

# Appendix A

## Matrix Algebra

### A.1 Notation

A **scalar**  $a$  is a single number.

A **vector**  $\mathbf{a}$  is a  $k \times 1$  list of numbers, typically arranged in a column. We write this as

$$\mathbf{a} = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_k \end{pmatrix}$$

Equivalently, a vector  $\mathbf{a}$  is an element of Euclidean  $k$  space, written as  $\mathbf{a} \in \mathbb{R}^k$ . If  $k = 1$  then  $\mathbf{a}$  is a scalar.

A **matrix**  $\mathbf{A}$  is a  $k \times r$  rectangular array of numbers, written as

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1r} \\ a_{21} & a_{22} & \cdots & a_{2r} \\ \vdots & \vdots & & \vdots \\ a_{k1} & a_{k2} & \cdots & a_{kr} \end{bmatrix}$$

By convention  $a_{ij}$  refers to the element in the  $i$ 'th row and  $j$ 'th column of  $\mathbf{A}$ . If  $r = 1$  then  $\mathbf{A}$  is a column vector. If  $k = 1$  then  $\mathbf{A}$  is a row vector. If  $r = k = 1$ , then  $\mathbf{A}$  is a scalar.

A standard convention (which we will follow in this text whenever possible) is to denote scalars by lower-case italics ( $a$ ), vectors by lower-case bold italics ( $\mathbf{a}$ ), and matrices by upper-case bold italics ( $\mathbf{A}$ ). Sometimes a matrix  $\mathbf{A}$  is denoted by the symbol  $(a_{ij})$ .

A matrix can be written as a set of column vectors or as a set of row vectors. That is,

$$\mathbf{A} = [ \mathbf{a}_1 \quad \mathbf{a}_2 \quad \cdots \quad \mathbf{a}_r ] = \begin{bmatrix} \boldsymbol{\alpha}_1 \\ \boldsymbol{\alpha}_2 \\ \vdots \\ \boldsymbol{\alpha}_k \end{bmatrix}$$

where

$$\mathbf{a}_i = \begin{bmatrix} a_{1i} \\ a_{2i} \\ \vdots \\ a_{ki} \end{bmatrix}$$

are column vectors and

$$\boldsymbol{\alpha}_j = [ a_{j1} \quad a_{j2} \quad \cdots \quad a_{jr} ]$$

are row vectors.

The **transpose** of a matrix  $\mathbf{A}$ , denoted  $\mathbf{A}'$ ,  $\mathbf{A}^\top$ , or  $\mathbf{A}^t$ , is obtained by flipping the matrix on its diagonal. Thus

$$\mathbf{A}' = \begin{bmatrix} a_{11} & a_{21} & \cdots & a_{k1} \\ a_{12} & a_{22} & \cdots & a_{k2} \\ \vdots & \vdots & & \vdots \\ a_{1r} & a_{2r} & \cdots & a_{kr} \end{bmatrix}$$

Alternatively, letting  $\mathbf{B} = \mathbf{A}'$ , then  $b_{ij} = a_{ji}$ . Note that if  $\mathbf{A}$  is  $k \times r$ , then  $\mathbf{A}'$  is  $r \times k$ . If  $\mathbf{a}$  is a  $k \times 1$  vector, then  $\mathbf{a}'$  is a  $1 \times k$  row vector.

A matrix is **square** if  $k = r$ . A square matrix is **symmetric** if  $\mathbf{A} = \mathbf{A}'$ , which requires  $a_{ij} = a_{ji}$ . A square matrix is **diagonal** if the off-diagonal elements are all zero, so that  $a_{ij} = 0$  if  $i \neq j$ . A square matrix is **upper (lower) diagonal** if all elements below (above) the diagonal equal zero.

An important diagonal matrix is the **identity matrix**, which has ones on the diagonal. The  $k \times k$  identity matrix is denoted as

$$\mathbf{I}_k = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix}.$$

A **partitioned matrix** takes the form

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} & \cdots & \mathbf{A}_{1r} \\ \mathbf{A}_{21} & \mathbf{A}_{22} & \cdots & \mathbf{A}_{2r} \\ \vdots & \vdots & & \vdots \\ \mathbf{A}_{k1} & \mathbf{A}_{k2} & \cdots & \mathbf{A}_{kr} \end{bmatrix}$$

where the  $A_{ij}$  denote matrices, vectors and/or scalars.

## A.2 Complex Matrices\*

Scalars, vectors and matrices may contain real or complex numbers as entries. (However, most econometric applications exclusively use real matrices.) If all elements of a vector  $\mathbf{x}$  are real we say that  $\mathbf{x}$  is a real vector, and similarly for matrices.

Recall that a complex number can be written as  $x = a + bi$  where where  $i = \sqrt{-1}$  and  $a$  and  $b$  are real numbers. Similarly a vector with complex elements can be written as

$$\mathbf{x} = \mathbf{a} + b\mathbf{i}$$

where  $\mathbf{a}$  and  $\mathbf{b}$  are real vectors, and a matrix with complex elements can be written as

$$\mathbf{X} = \mathbf{A} + B\mathbf{i}$$

where  $\mathbf{A}$  and  $\mathbf{B}$  are real matrices.

Recall that the complex conjugate of  $x = a + bi$  is  $x^* = a - bi$ . For matrices, the analogous concept is the conjugate transpose. The conjugate transpose of  $\mathbf{X} = \mathbf{A} + B\mathbf{i}$  is  $\mathbf{X}^* = \mathbf{A}' - B'\mathbf{i}$ . It is obtained by taking the transpose and taking the complex conjugate of each element.

## A.3 Matrix Addition

If the matrices  $\mathbf{A} = (a_{ij})$  and  $\mathbf{B} = (b_{ij})$  are of the same order, we define the sum

$$\mathbf{A} + \mathbf{B} = (a_{ij} + b_{ij}).$$

Matrix addition follows the commutative and associative laws:

$$\begin{aligned}\mathbf{A} + \mathbf{B} &= \mathbf{B} + \mathbf{A} \\ \mathbf{A} + (\mathbf{B} + \mathbf{C}) &= (\mathbf{A} + \mathbf{B}) + \mathbf{C}.\end{aligned}$$

## A.4 Matrix Multiplication

If  $\mathbf{A}$  is  $k \times r$  and  $c$  is real, we define their product as

$$\mathbf{A}c = c\mathbf{A} = (a_{ij}c).$$

If  $\mathbf{a}$  and  $\mathbf{b}$  are both  $k \times 1$ , then their inner product is

$$\mathbf{a}'\mathbf{b} = a_1b_1 + a_2b_2 + \cdots + a_kb_k = \sum_{j=1}^k a_jb_j.$$

Note that  $\mathbf{a}'\mathbf{b} = \mathbf{b}'\mathbf{a}$ . We say that two vectors  $\mathbf{a}$  and  $\mathbf{b}$  are **orthogonal** if  $\mathbf{a}'\mathbf{b} = 0$ .

