \mathbf{y}_1 \mp \mathbf{e}_1^T (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} (\mathbf{X}_1 \hat{\boldsymbol{\beta}} \pm \mathbf{H}_1 \mathbf{y}_1) \\ &= RSS \pm \mathbf{y}_1^T \mathbf{y}_1 \mp \hat{\boldsymbol{\beta}}^T \mathbf{X}_1^T \mathbf{y}_1 \mp \mathbf{e}_1^T (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} \mathbf{y}_1^T \\ &\quad \mp \mathbf{e}_1^T (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} [(\mathbf{I}_n \pm \mathbf{H}_1) \mathbf{y}_1 - \mathbf{y}_1] \\ &= RSS \pm (\mathbf{y}_1^T - \hat{\boldsymbol{\beta}}^T \mathbf{X}_1^T - \mathbf{e}_1^T) \mathbf{y}_1 \mp \mathbf{e}_1^T (\mathbf{I}_n \pm \mathbf{H}_1) (\mathbf{y}_1^T - \mathbf{y}_1) \\ &= RSS \pm \mathbf{e}_1^T (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} \mathbf{e}_1, \end{aligned}$$

since  $\mathbf{e}_1^T = \mathbf{y}_1^T - \hat{\boldsymbol{\beta}}^T \mathbf{X}_1^T$ . To obtain the second expression for  $RSS_N$ , note that

$$\begin{aligned} &\mathbf{e}_1^T [\mathbf{I}_n \mp \mathbf{X}_1 (\mathbf{X}^T \mathbf{X} \pm \mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T] \mathbf{e}_1 \\ &= \mathbf{e}_1^T [\mathbf{I}_n \mp \mathbf{X}_1 \{ \mathbf{A} \mp \mathbf{A} \mathbf{X}_1^T (\mathbf{I}_n \mp \mathbf{H}_1)^{-1} \mathbf{X}_1 \mathbf{A} \} \mathbf{X}_1^T] \mathbf{e}_1 \\ &= \mathbf{e}_1^T [\mathbf{I}_n \mp \mathbf{X}_1 \mathbf{A} \mathbf{X}_1^T + \mathbf{X}_1 \mathbf{A} \mathbf{X}_1^T (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} \mathbf{X}_1 \mathbf{A} \mathbf{X}_1^T] \mathbf{e}_1 \\ &= \mathbf{e}_1^T [\mathbf{I}_n \mp \mathbf{H}_1 \{ \mathbf{I}_n \mp (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} \mathbf{H}_1 \}] \mathbf{e}_1 \\ &= \mathbf{e}_1^T [\mathbf{I}_n \mp \mathbf{H}_1 (\mathbf{I}_n \pm \mathbf{H}_1)^{-1}] \mathbf{e}_1 \\ &= \mathbf{e}_1^T (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} \mathbf{e}_1. \end{aligned}$$

Finally, the first expression for  $\mathbf{V}_n$  follows from the Matrix Inversion Lemma (Theorem A.4.4) while the second is by definition.

Note that in each pair of equations for the updated quantities, the first equation is appropriate when  $n \leq p$  while the second is for  $n \geq p$ . The sweep operator can be used in either case. For  $n \leq p$ , we have

$$\begin{aligned} & \text{SWEEP}(1, \dots, p+n) \begin{bmatrix} \mathbf{X}^T \mathbf{X} & \mathbf{X}_1^T \\ \mp \mathbf{X}_1 & \mathbf{I}_n \end{bmatrix} = \begin{bmatrix} \mathbf{B} & \mathbf{C} \\ \mathbf{D} & \mathbf{E} \end{bmatrix} \\ & = \begin{bmatrix} \mathbf{A} \mp \mathbf{X}_1^T (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} \mathbf{X}_1 \mathbf{A} & -\mathbf{A} \mathbf{X}_1^T (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} \\ \pm (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} \mathbf{X}_1 \mathbf{A} & (\mathbf{I}_n \pm \mathbf{H}_1)^{-1} \end{bmatrix} \end{aligned}$$

and thus

$$\hat{\boldsymbol{\beta}}_N = \hat{\boldsymbol{\beta}} \mp \mathbf{C} \mathbf{e}_1, \quad RSS_N = RSS \pm \mathbf{e}_1^T \mathbf{E} \mathbf{e}_1, \quad \text{and} \quad \mathbf{V}_N = \mathbf{B}.$$

For  $n \geq p$

$$\begin{aligned} & \text{SWEEP}(1, \dots, p) \begin{bmatrix} \mathbf{X}^T \mathbf{X} \pm \mathbf{X}_1^T \mathbf{X}_1 & \mathbf{X}_1^T \mathbf{e}_1 \\ \pm \mathbf{e}_1^T \mathbf{X}_1 & \mathbf{e}_1^T \mathbf{e}_1 \end{bmatrix} = \begin{bmatrix} \mathbf{U} & \mathbf{v} \\ \mathbf{w}^T & z \end{bmatrix} \\ & \begin{bmatrix} (\mathbf{X}^T \mathbf{X} \pm \mathbf{X}_1^T \mathbf{X}_1)^{-1} & (\mathbf{X}^T \mathbf{X} \pm \mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{e}_1 \\ \mp \mathbf{e}_1^T \mathbf{X}_1 (\mathbf{X}^T \mathbf{X} \pm \mathbf{X}_1^T \mathbf{X}_1)^{-1} & \mathbf{e}_1^T \left\{ \mathbf{I}_n \mp \mathbf{X}_1 (\mathbf{X}^T \mathbf{X} \pm \mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1 \right\} \mathbf{e}_1 \end{bmatrix}. \end{aligned}$$

and thus

$$\hat{\boldsymbol{\beta}}_N = \hat{\boldsymbol{\beta}} \pm \mathbf{v}, \quad RSS_N = RSS \pm z, \quad \text{and} \quad \mathbf{V}_n = \mathbf{U}.$$

## A.6 Random Vectors and Multivariate Normal Distribution

In elementary statistics, we often avoid treating vectors of random variables since the variables are independent and can be treated separately. In time series analysis, this is usually not the case.

### A.6.1 Basic Definitions

We begin our discussion with a series of definitions.

1. A  **$d$ -dimension random vector**  $\mathbf{X}$  is a vector of  $d$  random variables; that is,  $\mathbf{X} = (X_1, \dots, X_d)^T$ , where  $X_1, \dots, X_d$  are random variables.
2. The **joint cumulative distribution function** (cdf)  $F_{\mathbf{X}}$  of the random vector  $\mathbf{X}$  is given by

$$F_{\mathbf{X}}(\mathbf{x}) = F_{\mathbf{X}}(x_1, \dots, x_d) = \Pr(X_1 \leq x_1, \dots, X_d \leq x_d).$$



3. If there exists a function  $f_X$  of  $d$  variables such that

$$F_X(\mathbf{x}) = \int_{-\infty}^{x_1} \cdots \int_{-\infty}^{x_d} f_X(\mathbf{y}) d\mathbf{y},$$

then  $\mathbf{X}$  is said to be a **continuous rand vector** and  $f_X$  is called the **joint probability density function** (pdf) of  $\mathbf{X}$ .