If  $\mathbf{A}$  is  $k \times r$  and  $\mathbf{B}$  is  $r \times s$ , so that the number of columns of  $\mathbf{A}$  equals the number of rows of  $\mathbf{B}$ , we say that  $\mathbf{A}$  and  $\mathbf{B}$  are **conformable**. In this event the matrix product  $\mathbf{AB}$  is defined. Writing  $\mathbf{A}$  as a set of row vectors and  $\mathbf{B}$  as a set of column vectors (each of length  $r$ ), then the matrix product

is defined as

$$\begin{aligned} \mathbf{AB} &= \begin{bmatrix} \mathbf{a}'_1 \\ \mathbf{a}'_2 \\ \vdots \\ \mathbf{a}'_k \end{bmatrix} \begin{bmatrix} \mathbf{b}_1 & \mathbf{b}_2 & \cdots & \mathbf{b}_s \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{a}'_1 \mathbf{b}_1 & \mathbf{a}'_1 \mathbf{b}_2 & \cdots & \mathbf{a}'_1 \mathbf{b}_s \\ \mathbf{a}'_2 \mathbf{b}_1 & \mathbf{a}'_2 \mathbf{b}_2 & \cdots & \mathbf{a}'_2 \mathbf{b}_s \\ \vdots & \vdots & & \vdots \\ \mathbf{a}'_k \mathbf{b}_1 & \mathbf{a}'_k \mathbf{b}_2 & \cdots & \mathbf{a}'_k \mathbf{b}_s \end{bmatrix}. \end{aligned}$$

Matrix multiplication is not commutative: in general  $\mathbf{AB} \neq \mathbf{BA}$ . However, it is associative and distributive:

$$\begin{aligned} \mathbf{A}(\mathbf{BC}) &= (\mathbf{AB})\mathbf{C} \\ \mathbf{A}(\mathbf{B} + \mathbf{C}) &= \mathbf{AB} + \mathbf{AC} \end{aligned}$$

An alternative way to write the matrix product is to use matrix partitions. For example,

$$\begin{aligned} \mathbf{AB} &= \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{A}_{11}\mathbf{B}_{11} + \mathbf{A}_{12}\mathbf{B}_{21} & \mathbf{A}_{11}\mathbf{B}_{12} + \mathbf{A}_{12}\mathbf{B}_{22} \\ \mathbf{A}_{21}\mathbf{B}_{11} + \mathbf{A}_{22}\mathbf{B}_{21} & \mathbf{A}_{21}\mathbf{B}_{12} + \mathbf{A}_{22}\mathbf{B}_{22} \end{bmatrix}. \end{aligned}$$

As another example,

$$\begin{aligned} \mathbf{AB} &= \begin{bmatrix} \mathbf{A}_1 & \mathbf{A}_2 & \cdots & \mathbf{A}_r \end{bmatrix} \begin{bmatrix} \mathbf{B}_1 \\ \mathbf{B}_2 \\ \vdots \\ \mathbf{B}_r \end{bmatrix} \\ &= \mathbf{A}_1\mathbf{B}_1 + \mathbf{A}_2\mathbf{B}_2 + \cdots + \mathbf{A}_r\mathbf{B}_r \\ &= \sum_{j=1}^r \mathbf{A}_j\mathbf{B}_j \end{aligned}$$

An important property of the identity matrix is that if  $\mathbf{A}$  is  $k \times r$ , then  $\mathbf{A}\mathbf{I}_r = \mathbf{A}$  and  $\mathbf{I}_k\mathbf{A} = \mathbf{A}$ .

The  $k \times r$  matrix  $\mathbf{A}$ ,  $r \leq k$ , is called **orthonormal** if  $\mathbf{A}'\mathbf{A} = \mathbf{I}_r$ .

## A.5 Trace

The **trace** of a  $k \times k$  square matrix  $\mathbf{A}$  is the sum of its diagonal elements

$$\text{tr}(\mathbf{A}) = \sum_{i=1}^k a_{ii}.$$

Some straightforward properties for square matrices  $\mathbf{A}$  and  $\mathbf{B}$  and real  $c$  are

$$\begin{aligned}\text{tr}(c\mathbf{A}) &= c \text{tr}(\mathbf{A}) \\ \text{tr}(\mathbf{A}') &= \text{tr}(\mathbf{A}) \\ \text{tr}(\mathbf{A} + \mathbf{B}) &= \text{tr}(\mathbf{A}) + \text{tr}(\mathbf{B}) \\ \text{tr}(\mathbf{I}_k) &= k.\end{aligned}$$

Also, for  $k \times r$   $\mathbf{A}$  and  $r \times k$   $\mathbf{B}$  we have

$$\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA}). \tag{A.1}$$

Indeed,

$$\begin{aligned}\text{tr}(\mathbf{AB}) &= \text{tr} \begin{bmatrix} \mathbf{a}'_1 \mathbf{b}_1 & \mathbf{a}'_1 \mathbf{b}_2 & \cdots & \mathbf{a}'_1 \mathbf{b}_k \\ \mathbf{a}'_2 \mathbf{b}_1 & \mathbf{a}'_2 \mathbf{b}_2 & \cdots & \mathbf{a}'_2 \mathbf{b}_k \\ \vdots & \vdots & & \vdots \\ \mathbf{a}'_k \mathbf{b}_1 & \mathbf{a}'_k \mathbf{b}_2 & \cdots & \mathbf{a}'_k \mathbf{b}_k \end{bmatrix} \\ &= \sum_{i=1}^k \mathbf{a}'_i \mathbf{b}_i \\ &= \sum_{i=1}^k \mathbf{b}'_i \mathbf{a}_i \\ &= \text{tr}(\mathbf{BA}).\end{aligned}$$

## A.6 Rank and Inverse

The rank of the  $k \times r$  matrix ( $r \leq k$ )

$$\mathbf{A} = [ \mathbf{a}_1 \quad \mathbf{a}_2 \quad \cdots \quad \mathbf{a}_r ]$$

is the number of linearly independent columns  $\mathbf{a}_j$ , and is written as  $\text{rank}(\mathbf{A})$ . We say that  $\mathbf{A}$  has full rank if  $\text{rank}(\mathbf{A}) = r$ .

A square  $k \times k$  matrix  $\mathbf{A}$  is said to be **nonsingular** if it has full rank, e.g.  $\text{rank}(\mathbf{A}) = k$ . This means that there is no  $k \times 1$   $\mathbf{c} \neq \mathbf{0}$  such that  $\mathbf{A}\mathbf{c} = \mathbf{0}$ .

If a square  $k \times k$  matrix  $\mathbf{A}$  is nonsingular then there exists a unique matrix  $k \times k$  matrix  $\mathbf{A}^{-1}$  called the **inverse** of  $\mathbf{A}$  which satisfies

$$\mathbf{A}\mathbf{A}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}_k.$$

For non-singular  $\mathbf{A}$  and  $\mathbf{C}$ , some important properties include

$$\begin{aligned} \mathbf{A}\mathbf{A}^{-1} &= \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}_k \\ (\mathbf{A}^{-1})' &= (\mathbf{A}')^{-1} \\ (\mathbf{A}\mathbf{C})^{-1} &= \mathbf{C}^{-1}\mathbf{A}^{-1} \\ (\mathbf{A} + \mathbf{C})^{-1} &= \mathbf{A}^{-1}(\mathbf{A}^{-1} + \mathbf{C}^{-1})^{-1}\mathbf{C}^{-1} \\ \mathbf{A}^{-1} - (\mathbf{A} + \mathbf{C})^{-1} &= \mathbf{A}^{-1}(\mathbf{A}^{-1} + \mathbf{C}^{-1})^{-1}\mathbf{A}^{-1} \end{aligned}$$

Also, if  $\mathbf{A}$  is an orthonormal matrix, then  $\mathbf{A}^{-1} = \mathbf{A}'$ .

Another useful result for non-singular  $\mathbf{A}$  is known as the **Woodbury matrix identity**

$$(\mathbf{A} + \mathbf{BCD})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{BC}(\mathbf{C} + \mathbf{CDA}^{-1}\mathbf{BC})^{-1}\mathbf{CDA}^{-1}. \quad (\text{A.2})$$

In particular, for  $\mathbf{C} = -1$ ,  $\mathbf{B} = \mathbf{b}$  and  $\mathbf{D} = \mathbf{b}'$  for vector  $\mathbf{b}$  we find what is known as the **Sherman–Morrison formula**

$$(\mathbf{A} - \mathbf{bb}')^{-1} = \mathbf{A}^{-1} + (1 - \mathbf{b}'\mathbf{A}^{-1}\mathbf{b})^{-1}\mathbf{A}^{-1}\mathbf{bb}'\mathbf{A}^{-1}. \quad (\text{A.3})$$

The following fact about inverting partitioned matrices is quite useful.

$$\begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}^{-1} \stackrel{def}{=} \begin{bmatrix} \mathbf{A}^{11} & \mathbf{A}^{12} \\ \mathbf{A}^{21} & \mathbf{A}^{22} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_{11.2}^{-1} & -\mathbf{A}_{11.2}^{-1}\mathbf{A}_{12}\mathbf{A}_{22}^{-1} \\ -\mathbf{A}_{22.1}^{-1}\mathbf{A}_{21}\mathbf{A}_{11}^{-1} & \mathbf{A}_{22.1}^{-1} \end{bmatrix} \quad (\text{A.4})$$

where  $\mathbf{A}_{11.2} = \mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21}$  and  $\mathbf{A}_{22.1} = \mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12}$ . There are alternative algebraic representations for the components. For example, using the Woodbury matrix identity you can show the following alternative expressions

$$\begin{aligned} \mathbf{A}^{11} &= \mathbf{A}_{11}^{-1} + \mathbf{A}_{11}^{-1}\mathbf{A}_{12}\mathbf{A}_{22.1}^{-1}\mathbf{A}_{21}\mathbf{A}_{11}^{-1} \\ \mathbf{A}^{22} &= \mathbf{A}_{22}^{-1} + \mathbf{A}_{22}^{-1}\mathbf{A}_{21}\mathbf{A}_{11.2}^{-1}\mathbf{A}_{12}\mathbf{A}_{22}^{-1} \\ \mathbf{A}^{12} &= -\mathbf{A}_{11}^{-1}\mathbf{A}_{12}\mathbf{A}_{22.1}^{-1} \\ \mathbf{A}^{21} &= -\mathbf{A}_{22}^{-1}\mathbf{A}_{21}\mathbf{A}_{11.2}^{-1} \end{aligned}$$

Even if a matrix  $\mathbf{A}$  does not possess an inverse, we can still define the **Moore-Penrose generalized inverse**  $\mathbf{A}^-$  as the matrix which satisfies

$$\begin{aligned} \mathbf{A}\mathbf{A}^- \mathbf{A} &= \mathbf{A} \\ \mathbf{A}^- \mathbf{A}\mathbf{A}^- &= \mathbf{A}^- \\ \mathbf{A}\mathbf{A}^- &\text{ is symmetric} \\ \mathbf{A}^- \mathbf{A} &\text{ is symmetric} \end{aligned}$$

For any matrix  $\mathbf{A}$ , the Moore-Penrose generalized inverse  $\mathbf{A}^-$  exists and is unique.

For example, if

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$$

and when  $\mathbf{A}_{11}^{-1}$  exists then

$$\mathbf{A}^- = \begin{bmatrix} \mathbf{A}_{11}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}.$$



## A.7 Determinant

The **determinant** is a measure of the volume of a square matrix. It is written as  $\det \mathbf{A}$  or  $|\mathbf{A}|$ .

While the determinant is widely used, its precise definition is rarely needed. However, we present the definition here for completeness. Let  $\mathbf{A} = (a_{ij})$  be a  $k \times k$  matrix. Let  $\pi = (j_1, \dots, j_k)$  denote a permutation of  $(1, \dots, k)$ . There are  $k!$  such permutations. There is a unique count of the number of inversions of the indices of such permutations (relative to the natural order  $(1, \dots, k)$ ), and let  $\varepsilon_\pi = +1$  if this count is even and  $\varepsilon_\pi = -1$  if the count is odd. Then the determinant of  $\mathbf{A}$  is defined as

$$\det \mathbf{A} = \sum_{\pi} \varepsilon_{\pi} a_{1j_1} a_{2j_2} \cdots a_{kj_k}.$$

For example, if  $\mathbf{A}$  is  $2 \times 2$ , then the two permutations of  $(1, 2)$  are  $(1, 2)$  and  $(2, 1)$ , for which  $\varepsilon_{(1,2)} = 1$  and  $\varepsilon_{(2,1)} = -1$ . Thus

$$\begin{aligned} \det \mathbf{A} &= \varepsilon_{(1,2)} a_{11} a_{22} + \varepsilon_{(2,1)} a_{21} a_{12} \\ &= a_{11} a_{22} - a_{12} a_{21}. \end{aligned}$$

For a square matrix  $\mathbf{A}$ , the **minor**  $M_{ij}$  of the  $ij^{\text{th}}$  element  $a_{ij}$  is the determinant of the matrix obtained by removing the  $i^{\text{th}}$  row and  $j^{\text{th}}$  column of  $\mathbf{A}$ . The **cofactor** of the  $ij^{\text{th}}$  element is  $C_{ij} = (-1)^{i+j} M_{ij}$ . An important representation known as Laplace's expansion, relates the determinant of  $\mathbf{A}$  to its cofactors:

$$\det \mathbf{A} = \sum_{j=1}^k a_{ij} C_{ij}$$

This holds for all  $i = 1, \dots, k$ . This is often presented as a method for computation of a determinant.

Some properties of the determinant include

- $\det(\mathbf{A}) = \det(\mathbf{A}')$
- $\det(c\mathbf{A}) = c^k \det \mathbf{A}$
- $\det(\mathbf{A}\mathbf{B}) = (\det \mathbf{A})(\det \mathbf{B})$

- $\det(\mathbf{A}^{-1}) = (\det \mathbf{A})^{-1}$
- $\det \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = (\det \mathbf{D}) \det(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})$  if  $\det \mathbf{D} \neq 0$
- $\det \mathbf{A} \neq 0$  if and only if  $\mathbf{A}$  is nonsingular
- If  $\mathbf{A}$  is triangular (upper or lower), then  $\det \mathbf{A} = \prod_{i=1}^k a_{ii}$
- If  $\mathbf{A}$  is orthogonal, then  $\det \mathbf{A} = \pm 1$
- $\mathbf{A}^{-1} = (\det \mathbf{A})^{-1} \mathbf{C}$  where  $\mathbf{C} = (C_{ij})$  is the matrix of cofactors

## A.8 Eigenvalues

The characteristic equation of a  $k \times k$  square matrix  $\mathbf{A}$  is

$$\det(\mathbf{A} - \lambda \mathbf{I}_k) = 0.$$

The left side is a polynomial of degree  $k$  in  $\lambda$  so it has exactly  $k$  roots, which are not necessarily distinct and may be real or complex. They are called the **latent roots** or **characteristic roots** or **eigenvalues** of  $\mathbf{A}$ . If  $\lambda_i$  is an eigenvalue of  $\mathbf{A}$ , then  $\mathbf{A} - \lambda_i \mathbf{I}_k$  is singular so there exists a non-zero vector  $\mathbf{h}_i$  such that

$$(\mathbf{A} - \lambda_i \mathbf{I}_k) \mathbf{h}_i = \mathbf{0}.$$

The vector  $\mathbf{h}_i$  is called a **latent vector** or **characteristic vector** or **eigen-vector** of  $\mathbf{A}$  corresponding to  $\lambda_i$ .

We now state some useful properties. Let  $\lambda_i$  and  $\mathbf{h}_i$ ,  $i = 1, \dots, k$  denote the  $k$  eigenvalues and eigenvectors of a square matrix  $\mathbf{A}$ . Let  $\mathbf{\Lambda}$  be a diagonal matrix with the characteristic roots in the diagonal, and let  $\mathbf{H} = [\mathbf{h}_1 \cdots \mathbf{h}_k]$ .

- $\det(\mathbf{A}) = \prod_{i=1}^k \lambda_i$
- $\text{tr}(\mathbf{A}) = \sum_{i=1}^k \lambda_i$
- $\mathbf{A}$  is non-singular if and only if all its characteristic roots are non-zero.
- If  $\mathbf{A}$  has distinct characteristic roots, there exists a nonsingular matrix  $\mathbf{P}$  such that  $\mathbf{A} = \mathbf{P}^{-1} \mathbf{\Lambda} \mathbf{P}$  and  $\mathbf{P} \mathbf{A} \mathbf{P}^{-1} = \mathbf{\Lambda}$ .