4. If  $\mathbf{Y}$  and  $\mathbf{X}$  are  $r$ -dimensional and  $d$ -dimensional random vectors, then  $\mathbf{Z} = (\mathbf{Y}^T, \mathbf{X}^T)^T$  is an  $(r + d)$ -dimensional random vector. If  $f_Z, f_Y$ , and  $f_X$  are the pdfs of  $\mathbf{Z}, \mathbf{Y}$ , and  $\mathbf{X}$ , respectively, then the functions

$$f_{Y|X} = \frac{f_Z}{f_X} \quad \text{and} \quad f_{X|Y} = \frac{f_Z}{f_Y}$$

are called the **conditional pdfs of  $\mathbf{Y}$  given  $\mathbf{X}$  and  $\mathbf{X}$  given  $\mathbf{Y}$** , respectively.

5. The **mean vector**  $\boldsymbol{\mu}_X$  and **covariance matrix**  $\boldsymbol{\Sigma}_X$  of  $\mathbf{X}$  are given by

$$\begin{aligned} \boldsymbol{\mu}_X &= E[\mathbf{X}] = \begin{bmatrix} E[X_1] \\ \vdots \\ E[X_d] \end{bmatrix} = \int \cdots \int \mathbf{x} f_X(\mathbf{x}) d\mathbf{x}, \\ \boldsymbol{\Sigma}_X &= \text{Var}(\mathbf{X}) = E[(\mathbf{X} - \boldsymbol{\mu}_X)(\mathbf{X} - \boldsymbol{\mu}_X)^T] \\ &= \int \cdots \int (\mathbf{x} - \boldsymbol{\mu}_X)(\mathbf{x} - \boldsymbol{\mu}_X)^T f_X(\mathbf{x}) d\mathbf{x} \\ &= \begin{bmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) & \text{Cov}(X_1, X_3) & \cdots & \text{Cov}(X_1, X_d) \\ & \text{Var}(X_2) & \text{Cov}(X_2, X_3) & \cdots & \text{Cov}(X_2, X_d) \\ & & \text{Var}(X_3) & \cdots & \vdots \\ & & & \ddots & \vdots \\ & & & & \text{Var}(X_d) \end{bmatrix}. \end{aligned}$$

If  $E[\mathbf{X}] = \boldsymbol{\mu}_X$  and  $\text{Var}(\mathbf{X}) = \boldsymbol{\Sigma}_X$  we will write

$$\mathbf{X} \sim (\boldsymbol{\mu}_X, \boldsymbol{\Sigma}_X).$$

The covariance matrix  $\boldsymbol{\Sigma}_{XY}$  of  $\mathbf{X}$  and  $\mathbf{Y}$  is the  $d \times r$  matrix whose  $(j, k)^{\text{th}}$  element is  $\text{Cov}(X_j, Y_k)$ . The conditional mean  $\boldsymbol{\mu}_{Y|X}$  and variance  $\boldsymbol{\Sigma}_{Y|X}$  of a random vector  $\mathbf{Y}$  given  $\mathbf{X}$  are defined in the same way as  $\boldsymbol{\mu}_Y$  and  $\boldsymbol{\Sigma}_Y$ , using the conditional pdf of  $\mathbf{Y}$  given  $\mathbf{X}$ .

6. If we partition the  $(r + d)$ -dimensional random vector  $\mathbf{Z}$  into the  $r$ -dimensional and  $d$ -dimensional random vectors  $\mathbf{Y}$  and  $\mathbf{X}$  consisting of the first  $r$  and last  $d$  elements of  $\mathbf{Z}$ , then we can partition  $\boldsymbol{\mu}_Z$  and  $\boldsymbol{\Sigma}_Z$  similarly as

$$\boldsymbol{\mu}_Z = \begin{bmatrix} \boldsymbol{\mu}_Y \\ \boldsymbol{\mu}_X \end{bmatrix} \quad \text{and} \quad \boldsymbol{\Sigma}_Z = \begin{bmatrix} \boldsymbol{\Sigma}_{YY} & \boldsymbol{\Sigma}_{YX} \\ \boldsymbol{\Sigma}_{XY} & \boldsymbol{\Sigma}_{XX} \end{bmatrix},$$

where in fact

$$\begin{aligned}\Sigma_{YY} &= \text{Var}(\mathbf{Y}), & \Sigma_{YX} &= \text{Cov}(\mathbf{Y}, \mathbf{X}) = \Sigma_{XY}^T, \\ \Sigma_{XX} &= \text{Var}(\mathbf{X}), & \Sigma_{XY} &= \text{Cov}(\mathbf{X}, \mathbf{Y}) = \Sigma_{YX}^T.\end{aligned}$$

7. The  $r$ -dimensional random vector  $\mathbf{Y}$  is called a **linear transformation of the  $d$ -dimensional random vector  $\mathbf{X}$**  if  $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{b}$  where  $\mathbf{A}$  is an  $r \times d$  matrix of constants  $\mathbf{b}$  is an  $r$ -dimensional vector of constants.
8. If  $\mathbf{X}$  is a  $d$ -dimensional random vector and  $\mathbf{A}$  is a  $d \times d$  matrix of constants, then the scalar random variable  $Q = \mathbf{X}^T \mathbf{A} \mathbf{X}$  is called a **quadratic form in  $\mathbf{X}$**  and  $\mathbf{A}$  is called the **matrix of the quadratic form**.

### A.6.2 Basic Facts

With these definitions in mind, we next state some theoretical results that are used throughout the book. We begin with the following theorem.

**Theorem A.6.1** MEAN AND VARIANCE OF A LINEAR FUNCTION.

If  $\mathbf{X} \sim (\boldsymbol{\mu}_X, \Sigma_X)$  then

$$\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{b} \sim (\mathbf{A}\boldsymbol{\mu}_X + \mathbf{b}, \mathbf{A}\Sigma_X\mathbf{A}^T),$$

that is,

$$\begin{aligned}E[\mathbf{A}\mathbf{X} + \mathbf{b}] &= \mathbf{A}E[\mathbf{X}] + \mathbf{b} \\ \text{Var}(\mathbf{A}\mathbf{X} + \mathbf{b}) &= \mathbf{A}\text{Var}(\mathbf{X})\mathbf{A}^T.\end{aligned}$$

Most of the theory and methods of time series analysis are built upon the properties of the multivariate normal distribution.

**Definition A.6.1** MULTIVARIATE NORMAL DISTRIBUTION.

The  $d$ -dimensional random vector  $\mathbf{X}$  having mean  $\boldsymbol{\mu}_X$  and variance  $\Sigma_X$  is said to have the  **$d$ -dimensional normal distribution** if its pdf is given by

$$f_X(\mathbf{x}) = (2\pi)^{-d/2} |\Sigma_X|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_X)^T \Sigma_X^{-1} (\mathbf{x} - \boldsymbol{\mu}_X) \right\},$$

for  $\mathbf{x} \in \Re^d$ , that is, in  $d$ -dimensional Euclidean space. Such a random vector is denoted by  $\mathbf{X} \sim N_d(\boldsymbol{\mu}_X, \Sigma_X)$ .

The important properties of the multivariate normal distribution are summaries in Theorem A.6.2 (see Rao (1973), Chapter 8, for example).

**Theorem A.6.2** PROPERTIES OF THE MULTIVARIATE NORMAL.

If

$$\mathbf{Z} = \begin{pmatrix} \mathbf{Y} \\ \mathbf{X} \end{pmatrix} \sim N_{r+d} \left( \begin{bmatrix} \boldsymbol{\mu}_Y \\ \boldsymbol{\mu}_X \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{YY} & \boldsymbol{\Sigma}_{YX} \\ \boldsymbol{\Sigma}_{XY} & \boldsymbol{\Sigma}_{XX} \end{bmatrix} \right),$$

then

(a) The characteristic function of  $\mathbf{Z}$  is given by

$$\phi_Z(\mathbf{t}) = E \left[ e^{i\mathbf{t}^T \mathbf{Z}} \right] = \exp \left( i\mathbf{t}^T \boldsymbol{\mu}_Z - \frac{1}{2} \mathbf{t}^T \boldsymbol{\Sigma}_Z \mathbf{t} \right).$$

(b) If  $\mathbf{A}$  is an  $s \times (r+d)$  matrix of constants and  $\mathbf{b}$  is an  $s$ -dimensional vector of constants, then

$$\mathbf{AZ} + \mathbf{b} \sim N_s \left( \mathbf{A}\boldsymbol{\mu}_Z + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}_Z \mathbf{A}^T \right);$$

that is, linear transformations of normal random vectors are also normally distributed.

(c) (i) The conditional distribution of  $\mathbf{Y}$  given  $\mathbf{X}$  is

$$\mathbf{Y} | \mathbf{X} \sim N_r \left( \boldsymbol{\mu}_Y + \boldsymbol{\Sigma}_{YX} \boldsymbol{\Sigma}_{XX}^{-1} (\mathbf{X} - \boldsymbol{\mu}_X), \boldsymbol{\Sigma}_{YY} - \boldsymbol{\Sigma}_{YX} \boldsymbol{\Sigma}_{XX}^{-1} \boldsymbol{\Sigma}_{XY} \right);$$

that is, the pdf of  $\mathbf{Y}$  given  $\mathbf{X}$  is that of an  $r$ -dimensional normal random vector with mean and variance given by

$$\begin{aligned} \boldsymbol{\mu}_{Y|X} &= E[\mathbf{Y} | \mathbf{X}] = \boldsymbol{\mu}_Y + \boldsymbol{\Sigma}_{YX} \boldsymbol{\Sigma}_{XX}^{-1} (\mathbf{X} - \boldsymbol{\mu}_X) \\ \boldsymbol{\Sigma}_{Y|X} &= \text{Var}(\mathbf{Y} | \mathbf{X}) = \boldsymbol{\Sigma}_{YY} - \boldsymbol{\Sigma}_{YX} \boldsymbol{\Sigma}_{XX}^{-1} \boldsymbol{\Sigma}_{XY}; \end{aligned}$$

that is the condition mean of  $\mathbf{Y}$  given  $\mathbf{X} = \mathbf{x}$  is a linear function of  $\mathbf{x}$ , while the conditional variance is the same, no matter what the value of  $\mathbf{x}$  is.

(ii) Conversely, if

$$\mathbf{Y} | \mathbf{X} \sim N_r \left( \boldsymbol{\mu}_Y + \boldsymbol{\Sigma}_{YX} \boldsymbol{\Sigma}_{XX}^{-1} (\mathbf{X} - \boldsymbol{\mu}_X), \boldsymbol{\Sigma}_{YY} - \boldsymbol{\Sigma}_{YX} \boldsymbol{\Sigma}_{XX}^{-1} \boldsymbol{\Sigma}_{XY} \right),$$

and  $\mathbf{X} \sim N_d(\boldsymbol{\mu}_X, \boldsymbol{\Sigma}_{XX})$ , then

$$\mathbf{Z} = \begin{pmatrix} \mathbf{Y} \\ \mathbf{X} \end{pmatrix} \sim N_{r+d} \left( \begin{bmatrix} \boldsymbol{\mu}_X \\ \boldsymbol{\mu}_Y \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{YY} & \boldsymbol{\Sigma}_{YX} \\ \boldsymbol{\Sigma}_{XY} & \boldsymbol{\Sigma}_{XX} \end{bmatrix} \right).$$

(d) If  $\mathbf{A}$  is a  $d \times d$  nonsingular matrix of constants, then the distribution of  $\mathbf{Y} | \mathbf{X}$  and  $\mathbf{Y} | \mathbf{AX}$  are the same.

(e)  $\text{Cov}(\mathbf{AX}, \mathbf{BY}) = \mathbf{A}\boldsymbol{\Sigma}_{XY} \mathbf{B}^T$ .

(f) The quadratic form

$$Q = (\mathbf{X} - \boldsymbol{\mu}_X)^T \boldsymbol{\Sigma}_{XX}^{-1} (\mathbf{X} - \boldsymbol{\mu}_X)$$

has the chi-square distribution with  $d$  degrees of freedom, and we write  $Q \sim \chi_d^2$ .

(g) The quadratic form

$$Q_A = (\mathbf{X} - \boldsymbol{\mu}_X)^T \mathbf{A} (\mathbf{X} - \boldsymbol{\mu}_X)$$

has the same distribution as the random variable  $\sum_{j=1}^d \lambda_j Z_j$ , where  $Z_1, \dots, Z_s$  are independent  $\chi_1^2$  random variables, and the  $\lambda_i$  are the non-zero eigenvalues of  $\mathbf{A}\boldsymbol{\Sigma}_{XX}$ . This reduces to the result of part (f) if  $\mathbf{A} = \boldsymbol{\Sigma}_{XX}^{-1}$  since the eigenvalues of the  $d$ -dimensional identity matrix are just one repeated  $d$  times, and the sum of  $d$  independent  $\chi_1^2$  random variables is  $\chi_d^2$ .

(h) The mean and variance of a quadratic form in normal variables are given by

$$\begin{aligned} E[\mathbf{X}^T \mathbf{A} \mathbf{X}] &= \text{tr}(\mathbf{A}\boldsymbol{\Sigma}_{XX}) + \boldsymbol{\mu}_X^T \mathbf{A} \boldsymbol{\mu}_X \\ \text{Var}(\mathbf{X}^T \mathbf{A} \mathbf{X}) &= 2\text{tr}([\mathbf{A}\boldsymbol{\Sigma}_{XX}]^2) + 4\boldsymbol{\mu}_X^T \mathbf{A} \boldsymbol{\Sigma}_{XX} \mathbf{A}^T \boldsymbol{\mu}_X. \end{aligned}$$

(i) The two quadratic forms  $Q_1 = \mathbf{X}^T \mathbf{A} \mathbf{X}$  and  $Q_2 = \mathbf{X}^T \mathbf{B} \mathbf{X}$  are independent if and only if  $\mathbf{A}\boldsymbol{\Sigma}_{XX} \mathbf{B}^T = \mathbf{0}$ .

(j) If  $Q_1$  and  $Q_2$  are two independent  $\chi^2$  random variables having  $n_1$  and  $n_2$  degrees of freedom, then the random variable

$$F = \frac{Q_1/n_1}{Q_2/n_2}$$

has the  $F$  distribution with  $n_1$  numerator degrees of freedom and  $n_2$  denominator degrees of freedom, and we write  $F \sim F_{n_1, n_2}$ .

(k) If  $Z \sim N_1(\mu, \sigma^2)$ , and  $Q \sim \chi_v^2$  are independent, then the random variable

$$t = \frac{(Z - \mu)/\sigma}{\sqrt{Q/v}}$$

has the Student's  $t$  distribution with  $v$  degrees of freedom and we write  $t \sim t_v$ .