- If  $k \times k$   $\mathbf{A}$  is symmetric, then  $\mathbf{A} = \mathbf{H}\mathbf{\Lambda}\mathbf{H}'$  and  $\mathbf{H}'\mathbf{A}\mathbf{H} = \mathbf{\Lambda}$ , where  $\mathbf{\Lambda}$  is a diagonal matrix with the eigenvalues on the diagonal, and  $\mathbf{H}'\mathbf{A}\mathbf{H} = \mathbf{I}_k$ . The characteristic roots are all real.  $\mathbf{A} = \mathbf{H}\mathbf{\Lambda}\mathbf{H}'$  is called the **spectral decomposition** of  $\mathbf{A}$ .
- When the eigenvalues of  $k \times k$   $\mathbf{A}$  are real it is conventional to write them in descending order  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k$ . We also write  $\lambda_{\min}(\mathbf{A}) = \lambda_k = \min\{\lambda_\ell\}$  and  $\lambda_{\max}(\mathbf{A}) = \lambda_1 = \max\{\lambda_\ell\}$ .
- For real symmetric  $\mathbf{A}$ ,  $\lambda_{\max}(\mathbf{A}) = \max_{\mathbf{x}'\mathbf{x}=1} \mathbf{x}'\mathbf{A}\mathbf{x}$
- For real symmetric  $\mathbf{A}$ ,  $\lambda_{\min}(\mathbf{A}) = \min_{\mathbf{x}'\mathbf{x}=1} \mathbf{x}'\mathbf{A}\mathbf{x}$
- The characteristic roots of  $\mathbf{A}^{-1}$  are  $\lambda_1^{-1}, \lambda_2^{-1}, \dots, \lambda_k^{-1}$ .
- The matrix  $\mathbf{H}$  has the **orthonormal** properties  $\mathbf{H}'\mathbf{H} = \mathbf{I}$  and  $\mathbf{H}\mathbf{H}' = \mathbf{I}$ .
- $\mathbf{H}^{-1} = \mathbf{H}'$  and  $(\mathbf{H}')^{-1} = \mathbf{H}$
- For any  $k \times 1$  vector  $\mathbf{a}$ ,  $\lambda_{\max}(\mathbf{a}\mathbf{a}') = \mathbf{a}'\mathbf{a}$

## A.9 Positive Definiteness

We say that a  $k \times k$  real symmetric square matrix  $\mathbf{A}$  is **positive semi-definite** if for all  $\mathbf{c} \neq \mathbf{0}$ ,  $\mathbf{c}'\mathbf{A}\mathbf{c} \geq 0$ . This is written as  $\mathbf{A} \geq 0$ . We say that  $\mathbf{A}$  is **positive definite** if for all  $\mathbf{c} \neq \mathbf{0}$ ,  $\mathbf{c}'\mathbf{A}\mathbf{c} > 0$ . This is written as  $\mathbf{A} > 0$ .

Some properties include:

- If  $\mathbf{A} = \mathbf{G}'\mathbf{B}\mathbf{G}$  with  $\mathbf{B} \geq 0$  and some matrix  $\mathbf{G}$ , then  $\mathbf{A}$  is positive semi-definite. (For any  $\mathbf{c} \neq \mathbf{0}$ ,  $\mathbf{c}'\mathbf{A}\mathbf{c} = \mathbf{\alpha}'\mathbf{B}\mathbf{\alpha} \geq 0$  where  $\mathbf{\alpha} = \mathbf{G}\mathbf{c}$ .) If  $\mathbf{G}$  has full column rank and  $\mathbf{B} > 0$ , then  $\mathbf{A}$  is positive definite.
- If  $\mathbf{A}$  is positive definite, then  $\mathbf{A}$  is non-singular and  $\mathbf{A}^{-1}$  exists. Furthermore,  $\mathbf{A}^{-1} > 0$ .
- $\mathbf{A} > 0$  if and only if it is symmetric and all its characteristic roots are positive.

- By the spectral decomposition,  $\mathbf{A} = \mathbf{H}\mathbf{\Lambda}\mathbf{H}'$  where  $\mathbf{H}'\mathbf{H} = \mathbf{I}$  and  $\mathbf{\Lambda}$  is diagonal with non-negative diagonal elements. All diagonal elements of  $\mathbf{\Lambda}$  are strictly positive if (and only if)  $\mathbf{A} > 0$ .
- The rank of  $\mathbf{A}$  equals the number of strictly positive characteristic roots.
- If  $\mathbf{A} > 0$  then  $\mathbf{A}^{-1} = \mathbf{H}\mathbf{\Lambda}^{-1}\mathbf{H}'$ .
- If  $\mathbf{A} \geq 0$  and  $\text{rank}(\mathbf{A}) = r < k$  then  $\mathbf{A}^- = \mathbf{H}\mathbf{\Lambda}^-\mathbf{H}'$  where  $\mathbf{A}^-$  is the Moore-Penrose generalized inverse, and  $\mathbf{\Lambda}^- = \text{diag}(\lambda_1^{-1}, \lambda_2^{-1}, \dots, \lambda_k^{-1}, 0, \dots, 0)$
- If  $\mathbf{A} \geq 0$  we can find a matrix  $\mathbf{B}$  such that  $\mathbf{A} = \mathbf{B}\mathbf{B}'$ . We call  $\mathbf{B}$  a **matrix square root** of  $\mathbf{A}$ . The matrix  $\mathbf{B}$  need not be unique. One way to construct  $\mathbf{B}$  is to use the spectral decomposition  $\mathbf{A} = \mathbf{H}\mathbf{\Lambda}\mathbf{H}'$  where  $\mathbf{\Lambda}$  is diagonal, and then set  $\mathbf{B} = \mathbf{H}\mathbf{\Lambda}^{1/2}$ . There is a unique square root  $\mathbf{B}$  which is also positive semi-definite  $\mathbf{B} \geq 0$ . If  $\mathbf{A} > 0$  then  $\mathbf{B} > 0$ .

A  $k \times k$  square matrix  $\mathbf{A}$  is **idempotent** if  $\mathbf{A}\mathbf{A} = \mathbf{A}$ . If  $\mathbf{A}$  is idempotent and symmetric with rank  $r$ , then it has  $r$  characteristic roots which equal 1 and  $k - r$  characteristic roots which equal 0. To see this, by the spectral decomposition that we can write  $\mathbf{A} = \mathbf{H}\mathbf{\Lambda}\mathbf{H}'$  where  $\mathbf{H}$  is orthogonal and  $\mathbf{\Lambda}$  contains the eigenvalues. Then

$$\mathbf{A} = \mathbf{A}\mathbf{A} = \mathbf{H}\mathbf{\Lambda}\mathbf{H}'\mathbf{H}\mathbf{\Lambda}\mathbf{H}' = \mathbf{H}\mathbf{\Lambda}^2\mathbf{H}'.$$

We deduce that  $\mathbf{\Lambda}^2 = \mathbf{\Lambda}$  and  $\lambda_i^2 = \lambda_i$  for  $i = 1, \dots, k$ . Hence the  $\lambda_i$  must equal either 0 or 1. Since the rank of  $\mathbf{A}$  is  $r$ , and the rank equals the number of positive characteristic roots, it follows that

$$\mathbf{\Lambda} = \begin{bmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0}_{k-r} \end{bmatrix}$$

and the spectral decomposition of idempotent  $\mathbf{A}$  takes the form

$$\mathbf{A} = \mathbf{H} \begin{bmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0}_{k-r} \end{bmatrix} \mathbf{H}' \tag{A.5}$$

with  $\mathbf{H}'\mathbf{H} = \mathbf{I}_k$ . Additionally,  $\text{tr}(\mathbf{A}) = \text{rank}(\mathbf{A})$  and  $\mathbf{A}$  is positive semi-definite.