## Computing the pdf, cdf, and Quantile Functions

The  $R$  software has functions for computing the pdf, cdf, and quantile functions for the  $N(\mu, \sigma^2)$ ,  $t_v$ ,  $\chi_v^2$ , and  $F_{v_1, v_2}$  distributions. They are

For the normal distribution with  $\mu = \text{mean}$  and standard deviation  $\sigma = \text{sd}$ , the pdf at  $\mathbf{x}$  is computed using `dnorm(x, mean, sd)`, the cdf at  $\mathbf{x}$  is computed using `pnorm(x, mean, sd)`, and the  $p^{\text{th}}$  quantile is computed using `qnorm(p, mean, sd)`. For all three functions, the default values for the mean and standard deviation are `mean = 0` and `sd = 1`.

- For the Student's  $t_v$  distribution with degrees of freedom  $v = \text{df}$ , the pdf at  $\mathbf{x}$  is computed using `dt(x, df)`, the cdf at  $\mathbf{x}$  is computed using `pt(x, df)`, and the  $p^{\text{th}}$  quantile is computed using `qt(p, df)`.

- For the  $\chi_v^2$  with degrees of freedom  $v = \text{df}$ , the pdf at  $\mathbf{x}$  is computed using `dchisq(x, df)`, the cdf at  $\mathbf{x}$  is computed using `pchisq(x, df)`, and the  $p^{\text{th}}$  quantile is computed using `qchisq(p, df)`.
- For the  $F_{v_1, v_2}$  with numerator degrees of freedom  $v_1 = \text{df1}$ , and denominator degrees of freedom  $v_2 = \text{df2}$  the pdf at  $\mathbf{x}$  is computed using `df(x, df1, df2)`, the cdf at  $\mathbf{x}$  is computed using `pf(x, df1, df2)`, and the  $p^{\text{th}}$  quantile is computed using `qf(p, df1, df2)`.

These functions may be computed for many other distributions in *R* using the same type of syntax. For a comprehensive list and documentations, issue the command `help("Distributions")` at the *R* prompt.

### Simultaneous Confidence Intervals

A result that is useful in finding simultaneous confidence bands for functions is given in the next theorem, due to Scheffé (1959), p. 406.

**Theorem A.6.3** SIMULTANEOUS CONFIDENCE INTERVALS.

*Let  $\mathbf{X}$  be an arbitrary  $p$ -dimensional random vector and let  $\mathbf{V}$  be a positive definite  $p \times p$  matrix of constants. Then the  $p$ -dimensional vector  $\mathbf{a}$  satisfies*

$$(\mathbf{X} - \mathbf{a})^T \mathbf{V} (\mathbf{X} - \mathbf{a}) \leq 1$$

*if and only if*

$$|\mathbf{h}^T (\mathbf{X} - \mathbf{a})| \leq (\mathbf{h}^T \mathbf{V}^{-1} \mathbf{h})^{-1/2} = D$$

*for all  $p$ -dimensional vectors  $\mathbf{h}$ ; that is,*

$$\mathbf{h}^T \mathbf{X} - D \leq \mathbf{h}^T \mathbf{a} \leq \mathbf{h}^T \mathbf{X} + D.$$

### A.6.3 Prediction of a Random Vector

Consider random vectors  $\mathbf{Y}$  and  $\mathbf{X}$  such that  $\mathbf{Y}$  has  $d$  elements,  $\mathbf{X}$  has  $r$  elements, and

$$E \left[ \begin{pmatrix} \mathbf{Y} \\ \mathbf{X} \end{pmatrix} \right] = \begin{bmatrix} \boldsymbol{\mu}_Y \\ \boldsymbol{\mu}_X \end{bmatrix} \quad \text{Var} \left( \begin{bmatrix} \mathbf{Y} \\ \mathbf{X} \end{bmatrix} \right) = \begin{bmatrix} \boldsymbol{\Sigma}_{YY} & \boldsymbol{\Sigma}_{YX} \\ \boldsymbol{\Sigma}_{XY} & \boldsymbol{\Sigma}_{XX} \end{bmatrix}.$$

We consider to criteria for predicting  $\mathbf{Y}$  from  $\mathbf{X}$ : (1) best unbiased prediction and (2) best linear unbiased prediction.

**Definition A.6.2** BEST UNBIASED PREDICTOR & BEST LINEAR UNBIASED PREDICTOR.

*The best unbiased predictor of  $\mathbf{Y}$  from  $\mathbf{X}$  is that function  $\tilde{\mathbf{Y}}$  of  $\mathbf{X}$  that has the same expectation as  $\mathbf{Y}$  and is closest to  $\mathbf{Y}$  in expected mean square; that is,*

$$E[(\mathbf{Y} - \mathbf{Y}^*)(\mathbf{Y} - \mathbf{Y}^*)^T] - E[(\mathbf{Y} - \tilde{\mathbf{Y}})(\mathbf{Y} - \tilde{\mathbf{Y}})^T]$$

is positive semidefinite for any other unbiased function  $\mathbf{Y}^*$  of  $\mathbf{X}$ . The **best linear unbiased predictor**  $\hat{\mathbf{Y}}$  of  $\mathbf{Y}$  given  $\mathbf{X}$  is that linear function of  $\mathbf{X}$  having the same expectation of  $\mathbf{Y}$  and is closest to  $\mathbf{Y}$  in mean squared error; that is,  $\hat{\mathbf{Y}} = \hat{\mathbf{A}}\mathbf{X}$  where

$$E[(\mathbf{Y} - \mathbf{A}^*\mathbf{X})(\mathbf{Y} - \mathbf{A}^*\mathbf{X})^T] - E[(\mathbf{Y} - \hat{\mathbf{A}}\mathbf{X})(\mathbf{Y} - \hat{\mathbf{A}}\mathbf{X})^T]$$

is positive semidefinite for any other linear function  $\mathbf{A}^*\mathbf{X}$  that is also unbiased for  $\mathbf{Y}$ .

It will turn out that  $\tilde{\mathbf{Y}}$  is the conditional expectation of  $\mathbf{Y}$  given  $\mathbf{X}$ . Conditional expectation is a crucial concept in time series analysis and can be defined at many levels of generality. A more general and useful definition of conditional expectation that is given in Definition 5 of Appendix A.6.1 is given by the following (see Parzen (1962b), for example).

**Definition A.6.3** CONDITIONAL EXPECTATION OF A RANDOM VECTOR.

Let  $\mathbf{Y}$  be a random vector having finite mean. Then the **conditional expectation**  $E[\mathbf{Y} | \mathbf{X}]$  of  $\mathbf{Y}$  given  $\mathbf{X}$  is that function  $\phi$  of  $\mathbf{X}$  such that

$$E[\{\mathbf{Y} - \phi(\mathbf{X})\}\mathbf{g}(\mathbf{X})] = \mathbf{0}$$

for any bounded function  $\mathbf{g}$  of  $\mathbf{X}$ .

This definition is called descriptive (as opposed to the constructive definition of Appendix A.6.1) as it says that the “residual”  $\mathbf{Y} - E[\mathbf{Y} | \mathbf{X}]$  is uncorrelated with any other function of  $\mathbf{X}$ . In this sense, there is no more information about  $\mathbf{Y}$  in  $\mathbf{X}$  after having removed  $E[\mathbf{Y} | \mathbf{X}]$ . We note that the conditional variance  $\text{Var}(\mathbf{Y} | \mathbf{X})$  of  $\mathbf{Y}$  given  $\mathbf{X}$  is just the conditional expectation of the random variable  $[(\mathbf{Y} - E[\mathbf{Y} | \mathbf{X}])(\mathbf{Y} - E[\mathbf{Y} | \mathbf{X}])^T]$  given  $\mathbf{X}$ .

Letting  $\mathbf{g}(\mathbf{X}) = \mathbf{1}$  in this definition yields  $E[\phi(\mathbf{X})] = E[\mathbf{Y}]$ , and thus

$$\text{Cov}(\mathbf{Y} - \phi(\mathbf{X}), \mathbf{g}(\mathbf{X})) = \mathbf{0}.$$

With these definitions in mind, we have the following theorem. For simplicity, we assume that the random vectors involved have zero mean.

**Theorem A.6.4** GENERAL PREDICTION THEORY.

Let  $\mathbf{Y}$  and  $\mathbf{X}$  be jointly distributed random vectors having zero mean. Then

- (a) The best unbiased predictor  $\tilde{\mathbf{Y}}$  of  $\mathbf{Y}$  given  $\mathbf{X}$  and its error variance  $\Sigma_{\tilde{\mathbf{Y}}} = \text{Var}(\mathbf{Y} - \tilde{\mathbf{Y}})$  are given by

$$\tilde{\mathbf{Y}} = E[\mathbf{Y} | \mathbf{X}] \quad \text{and} \quad \Sigma_{\tilde{\mathbf{Y}}} = \text{Var}(\mathbf{Y} | \mathbf{X}).$$

- (b) If  $\Sigma_{XX}$  is nonsingular, then the best linear unbiased predictor  $\hat{\mathbf{Y}}$  of  $\mathbf{Y}$  given  $\mathbf{X}$  and its error variance  $\Sigma_{\hat{\mathbf{Y}}} = \text{Var}(\mathbf{Y} - \hat{\mathbf{Y}})$  are given by

$$\hat{\mathbf{Y}} = \hat{\mathbf{A}}\mathbf{X} \quad \text{and} \quad \Sigma_{\hat{\mathbf{Y}}} = \Sigma_{YY} - \Sigma_{YX}\Sigma_{XX}^{-1}\Sigma_{XY},$$

where  $\hat{\mathbf{A}}$  satisfies the prediction normal equations

$$\Sigma_{XX}\hat{\mathbf{A}}^T = \Sigma_{XY}.$$

(c) If  $\mathbf{Y}$  and  $\mathbf{X}$  are jointly multivariate normal, then  $\tilde{\mathbf{Y}} = \hat{\mathbf{Y}}$  and  $\Sigma_{\tilde{\mathbf{Y}}} = \Sigma_{\hat{\mathbf{Y}}}$ .

**Proof:** To prove part (a), let  $\mathbf{Y}^*$  be an arbitrary function of  $\mathbf{X}$  having mean  $\mathbf{0}$ , and let

$$\boldsymbol{\delta} = \mathbf{Y} - \mathbf{E}[\mathbf{Y} | \mathbf{X}], \quad \boldsymbol{\gamma} = \mathbf{E}[\mathbf{Y} | \mathbf{X}] - \mathbf{Y}^*, \quad \text{and} \quad \boldsymbol{\alpha} = \mathbf{Y} - \mathbf{Y}^* = \boldsymbol{\delta} + \boldsymbol{\gamma}.$$

We need to show that  $\mathbf{E}[\boldsymbol{\alpha}\boldsymbol{\alpha}^T] - \mathbf{E}[\boldsymbol{\delta}\boldsymbol{\delta}^T]$  is positive semidefinite. We have

$$\begin{aligned} \mathbf{E}[\boldsymbol{\alpha}\boldsymbol{\alpha}^T] - \mathbf{E}[\boldsymbol{\delta}\boldsymbol{\delta}^T] &= \mathbf{E}[(\boldsymbol{\delta} + \boldsymbol{\gamma})(\boldsymbol{\delta} + \boldsymbol{\gamma})^T] - \mathbf{E}[\boldsymbol{\delta}\boldsymbol{\delta}^T] \\ &= \mathbf{E}[\boldsymbol{\gamma}\boldsymbol{\delta}^T] + \mathbf{E}[\boldsymbol{\delta}\boldsymbol{\gamma}^T] + \mathbf{E}[\boldsymbol{\gamma}\boldsymbol{\gamma}^T], \end{aligned}$$

the third term of which is positive semidefinite since for any vector  $\mathbf{h}$  not identically zero, we have

$$\mathbf{h}^T \mathbf{E}[\boldsymbol{\gamma}\boldsymbol{\gamma}^T] \mathbf{h} = \mathbf{E}[\mathbf{h}^T \boldsymbol{\gamma}\boldsymbol{\gamma}^T \mathbf{h}] = \mathbf{E}[(\mathbf{h}^T \boldsymbol{\gamma})^2] \geq 0.$$

Now  $\mathbf{E}[\boldsymbol{\gamma}] = \mathbf{E}[\boldsymbol{\delta}] = \mathbf{0}$  and thus  $\mathbf{E}[\boldsymbol{\gamma}\boldsymbol{\delta}^T] = \text{Cov}(\boldsymbol{\gamma}, \boldsymbol{\delta})$  and  $\mathbf{E}[\boldsymbol{\delta}\boldsymbol{\gamma}^T] = \text{Cov}(\boldsymbol{\delta}, \boldsymbol{\gamma})$ . But both of these covariances are zero since  $\boldsymbol{\gamma} = \mathbf{E}[\mathbf{Y} | \mathbf{X}] - \mathbf{Y}^*$  is a function of  $\mathbf{X}$  since  $\mathbf{E}[\mathbf{Y} | \mathbf{X}]$  and  $\mathbf{Y}^*$  are both functions of  $\mathbf{X}$ .

To prove part (b), let  $\mathbf{A}$  be an arbitrary  $r \times d$  matrix. Then

$$\begin{aligned} S(\mathbf{A}) &= \mathbf{E}[(\mathbf{Y} - \mathbf{A}\mathbf{X})(\mathbf{Y} - \mathbf{A}\mathbf{X})^T] \\ &= \mathbf{E}[\mathbf{Y}\mathbf{Y}^T] - \mathbf{A}\mathbf{E}[\mathbf{X}\mathbf{Y}^T] - \mathbf{E}[\mathbf{Y}\mathbf{X}^T]\mathbf{A}^T + \mathbf{A}\mathbf{E}[\mathbf{X}\mathbf{X}^T]\mathbf{A}^T \\ &= \Sigma_{YY} - \mathbf{A}\Sigma_{XY} - \Sigma_{YX}\mathbf{A}^T + \mathbf{A}\Sigma_{XX}\mathbf{A}^T \\ &= (\mathbf{A} - \Sigma_{YX}\Sigma_{XX}^{-1})\Sigma_{XX}(\mathbf{A} - \Sigma_{YX}\Sigma_{XX}^{-1})^T \\ &\quad \Sigma_{YY} - \Sigma_{YX}\Sigma_{XX}^{-1}\Sigma_{XY}. \end{aligned}$$

Thus  $\Sigma_{\hat{\mathbf{Y}}} = S(\hat{\mathbf{A}}) = \Sigma_{YY} - \Sigma_{YX}\Sigma_{XX}^{-1}\Sigma_{XY}$  and

$$S(\mathbf{A}^*) - S(\hat{\mathbf{A}}) = (\mathbf{A}^* - \Sigma_{YX}\Sigma_{XX}^{-1})\Sigma_{XX}(\mathbf{A}^* - \Sigma_{YX}\Sigma_{XX}^{-1})^T$$

which is clearly positive definite for any matrix  $\mathbf{A}^*$ .