If  $\mathbf{A}$  is **idempotent** then  $\mathbf{I} - \mathbf{A}$  is also idempotent.

One useful fact is that if  $\mathbf{A}$  is idempotent then for any conformable vector  $\mathbf{c}$ ,

$$\mathbf{c}'\mathbf{A}\mathbf{c} \leq \mathbf{c}'\mathbf{c} \tag{A.6}$$

$$\mathbf{c}'(\mathbf{I} - \mathbf{A})\mathbf{c} \leq \mathbf{c}'\mathbf{c} \tag{A.7}$$

To see this, note that

$$\mathbf{c}'\mathbf{c} = \mathbf{c}'\mathbf{A}\mathbf{c} + \mathbf{c}'(\mathbf{I} - \mathbf{A})\mathbf{c}.$$

Since  $\mathbf{A}$  and  $\mathbf{I} - \mathbf{A}$  are idempotent, they are both positive semi-definite, so both  $\mathbf{c}'\mathbf{A}\mathbf{c}$  and  $\mathbf{c}'(\mathbf{I} - \mathbf{A})\mathbf{c}$  are non-negative. Thus they must satisfy (A.6)-(A.7).

## A.10 Singular Values

The singular values of a  $k \times r$  real matrix  $\mathbf{A}$  are the square roots of the eigenvalues of  $\mathbf{A}'\mathbf{A}$ . Thus for  $j = 1, \dots, r$

$$s_j = \sqrt{\lambda_j(\mathbf{A}'\mathbf{A})}$$

Since  $\mathbf{A}'\mathbf{A}$  is positive semi-definite, its eigenvalues are non-negative so the singular values are real and non-negative.

The non-zero singular values of  $\mathbf{A}$  and  $\mathbf{A}'$  are the same.

When  $\mathbf{A}$  is positive semi-definite then the singular values of  $\mathbf{A}$  correspond to its eigenvalues.

The singular value decomposition of a  $k \times r$  real matrix  $\mathbf{A}$  takes the form  $\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{V}'$  where  $\mathbf{U}$  is  $k \times k$ ,  $\mathbf{\Lambda}$  is  $k \times r$  and  $\mathbf{V}$  is  $r \times r$ , with  $\mathbf{U}$  and  $\mathbf{V}$  orthonormal ( $\mathbf{U}'\mathbf{U} = \mathbf{I}_k$  and  $\mathbf{V}'\mathbf{V} = \mathbf{I}_r$ ) and  $\mathbf{\Lambda}$  is a diagonal matrix with the singular values of  $\mathbf{A}$  on the diagonal.

It is convention to write the singular values in descending order  $s_1 \geq s_2 \geq \dots \geq s_r$ .

## A.11 Matrix Calculus

Let  $\mathbf{x} = (x_1, \dots, x_k)$  be  $k \times 1$  and  $g(\mathbf{x}) = g(x_1, \dots, x_k) : \mathbb{R}^k \rightarrow \mathbb{R}$ . The vector derivative is

$$\frac{\partial}{\partial \mathbf{x}} g(\mathbf{x}) = \begin{pmatrix} \frac{\partial}{\partial x_1} g(\mathbf{x}) \\ \vdots \\ \frac{\partial}{\partial x_k} g(\mathbf{x}) \end{pmatrix}$$

and

$$\frac{\partial}{\partial \mathbf{x}'} g(\mathbf{x}) = \left( \frac{\partial}{\partial x_1} g(\mathbf{x}) \quad \cdots \quad \frac{\partial}{\partial x_k} g(\mathbf{x}) \right).$$

Some properties are now summarized.

- $\frac{\partial}{\partial \mathbf{x}} (\mathbf{a}'\mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} (\mathbf{x}'\mathbf{a}) = \mathbf{a}$
- $\frac{\partial}{\partial \mathbf{x}'} (\mathbf{A}\mathbf{x}) = \mathbf{A}$
- $\frac{\partial}{\partial \mathbf{x}} (\mathbf{x}'\mathbf{A}\mathbf{x}) = (\mathbf{A} + \mathbf{A}')\mathbf{x}$
- $\frac{\partial^2}{\partial \mathbf{x} \partial \mathbf{x}'} (\mathbf{x}'\mathbf{A}\mathbf{x}) = \mathbf{A} + \mathbf{A}'$
- $\frac{\partial}{\partial \mathbf{A}} \text{tr}(\mathbf{B}\mathbf{A}) = \mathbf{B}'$
- $\frac{\partial}{\partial \mathbf{A}} \log \det(\mathbf{A}) = (\mathbf{A}^{-1})'$

The final two results require some justification. Recall from Section A.5 that we can write out explicitly

$$\text{tr}(\mathbf{B}\mathbf{A}) = \sum_i \sum_j a_{ij} b_{ji}.$$

Thus if we take the derivative with respect to  $a_{ij}$  we find

$$\frac{\partial}{\partial a_{ij}} \text{tr}(\mathbf{B}\mathbf{A}) = b_{ji}.$$

which is the  $ij^{\text{th}}$  element of  $\mathbf{B}'$ , establishing the second-to-last result.

For the final result, recall Laplace's expansion

$$\det \mathbf{A} = \sum_{j=1}^k a_{ij} C_{ij}$$

where  $C_{ij}$  is the  $ij^{\text{th}}$  cofactor of  $\mathbf{A}$ . Set  $\mathbf{C} = (C_{ij})$ . Observe that  $C_{ij}$  for  $j = 1, \dots, k$  are not functions of  $a_{ij}$ . Thus the derivative with respect to  $a_{ij}$  is

$$\frac{\partial}{\partial a_{ij}} \log \det (\mathbf{A}) = (\det \mathbf{A})^{-1} \frac{\partial}{\partial a_{ij}} \det \mathbf{A} = (\det \mathbf{A})^{-1} C_{ij}$$

Together this implies

$$\frac{\partial}{\partial \mathbf{A}} \log \det (\mathbf{A}) = (\det \mathbf{A})^{-1} \mathbf{C} = \mathbf{A}^{-1}$$

where the second equality is a property of the inverse from Section A.7.

## A.12 Kronecker Products and the Vec Operator

Let  $\mathbf{A} = [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_n]$  be  $m \times n$ . The **vec** of  $\mathbf{A}$ , denoted by  $\text{vec}(\mathbf{A})$ , is the  $mn \times 1$  vector

$$\text{vec}(\mathbf{A}) = \begin{pmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \\ \vdots \\ \mathbf{a}_n \end{pmatrix}.$$

Let  $\mathbf{A} = (a_{ij})$  be an  $m \times n$  matrix and let  $\mathbf{B}$  be any matrix. The **Kronecker product** of  $\mathbf{A}$  and  $\mathbf{B}$ , denoted  $\mathbf{A} \otimes \mathbf{B}$ , is the matrix

$$\mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} & \cdots & a_{1n}\mathbf{B} \\ a_{21}\mathbf{B} & a_{22}\mathbf{B} & \cdots & a_{2n}\mathbf{B} \\ \vdots & \vdots & \cdots & \vdots \\ a_{m1}\mathbf{B} & a_{m2}\mathbf{B} & \cdots & a_{mn}\mathbf{B} \end{bmatrix}.$$

Some important properties are now summarized. These results hold for matrices for which all matrix multiplications are conformable.