Finally, part (c) follows from part c of Theorem A.6.2.

## The Kalman Filter Algorithm

If we have a sequence of random vectors that are multivariate normal and satisfy a recursive relationship of a particular form, then predictors can also be found recursively. This type of recursion was used extensively by Kalman (1960) and thus the resulting algorithm is usually referred to as the Kalman Filter Algorithm.

**Theorem A.6.5 THE KALMAN FILTER.**

Let  $\mathbf{X}_0, \mathbf{X}_1, \dots$  and  $\mathbf{Z}_1, \mathbf{Z}_2, \dots$  be sequences of  $r$ - and  $d$ -dimensional random vectors satisfying the two equations (called the **state** and **observation equations** respectively) for  $t \geq 1$

$$\begin{aligned}\mathbf{X}_t &= \mathbf{A}_t \mathbf{X}_{t-1} + \mathbf{w}_t \\ \mathbf{Z}_t &= \mathbf{H}_t \mathbf{X}_t + \mathbf{v}_t,\end{aligned}$$

where the  $\mathbf{A}_t$  and  $\mathbf{H}_t$  are full rank  $r \times r$  and  $d \times r$  matrices of constants,

$$\begin{aligned}\mathbf{w}_t &\sim N_r(\mathbf{0}_r, \mathbf{W}_t) \\ \mathbf{v}_t &\sim N_d(\mathbf{0}_d, \mathbf{V}_t)\end{aligned}$$

all of the  $\mathbf{w}_t$  and  $\mathbf{v}_t$  are independent, and

$$\mathbf{X}_0 \sim N_r(\boldsymbol{\mu}_0, \boldsymbol{\Gamma}_0).$$

Let  $\hat{\mathbf{X}}_t$  and  $\hat{\boldsymbol{\Sigma}}_t$  denote the mean and variance of the conditional distribution of  $\mathbf{X}_t$  given  $\mathbf{Z}_1, \dots, \mathbf{Z}_t$ . Then

(a) For  $t \geq 1$

$$\begin{aligned}\hat{\mathbf{X}}_t &= \mathbf{A}_t \hat{\mathbf{X}}_{t-1} + \mathbf{R}_t \mathbf{H}_t^T (\mathbf{H}_t \mathbf{R}_t \mathbf{H}_t^T + \mathbf{V}_t)^{-1} \mathbf{e}_t \\ \hat{\boldsymbol{\Sigma}}_t &= \mathbf{R}_t - \mathbf{R}_t \mathbf{H}_t^T (\mathbf{H}_t \mathbf{R}_t \mathbf{H}_t^T + \mathbf{V}_t)^{-1} \mathbf{H}_t \mathbf{R}_t,\end{aligned}$$

where

$$\begin{aligned}\mathbf{R}_t &= \mathbf{A}_t \hat{\boldsymbol{\Sigma}}_{t-1} \mathbf{A}_t^T + \mathbf{W}_t \\ \mathbf{e}_t &= \mathbf{Z}_t - \mathbf{H}_t \mathbf{A}_t \hat{\mathbf{X}}_{t-1},\end{aligned}$$

and

$$\hat{\boldsymbol{\Sigma}}_0 = \boldsymbol{\Sigma}_0 \quad \text{and} \quad \hat{\mathbf{X}}_0 = \boldsymbol{\mu}_0.$$

(b) For  $t > 1$  and  $h \geq 0$ , let  $\tilde{\mathbf{X}}_{t+h-1|t-1}$  and  $\tilde{\boldsymbol{\Sigma}}_{t+h-1|t-1}$  denote the mean and variance of  $\mathbf{X}_{t+h-1}$  given  $\mathbf{Z}_1, \dots, \mathbf{Z}_{t-1}$ . Then

$$\begin{aligned}\tilde{\mathbf{X}}_{t+h-1|t-1} &= \begin{cases} \hat{\mathbf{X}}_{t-1}, & h = 0 \\ \mathbf{A}_{t+h-1} \tilde{\mathbf{X}}_{t+h-2|t-1}, & h \geq 1 \end{cases} \\ \tilde{\boldsymbol{\Sigma}}_{t+h-1|t-1} &= \begin{cases} \hat{\boldsymbol{\Sigma}}_{t-1}, & h = 0 \\ \mathbf{A}_{t+h-1} \tilde{\boldsymbol{\Sigma}}_{t+h-2|t-1} \mathbf{A}_{t+h-1}^T + \mathbf{W}_t, & h \geq 1. \end{cases}\end{aligned}$$

**Proof:** In our proof, we follow Meinhold & Singpurwalla (1983). Let  $\mathbf{Z}^t$  represent the set of variables  $\{\mathbf{Z}_1, \dots, \mathbf{Z}_t\}$ . We know that the distribution of  $\mathbf{X}_t | \mathbf{Z}^t$ , which is what we want, is normal since  $\mathbf{X}_0$  and all of the  $\mathbf{w}_t$ s and  $\mathbf{v}_t$ s are normal, and thus all of the  $\mathbf{X}_t$  and  $\mathbf{Z}_t$



are jointly normal. We want to show that the mean and variance of  $\mathbf{X}_t | \mathbf{Z}^t$  satisfy the relations in part (a). Note that  $\hat{\mathbf{X}}_{t-1} = \mathbb{E}[\mathbf{X}_{t-1} | \mathbf{Z}^{t-1}]$  is a linear function of  $\mathbf{Z}^{t-1}$  and thus  $\mathbf{e}_t = \mathbf{Z}_t - \mathbf{H}_t \mathbf{Z}_t \hat{\mathbf{X}}_{t-1}$  and  $\mathbf{Z}^{t-1}$  are a linear transformation of  $\mathbf{Z}^t$ . Thus, by part (d) of Theorem A.6.2, we know that the conditional distribution of  $\mathbf{X}_t | \mathbf{Z}^t$  is the same as that of  $\mathbf{X}_t | \mathbf{e}_t, \mathbf{Z}^{t-1}$ . The distribution of  $\mathbf{X}_t | \mathbf{e}_t, \mathbf{Z}^{t-1}$  can be found from that of  $\begin{pmatrix} \mathbf{X}_t \\ \mathbf{e}_t \end{pmatrix} | \mathbf{Z}^{t-1}$  by conditioning on  $\mathbf{e}_t$  (see part (ci) of Theorem A.6.2), and the distribution of  $\begin{pmatrix} \mathbf{X}_t \\ \mathbf{e}_t \end{pmatrix} | \mathbf{Z}^{t-1}$  can be found from  $\mathbf{X}_t | \mathbf{Z}^{t-1}$  and  $\mathbf{e}_t | \mathbf{X}_t, \mathbf{Z}^{t-1}$  by part (cii) of Theorem A.6.2. When  $t = 1$  we want the distribution of  $\mathbf{X}_1 | \mathbf{Z}_1$ , so we argue in the same way except there is no  $\mathbf{Z}^{t-1}$ .

Since  $\mathbf{X}_t = \mathbf{A}_t \mathbf{X}_{t-1} + \mathbf{w}_t$ , we have

$$\begin{aligned} \mathbb{E}[\mathbf{X}_t | \mathbf{Z}^{t-1}] &= \mathbf{A}_t \mathbb{E}[\mathbf{X}_{t-1} | \mathbf{Z}^{t-1}] + \mathbb{E}[\mathbf{w}_t | \mathbf{Z}^{t-1}] = \mathbf{A}_t \hat{\mathbf{X}}_{t-1} \\ \text{Var}(\mathbf{X}_t | \mathbf{Z}^{t-1}) &= \mathbf{A}_t \text{Var}(\mathbf{X}_{t-1} | \mathbf{Z}^{t-1}) \mathbf{A}_t^T + \text{Var}(\mathbf{w}_t | \mathbf{Z}^{t-1}) = \mathbf{A}_t \hat{\Sigma}_{t-1} \mathbf{A}_t^T + \mathbf{W}_t, \end{aligned}$$

that is

$$\mathbf{X}_t | \mathbf{Z}^{t-1} \sim N_r(\mathbf{A}_t \hat{\mathbf{X}}_{t-1}, \mathbf{A}_t \hat{\Sigma}_{t-1} \mathbf{A}_t^T + \mathbf{W}_t).$$

Since  $\mathbf{e}_t = \mathbf{Z}_t - \mathbf{H}_t \mathbf{A}_t \hat{\mathbf{X}}_{t-1}$  and  $\mathbf{Z}_t = \mathbf{H}_t \mathbf{X}_t + \mathbf{v}_t$ , we have  $\mathbf{e}_t = \mathbf{H}_t [\mathbf{X}_t - \mathbf{A}_t \hat{\mathbf{X}}_{t-1}] + \mathbf{v}_t$ , and since  $\hat{\mathbf{X}}_{t-1}$  is a function of  $\mathbf{Z}^{t-1}$ , we have

$$\mathbf{e}_t | \mathbf{X}_t, \mathbf{Z}^{t-1} \sim N_d(\mathbf{H}_t \{\mathbf{X}_t - \mathbf{A}_t \hat{\mathbf{X}}_{t-1}\}, \mathbf{V}_t).$$

Let  $\mathbf{R}_t = \mathbf{A}_t \hat{\Sigma}_{t-1} \mathbf{A}_t^T + \mathbf{W}_t$ , and to use part (cii) of Theorem A.6.2, make the identifications

$$\begin{aligned} \mathbf{X} &\leftrightarrow \mathbf{e}_t, & \mathbf{Y} &\leftrightarrow \mathbf{X}_t, & \boldsymbol{\mu}_Y &\leftrightarrow \mathbf{A}_t \hat{\mathbf{X}}_{t-1}, & \boldsymbol{\mu}_X &\leftrightarrow \mathbf{0}, \\ \Sigma_{XX} &\leftrightarrow \mathbf{H}_t \mathbf{R}_t \mathbf{H}_t^T + \mathbf{V}_t, & \Sigma_{XX} &\leftrightarrow \mathbf{H}_t \mathbf{R}_t, & \Sigma_{YY} &\leftrightarrow \mathbf{R}_t. \end{aligned}$$

To verify that this identification is valid, we suppress reference to  $\mathbf{Z}^{t-1}$  and note that

$$\begin{aligned} \mathbf{e}_t | \mathbf{X}_t &\leftrightarrow \mathbf{X} | \mathbf{Y} \\ &\sim N_d(\boldsymbol{\mu}_X + \Sigma_{XY} \Sigma_{YY}^{-1} \{\mathbf{Y} - \boldsymbol{\mu}_Y\}, \Sigma_{XX} - \Sigma_{XY} \Sigma_{YY}^{-1} \Sigma_{YX}) \\ &= N_d(\mathbf{H}_t \mathbf{R}_t \mathbf{R}_t^{-1} \{\mathbf{X}_t - \mathbf{A}_t \hat{\mathbf{X}}_{t-1}\}, \mathbf{H}_t \mathbf{R}_t \mathbf{H}_t^T + \mathbf{V}_t - \mathbf{H}_t \mathbf{R}_t \mathbf{R}_t^{-1} \mathbf{R}_t \mathbf{H}_t^T) \\ &= N_d(\mathbf{H}_t \{\mathbf{X}_t - \mathbf{A}_t \hat{\mathbf{X}}_{t-1}\}, \mathbf{V}_t) \end{aligned}$$

as required. Thus we have by part (cii) of Theorem A.6.2,

$$\begin{pmatrix} \mathbf{X}_t \\ \mathbf{e}_t \end{pmatrix} | \mathbf{Z}^{t-1} \sim N_{r+d} \left( \begin{bmatrix} \mathbf{0} \\ \mathbf{A}_t \hat{\mathbf{X}}_{t-1} \end{bmatrix}, \begin{bmatrix} \mathbf{H}_t \mathbf{R}_t \mathbf{H}_t^T + \mathbf{V}_t & \mathbf{H}_t \mathbf{R}_t \\ \mathbf{R}_t^T \mathbf{H}_t^T & \mathbf{R}_t \end{bmatrix} \right),$$

which in turn gives the desired recursion for  $\hat{\mathbf{X}}_t$  and  $\hat{\Sigma}_t$  by part (ci) of Theorem A.6.2.

Part (b) for  $h = 0$  follows by definition, while  $h > 0$  the results come from the fact that  $\mathbf{X}_{t+h-1} = \mathbf{A}_{t+h-2} \mathbf{X}_{t+h-2} + \mathbf{w}_{t+h-1}$ .

## Alternate Expression of the Kalman Filter Algorithm

There are many algebraically equivalent ways to express the Kalman Filter Algorithm (KFA). In Theorem A.6.5 we wrote it in a way making it clear how the recursion goes from one step to the next, and how it starts. To obtain a more intuitively appealing expressions, define

$$\tilde{\mathbf{X}}_t = \mathbf{E} [\mathbf{X}_t | \mathbf{Z}^{t-1}] \quad \text{and} \quad \tilde{\Sigma}_t = \text{Var} (\mathbf{X}_t | \mathbf{Z}^{t-1}),$$

that is, the mean and variance of  $\mathbf{X}_t$  having only observed  $\mathbf{Z}$  up to time  $t - 1$ . Then clearly

$$\tilde{\mathbf{X}}_t = \mathbf{A}_t \hat{\mathbf{X}}_{t-1} \quad \text{and} \quad \tilde{\Sigma}_t = \mathbf{R}_t;$$

that is, our best “estimate” of  $\mathbf{X}_t$  without having observed  $\mathbf{Z}_t$  is  $\mathbf{A}_t \hat{\mathbf{X}}_{t-1}$ . If we define the Kalman gain matrix  $\mathbf{K}_t$  by

$$\mathbf{K}_t = \tilde{\Sigma}_t \mathbf{H}_t^T (\mathbf{H}_t \tilde{\Sigma}_t \mathbf{H}_t^T + \mathbf{V}_t)^{-1},$$

we can write

$$\hat{\mathbf{X}}_t = \tilde{\mathbf{X}}_t + \mathbf{K}_t (\mathbf{Z}_t - \mathbf{H}_t \tilde{\mathbf{X}}_t);$$

that is, our best guess of  $\mathbf{X}_t$  given  $\mathbf{Z}^t$  is a linear combination of  $\tilde{\mathbf{X}}_t$  (our best guess of  $\mathbf{X}_t$  given  $\mathbf{Z}^{t-1}$ ) and a variable expressing how well  $\mathbf{H}_t \tilde{\mathbf{X}}_t$  predicts the new observation  $\mathbf{Z}_t$ . Note how similar this is to recursive regression (see Theorem A.5.1). We can summarize the alternative expression for the KFA as  $\hat{\mathbf{X}}_0 = \mathbf{E} [\mathbf{X}_0]$ ,  $\hat{\Sigma}_0 = \text{Var} (\mathbf{X}_0)$  and