- $(\mathbf{A} + \mathbf{B}) \otimes \mathbf{C} = \mathbf{A} \otimes \mathbf{C} + \mathbf{B} \otimes \mathbf{C}$

- $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = \mathbf{AC} \otimes \mathbf{BD}$
- $\mathbf{A} \otimes (\mathbf{B} \otimes \mathbf{C}) = (\mathbf{A} \otimes \mathbf{B}) \otimes \mathbf{C}$
- $(\mathbf{A} \otimes \mathbf{B})' = \mathbf{A}' \otimes \mathbf{B}'$
- $\text{tr}(\mathbf{A} \otimes \mathbf{B}) = \text{tr}(\mathbf{A}) \text{tr}(\mathbf{B})$
- If  $\mathbf{A}$  is  $m \times m$  and  $\mathbf{B}$  is  $n \times n$ ,  $\det(\mathbf{A} \otimes \mathbf{B}) = (\det(\mathbf{A}))^n (\det(\mathbf{B}))^m$
- $(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}$
- If  $\mathbf{A} > 0$  and  $\mathbf{B} > 0$  then  $\mathbf{A} \otimes \mathbf{B} > 0$
- $\text{vec}(\mathbf{ABC}) = (\mathbf{C}' \otimes \mathbf{A}) \text{vec}(\mathbf{B})$
- $\text{tr}(\mathbf{ABCD}) = \text{vec}(\mathbf{D}')' (\mathbf{C}' \otimes \mathbf{A}) \text{vec}(\mathbf{B})$

## A.13 Vector Norms

Given any vector space  $V$  (such as Euclidean space  $\mathbb{R}^m$ ) a **norm** on  $V$  is a function  $\rho : V \rightarrow \mathbb{R}$  with the properties

1.  $\rho(c\mathbf{a}) = |c| \rho(\mathbf{a})$  for any complex number  $c$  and  $\mathbf{a} \in V$
2.  $\rho(\mathbf{a} + \mathbf{b}) \leq \rho(\mathbf{a}) + \rho(\mathbf{b})$
3. If  $\rho(\mathbf{a}) = 0$  then  $\mathbf{a} = \mathbf{0}$

A seminorm on  $V$  is a function which satisfies the first two properties. The second property is known as the triangle inequality, and it is the one property which typically needs a careful demonstration (as the other two properties typically hold by inspection).

The typical norm used for Euclidean space  $\mathbb{R}^m$  is the **Euclidean norm**

$$\begin{aligned} \|\mathbf{a}\| &= (\mathbf{a}'\mathbf{a})^{1/2} \\ &= \left( \sum_{i=1}^m a_i^2 \right)^{1/2}. \end{aligned}$$



Alternative norms include the  $p$ -norm (for  $p \geq 1$ )

$$\|\mathbf{a}\|_p = \left( \sum_{i=1}^m |a_i|^p \right)^{1/p}$$

Special cases include the Euclidean norm ( $p = 2$ ), the 1-norm

$$\|\mathbf{a}\|_1 = \sum_{i=1}^m |a_i|$$

and the sup-norm

$$\|\mathbf{a}\|_\infty = \max(|a_1|, \dots, |a_m|).$$

For real numbers ( $m = 1$ ) these norms coincide.

Some standard inequalities for Euclidean space are now given. The Minkowski inequality given below establishes that any  $p$ -norm with  $p \geq 1$  (including the Euclidean norm) satisfies the triangle inequality and is thus a valid norm.

**Jensen's Inequality.** If  $g(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$  is convex, then for any non-negative weights  $a_j$  such that  $\sum_{j=1}^m a_j = 1$ , and any real numbers  $x_j$

$$g\left(\sum_{j=1}^m a_j x_j\right) \leq \sum_{j=1}^m a_j g(x_j). \quad (\text{A.8})$$

In particular, setting  $a_j = 1/m$ , then

$$g\left(\frac{1}{m} \sum_{j=1}^m x_j\right) \leq \frac{1}{m} \sum_{j=1}^m g(x_j). \quad (\text{A.9})$$

If  $g(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$  is concave then the inequalities in (A.8) and (A.9) are reversed.

**Weighted Geometric Mean Inequality.** For any non-negative real weights  $a_j$  such that  $\sum_{j=1}^m a_j = 1$ , and any non-negative real numbers  $x_j$

$$x_1^{a_1} x_2^{a_2} \cdots x_m^{a_m} \leq \sum_{j=1}^m a_j x_j \quad (\text{A.10})$$

**Loève's  $c_r$  Inequality.** For  $r > 0$ ,

$$\left| \sum_{j=1}^m a_j \right|^r \leq c_r \sum_{j=1}^m |a_j|^r \quad (\text{A.11})$$

where  $c_r = 1$  when  $r \leq 1$  and  $c_r = m^{r-1}$  when  $r \geq 1$ .

**$c_2$  Inequality.** For any  $m \times 1$  vectors  $\mathbf{a}$  and  $\mathbf{b}$ ,

$$(\mathbf{a} + \mathbf{b})' (\mathbf{a} + \mathbf{b}) \leq 2\mathbf{a}'\mathbf{a} + 2\mathbf{b}'\mathbf{b} \quad (\text{A.12})$$

**Hölder's Inequality.** If  $p > 1$ ,  $q > 1$ , and  $1/p + 1/q = 1$ , then for any  $m \times 1$  vectors  $\mathbf{a}$  and  $\mathbf{b}$ ,

$$\sum_{j=1}^m |a_j b_j| \leq \|\mathbf{a}\|_p \|\mathbf{b}\|_q \quad (\text{A.13})$$

**Minkowski's Inequality.** For any  $m \times 1$  vectors  $\mathbf{a}$  and  $\mathbf{b}$ , if  $p \geq 1$ , then

$$\|\mathbf{a} + \mathbf{b}\|_p \leq \|\mathbf{a}\|_p + \|\mathbf{b}\|_p \quad (\text{A.14})$$

**Schwarz Inequality.** For any  $m \times 1$  vectors  $\mathbf{a}$  and  $\mathbf{b}$ ,

$$|\mathbf{a}'\mathbf{b}| \leq \|\mathbf{a}\| \|\mathbf{b}\|. \quad (\text{A.15})$$

**Proof of Jensen's Inequality (A.8).** By the definition of convexity, for any  $\lambda \in [0, 1]$

$$g(\lambda x_1 + (1 - \lambda) x_2) \leq \lambda g(x_1) + (1 - \lambda) g(x_2). \quad (\text{A.16})$$

This implies

$$\begin{aligned} g\left(\sum_{j=1}^m a_j x_j\right) &= g\left(a_1 x_1 + (1 - a_1) \sum_{j=2}^m \frac{a_j}{1 - a_1} x_j\right) \\ &\leq a_1 g(x_1) + (1 - a_1) g\left(\sum_{j=2}^m b_j x_j\right). \end{aligned}$$

where  $b_j = a_j/(1 - a_1)$  and  $\sum_{j=2}^m b_j = 1$ . By another application of (A.16) this is bounded by

$$\begin{aligned} & a_1 g(x_1) + (1 - a_1) \left( b_2 g(x_2) + (1 - b_2) g \left( \sum_{j=2}^m c_j x_j \right) \right) \\ &= a_1 g(x_1) + a_2 g(x_2) + (1 - a_1)(1 - b_2) g \left( \sum_{j=2}^m c_j x_j \right) \end{aligned}$$

where  $c_j = b_j/(1 - b_2)$ . By repeated application of (A.16) we obtain (A.8).  $\blacksquare$

**Proof of Weighted Geometric Mean Inequality.** Since the logarithm is strictly concave, by Jensen's inequality

$$\log(x_1^{a_1} x_2^{a_2} \cdots x_m^{a_m}) = \sum_{j=1}^m a_j \log x_j \leq \log \left( \sum_{j=1}^m a_j x_j \right).$$