$$\hat{\mathbf{X}}_t = \tilde{\mathbf{X}}_t + \mathbf{K}_t \mathbf{e}_t \quad \text{and} \quad \hat{\Sigma}_t = [\mathbf{I}_t - \mathbf{K}_t \mathbf{H}_t] \tilde{\Sigma}_t,$$

where

$$\begin{aligned} \tilde{\mathbf{X}}_t &= \mathbf{A}_t \hat{\mathbf{X}}_{t-1} \\ \mathbf{e}_t &= \mathbf{Z}_t - \mathbf{H}_t \tilde{\mathbf{X}}_t \\ \tilde{\Sigma}_t &= \mathbf{A}_t \hat{\Sigma}_{t-1} \mathbf{A}_t^T + \mathbf{W}_t \\ \mathbf{K}_t &= \tilde{\Sigma}_t \mathbf{H}_t^T (\mathbf{H}_t \tilde{\Sigma}_t \mathbf{H}_t^T + \mathbf{V}_t)^{-1}. \end{aligned}$$

### A.6.4 Convergence of Random Variables

Many of the theoretical results for statistics calculated in time series analysis are stated as limiting results as the length of the time series gets large. In this section we discuss the various modes of convergence of random variables.

**Definition A.6.4** TYPES OF CONVERGENCE. *Let  $\{X_n, n = 1, 2, \dots\}$  be a sequence of random variables such that the cdf of  $X_n$  is  $F_n$ . Let  $X$  be another random variable, and suppose that the cdf of  $X$  is  $F$ . Then*

(a)  $X_n$  converges to the constant  $c$  **in probability** if for every  $\epsilon > 0$ ,

$$\lim_{n \rightarrow \infty} \Pr(|X_n - c| > \epsilon) = 0.$$

We denote this by  $X_n \xrightarrow{\mathcal{P}} c$ , and say that  $X_n \xrightarrow{\mathcal{P}} X$  if  $X_n - X \xrightarrow{\mathcal{P}} 0$ .

(b)  $X_n$  converges to the constant  $c$  **with probability one** (or **almost surely**) if for every  $\epsilon > 0$ ,

$$\lim_{N \rightarrow \infty} \Pr\left(\sup_{n \geq N} |X_n - c| > \epsilon\right) = 0.$$

We denote this by  $X_n \xrightarrow{a.s.} c$ , and say that  $X_n \xrightarrow{a.s.} X$  if  $X_n - X \xrightarrow{a.s.} 0$ .

(c)  $X_n$  converges to  $X$  **in mean square** if

$$\lim_{n \rightarrow \infty} E[|X_n - X|^2] = 0.$$

We denote this by  $X_n \xrightarrow{m.s.} X$ .

(d)  $X_n$  converges to  $X$  **in distribution** if

$$\lim_{n \rightarrow \infty} F_n(x) = F(x)$$

for every  $x$  where  $F$  is continuous. We denote this by  $X_n \xrightarrow{\mathcal{L}} X$ .

(e) If  $\{Z_n, n = 0, \pm 1, \pm 2, \dots\}$  is a doubly infinite sequence of random variables and  $\beta_j, j = 0, \pm 1, \pm 2, \dots\}$  is a sequence of constants, we say that

$$Y_n = \sum_{j=-\infty}^{\infty} \beta_j Z_{n-j}$$

is a **limit in mean square** if

$$Y_{n,M} = \sum_{j=-M}^M \beta_j Z_{n-j} \xrightarrow{m.s.} Y_n.$$

We note that a statistic  $\hat{\theta}_n$ , calculated from a realization of length  $n$ , is said to be a weakly (strongly) consistent estimator of a parameter  $\theta$  if  $\hat{\theta}_n \xrightarrow{\mathcal{P}} \theta$  ( $\hat{\theta}_n \xrightarrow{a.s.} \theta$ ). Also, a random variable  $X_n$  is said to be asymptotically normally distributed with asymptotic mean  $\mu_n$  and asymptotic variance  $\sigma_n^2$  if

$$\frac{X_n - \mu_n}{\sigma_n} \xrightarrow{\mathcal{L}} X \sim N(0, 1).$$

We denote such a random variable by  $X_n \sim AN(\mu_n, \sigma_n^2)$ . For example, the usual Central Limit Theorem of statistics says that  $\bar{X}_n \sim AN(\mu, n^{-1}\sigma^2)$ .

## Convergence of Random Vectors

For a sequence of random vectors  $\{\mathbf{X}_n\}$ , we can define concepts of convergence analogous to those for a sequence of single random variables. Converges in probability, with probability one, and in mean square are defined pointwise; that is, we say that  $\mathbf{X}_n \xrightarrow{\mathcal{P}} \mathbf{X}$ ,  $\mathbf{X}_n \xrightarrow{a.s.} \mathbf{X}$ , or  $\mathbf{X}_n \xrightarrow{m.s.} \mathbf{X}$  if each element of  $\mathbf{X}_n$  converges to the corresponding element of  $\mathbf{X}$ . For a random vector to be asymptotically normal, we will write  $\mathbf{X}_n \sim AN_r(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$  if  $\boldsymbol{\Sigma}_n$  has nonzero diagonal elements for all  $n$  sufficiently large, and if for every vector  $\mathbf{h}$  such that  $\mathbf{h}^T \boldsymbol{\Sigma}_n \mathbf{h} > 0$  for all  $n$  sufficiently large, then

$$\mathbf{h}^T \mathbf{X}_n \sim AN(\mathbf{h}^T \boldsymbol{\mu}_n, \mathbf{h}^T \boldsymbol{\Sigma}_n \mathbf{h}).$$

Thus a random vector is said to be asymptotically normal if all linear functions of it are asymptotically normal in the scalar sense.

Several times in the book (see Exercise T4.8 for example), we make use of the fact that a continuous function of an asymptotically normal sequence of random vectors is also asymptotically normal. We state this formally in the following theorem (see Serfling (1980), p. 122, for example).

### Theorem A.6.6 CONTINUOUS FUNCTION THEOREM.

Suppose that  $\mathbf{X}_n \sim AN_r(\boldsymbol{\mu}, b_n \boldsymbol{\Sigma})$ , where  $\boldsymbol{\Sigma}$  is a covariance matrix and  $b_n \rightarrow 0$  as  $n \rightarrow \infty$ . Let

$$\mathbf{g}(\mathbf{x}) = (g_1(\mathbf{x}), \dots, g_p(\mathbf{x}))^T, \quad \mathbf{x} = (x_1, \dots, x_p)^T,$$

be a vector-valued function for which each component function  $g_i(\mathbf{x})$  is real-valued and totally differentiable. Let  $\mathbf{G}$  be the  $p \times r$  matrix whose  $(i, j)^{th}$  element is the partial derivative of  $g_i$  with respect to  $x_j$ , evaluated at the vector  $\boldsymbol{\mu}$ . Then

$$\mathbf{g}(\mathbf{X}_n) \sim AN(\mathbf{g}(\boldsymbol{\mu}), b_n \mathbf{G} \boldsymbol{\Sigma} \mathbf{G}^T).$$