Applying the exponential yields (A.10).  $\blacksquare$

**Proof of Loève's  $c_r$  Inequality.** For  $r \geq 1$  this is simply a rewriting of the finite form Jensen's inequality (A.9) with  $g(u) = u^r$ . For  $r < 1$ , define  $b_j = |a_j| / \left( \sum_{j=1}^m |a_j| \right)$ . The facts that  $0 \leq b_j \leq 1$  and  $r < 1$  imply  $b_j \leq b_j^r$  and thus

$$1 = \sum_{j=1}^m b_j \leq \sum_{j=1}^m b_j^r$$

which implies

$$\left( \sum_{j=1}^m |a_j| \right)^r \leq \sum_{j=1}^m |a_j|^r.$$

The proof is completed by observing that

$$\left( \sum_{j=1}^m a_j \right)^r \leq \left( \sum_{j=1}^m |a_j| \right)^r.$$

■  
**Proof of  $c_2$  Inequality.** By the  $c_r$  inequality,  $(a_j + b_j)^2 \leq 2a_j^2 + 2b_j^2$ . Thus

$$\begin{aligned}(\mathbf{a} + \mathbf{b})'(\mathbf{a} + \mathbf{b}) &= \sum_{j=1}^m (a_j + b_j)^2 \\ &\leq 2 \sum_{j=1}^m a_j^2 + 2 \sum_{j=1}^m b_j^2 \\ &= 2\mathbf{a}'\mathbf{a} + 2\mathbf{b}'\mathbf{b}\end{aligned}$$

■  
**Proof of Hölder's Inequality.** Set  $u_j = |a_j|^p / \|\mathbf{a}\|_p^p$  and  $v_j = |b_j|^q / \|\mathbf{b}\|_q^q$  and observe that  $\sum_{j=1}^m u_j = 1$  and  $\sum_{j=1}^m v_j = 1$ . By the weighted geometric mean inequality,

$$u_j^{1/p} v_j^{1/q} \leq \frac{u_j}{p} + \frac{v_j}{q}.$$

Then since  $\sum_{j=1}^m u_j = 1$ ,  $\sum_{j=1}^m v_j = 1$  and  $1/p + 1/q = 1$

$$\frac{\sum_{j=1}^m |a_j b_j|}{\|\mathbf{a}\|_p \|\mathbf{b}\|_q} = \sum_{j=1}^m u_j^{1/p} v_j^{1/q} \leq \sum_{j=1}^m \left( \frac{u_j}{p} + \frac{v_j}{q} \right) = 1$$

which is (A.13). ■

**Proof of Minkowski's Inequality.** Set  $q = p/(p-1)$  so that  $1/p + 1/q = 1$ . Using the triangle inequality for real numbers and two applications of Hölder's

inequality

$$\begin{aligned}\|\mathbf{a} + \mathbf{b}\|_p^p &= \sum_{j=1}^m |a_j + b_j|^p \\ &= \sum_{j=1}^m |a_j + b_j| |a_j + b_j|^{p-1} \\ &\leq \sum_{j=1}^m |a_j| |a_j + b_j|^{p-1} + \sum_{j=1}^m |b_j| |a_j + b_j|^{p-1} \\ &\leq \|\mathbf{a}\|_p \left( \sum_{j=1}^m |a_j + b_j|^{(p-1)q} \right)^{1/q} + \|\mathbf{b}\|_p \left( \sum_{j=1}^m |a_j + b_j|^{(p-1)q} \right)^{1/q} \\ &= \left( \|\mathbf{a}\|_p + \|\mathbf{b}\|_p \right) \|\mathbf{a} + \mathbf{b}\|_p^{p-1}\end{aligned}$$

Solving, we find (A.14). ■

**Proof of Schwarz Inequality.** Using Hölder's inequality with  $p = q = 2$

$$|\mathbf{a}'\mathbf{b}| \leq \sum_{j=1}^m |a_j b_j| \leq \|\mathbf{a}\| \|\mathbf{b}\|$$

■

## A.14 Matrix Norms

Two common norms used for matrix spaces are the **Frobenius norm** and the **spectral norm**. We can write either as  $\|\mathbf{A}\|$ , but may write  $\|\mathbf{A}\|_F$  or  $\|\mathbf{A}\|_2$  when we want to be specific.

The **Frobenius norm** of an  $m \times k$  matrix  $\mathbf{A}$  is the Euclidean norm applied

to its elements

$$\begin{aligned}\|\mathbf{A}\|_F &= \|\text{vec}(\mathbf{A})\| \\ &= (\text{tr}(\mathbf{A}'\mathbf{A}))^{1/2} \\ &= \left( \sum_{i=1}^m \sum_{j=1}^k a_{ij}^2 \right)^{1/2}.\end{aligned}$$

If an  $m \times m$  real matrix  $\mathbf{A}$  is symmetric with eigenvalues  $\lambda_\ell$ ,  $\ell = 1, \dots, m$ , then

$$\|\mathbf{A}\|_F = \left( \sum_{\ell=1}^m \lambda_\ell^2 \right)^{1/2}.$$

To see this, by the spectral decomposition  $\mathbf{A} = \mathbf{H}\mathbf{\Lambda}\mathbf{H}'$  with  $\mathbf{H}'\mathbf{H} = \mathbf{I}$  and  $\mathbf{\Lambda} = \text{diag}\{\lambda_1, \dots, \lambda_m\}$ , so

$$\|\mathbf{A}\|_F = (\text{tr}(\mathbf{H}\mathbf{\Lambda}\mathbf{H}'\mathbf{H}\mathbf{\Lambda}\mathbf{H}'))^{1/2} = (\text{tr}(\mathbf{\Lambda}\mathbf{\Lambda}))^{1/2} = \left( \sum_{\ell=1}^m \lambda_\ell^2 \right)^{1/2}. \quad (\text{A.17})$$

A useful calculation is for any  $m \times 1$  vectors  $\mathbf{a}$  and  $\mathbf{b}$ , using (A.1),

$$\|\mathbf{a}\mathbf{b}'\|_F = \text{tr}(\mathbf{b}\mathbf{a}'\mathbf{a}\mathbf{b}')^{1/2} = (\mathbf{b}'\mathbf{b}\mathbf{a}'\mathbf{a})^{1/2} = \|\mathbf{a}\| \|\mathbf{b}\| \quad (\text{A.18})$$

and in particular

$$\|\mathbf{a}\mathbf{a}'\|_F = \|\mathbf{a}\|^2. \quad (\text{A.19})$$

The **spectral norm** of an  $m \times k$  real matrix  $\mathbf{A}$  is its largest singular value

$$\|\mathbf{A}\|_2 = s_{\max}(\mathbf{A}) = (\lambda_{\max}(\mathbf{A}'\mathbf{A}))^{1/2}$$

where  $\lambda_{\max}(\mathbf{B})$  denotes the largest eigenvalue of the matrix  $\mathbf{B}$ . Notice that

$$\lambda_{\max}(\mathbf{A}'\mathbf{A}) = \|\mathbf{A}'\mathbf{A}\|_2$$

so

$$\|\mathbf{A}\|_2 = \|\mathbf{A}'\mathbf{A}\|_2^{1/2}.$$

If  $\mathbf{A}$  is  $m \times m$  and symmetric with eigenvalues  $\lambda_j$  then

$$\|\mathbf{A}\|_2 = \max_{j \leq m} |\lambda_j|.$$

The Frobenius and spectral norms are closely related. They are equivalent when applied to a matrix of rank 1, since  $\|\mathbf{a}\mathbf{b}'\|_2 = \|\mathbf{a}\| \|\mathbf{b}\| = \|\mathbf{a}\mathbf{b}'\|_F$ . In general, for  $m \times k$  matrix  $\mathbf{A}$  with rank  $r$

$$\|\mathbf{A}\|_2 = (\lambda_{\max}(\mathbf{A}'\mathbf{A}))^{1/2} \leq \left( \sum_{j=1}^k \lambda_j(\mathbf{A}'\mathbf{A}) \right)^{1/2} = \|\mathbf{A}\|_F.$$

Since  $\mathbf{A}'\mathbf{A}$  also has rank at most  $r$ , it has at most  $r$  non-zero eigenvalues, and hence

$$\|\mathbf{A}\|_F = \left( \sum_{j=1}^k \lambda_j(\mathbf{A}'\mathbf{A}) \right)^{1/2} = \left( \sum_{j=1}^r \lambda_j(\mathbf{A}'\mathbf{A}) \right)^{1/2} \leq (r \lambda_{\max}(\mathbf{A}'\mathbf{A}))^{1/2} = \sqrt{r} \|\mathbf{A}\|_2.$$

Given any vector norm  $\|\mathbf{a}\|$  the **induced matrix norm** is defined as

$$\|\mathbf{A}\| = \sup_{\mathbf{x}'\mathbf{x}=1} \|\mathbf{A}\mathbf{x}\| = \sup_{\mathbf{x} \neq 0} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|}.$$

To see that this is a norm we need to check that it satisfies the triangle inequality. Indeed

$$\|\mathbf{A} + \mathbf{B}\| = \sup_{\mathbf{x}'\mathbf{x}=1} \|\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{x}\| \leq \sup_{\mathbf{x}'\mathbf{x}=1} \|\mathbf{A}\mathbf{x}\| + \sup_{\mathbf{x}'\mathbf{x}=1} \|\mathbf{B}\mathbf{x}\| = \|\mathbf{A}\| + \|\mathbf{B}\|.$$

For any vector  $\mathbf{x}$ , by the definition of the induced norm

$$\|\mathbf{A}\mathbf{x}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|$$

a property which is called consistent norms.

Let  $\mathbf{A}$  and  $\mathbf{B}$  be conformable and  $\|\mathbf{A}\|$  an induced matrix norm. Then using the property of consistent norms

$$\|\mathbf{A}\mathbf{B}\| = \sup_{\mathbf{x}'\mathbf{x}=1} \|\mathbf{A}\mathbf{B}\mathbf{x}\| \leq \sup_{\mathbf{x}'\mathbf{x}=1} \|\mathbf{A}\| \|\mathbf{B}\mathbf{x}\| = \|\mathbf{A}\| \|\mathbf{B}\|,$$

A matrix norm which satisfies this property is called a **sub-multiplicative norm**, and is a matrix form of the Schwarz inequality.

Of particular interest, the matrix norm induced by the Euclidean vector norm is the spectral norm. Indeed,

$$\sup_{\mathbf{x}'\mathbf{x}=1} \|\mathbf{Ax}\|^2 = \sup_{\mathbf{x}'\mathbf{x}=1} \mathbf{x}'\mathbf{A}'\mathbf{Ax} = \lambda_{\max}(\mathbf{A}'\mathbf{A}) = \|\mathbf{A}\|_2^2.$$

It follows that the spectral norm is consistent with the Euclidean norm, and is sub-multiplicative.

## A.15 Matrix Inequalities

**Schwarz Matrix Inequality:** For any  $m \times k$  and  $k \times m$  matrices  $\mathbf{A}$  and  $\mathbf{B}$ , and either the Frobenius or spectral norm,

$$\|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|. \quad (\text{A.20})$$

**Triangle Inequality:** For any  $m \times k$  matrices  $\mathbf{A}$  and  $\mathbf{B}$ , and either the Frobenius or spectral norm,

$$\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|. \quad (\text{A.21})$$

**Trace Inequality.** For any  $m \times m$  matrices  $\mathbf{A}$  and  $\mathbf{B}$  such that  $\mathbf{A}$  is symmetric and  $\mathbf{B} \geq 0$

$$\text{tr}(\mathbf{AB}) \leq \|\mathbf{A}\|_2 \text{tr}(\mathbf{B}). \quad (\text{A.22})$$

**Quadratic Inequality.** For any  $m \times 1$   $\mathbf{b}$  and  $m \times m$  symmetric matrix  $\mathbf{A}$

$$\mathbf{b}'\mathbf{Ab} \leq \|\mathbf{A}\|_2 \mathbf{b}'\mathbf{b} \quad (\text{A.23})$$

**Strong Schwarz Matrix Inequality.** For any conformable matrices  $\mathbf{A}$  and  $\mathbf{B}$

$$\|\mathbf{AB}\|_F \leq \|\mathbf{A}\|_2 \|\mathbf{B}\|_F. \quad (\text{A.24})$$

**Norm Equivalence.** For any  $m \times k$  matrix  $\mathbf{A}$  of rank  $r$

$$\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_F \leq \sqrt{r} \|\mathbf{A}\|_2. \quad (\text{A.25})$$



**Eigenvalue Product Inequality.** For any  $m \times m$  real symmetric matrices  $\mathbf{A} \geq 0$  and  $\mathbf{B} \geq 0$ , the eigenvalues  $\lambda_\ell(\mathbf{AB})$  are real and satisfy

$$\lambda_{\min}(\mathbf{A}) \lambda_{\min}(\mathbf{B}) \leq \lambda_\ell(\mathbf{AB}) = \lambda_\ell(\mathbf{A}^{1/2} \mathbf{B} \mathbf{A}^{1/2}) \leq \lambda_{\max}(\mathbf{A}) \lambda_{\max}(\mathbf{B}) \quad (\text{A.26})$$

(Zhang and Zhang, 2006, Corollary 11)

**Proof of Schwarz Matrix Inequality:** The inequality holds for the spectral norm since it is an induced norm. Now consider the Frobenius norm. Partition  $\mathbf{A}' = [\mathbf{a}_1, \dots, \mathbf{a}_n]$  and  $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_n]$ . Then by partitioned matrix multiplication, the definition of the Frobenius norm and the Schwarz inequality for vectors

$$\begin{aligned} \|\mathbf{AB}\|_F &= \left\| \begin{bmatrix} \mathbf{a}'_1 \mathbf{b}_1 & \mathbf{a}'_1 \mathbf{b}_2 & \cdots \\ \mathbf{a}'_2 \mathbf{b}_1 & \mathbf{a}'_2 \mathbf{b}_2 & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix} \right\|_F \\ &\leq \left\| \begin{bmatrix} \|\mathbf{a}_1\| \|\mathbf{b}_1\| & \|\mathbf{a}_1\| \|\mathbf{b}_2\| & \cdots \\ \|\mathbf{a}_2\| \|\mathbf{b}_1\| & \|\mathbf{a}_2\| \|\mathbf{b}_2\| & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix} \right\|_F \\ &= \left( \sum_{i=1}^m \sum_{j=1}^m \|\mathbf{a}_i\|^2 \|\mathbf{b}_j\|^2 \right)^{1/2} \\ &= \left( \sum_{i=1}^m \|\mathbf{a}_i\|^2 \right)^{1/2} \left( \sum_{i=1}^m \|\mathbf{b}_i\|^2 \right)^{1/2} \\ &= \left( \sum_{i=1}^k \sum_{j=1}^m \mathbf{a}_{ji}^2 \right)^{1/2} \left( \sum_{i=1}^m \sum_{j=1}^k \|\mathbf{b}_{ji}\|^2 \right)^{1/2} \\ &= \|\mathbf{A}\|_F \|\mathbf{B}\|_F \end{aligned}$$

■

**Proof of Triangle Inequality:** The inequality holds for the spectral norm since it is an induced norm. Now consider the Frobenius norm. Let  $\mathbf{a} =$

$\text{vec}(\mathbf{A})$  and  $\mathbf{b} = \text{vec}(\mathbf{B})$ . Then by the definition of the Frobenius norm and the Schwarz Inequality for vectors

$$\begin{aligned}\|\mathbf{A} + \mathbf{B}\|_F &= \|\text{vec}(\mathbf{A} + \mathbf{B})\|_F \\ &= \|\mathbf{a} + \mathbf{b}\| \\ &\leq \|\mathbf{a}\| + \|\mathbf{b}\| \\ &= \|\mathbf{A}\|_F + \|\mathbf{B}\|_F\end{aligned}$$

■

**Proof of Trace Inequality.** By the spectral decomposition for symmetric matrices,  $\mathbf{A} = \mathbf{H}\mathbf{\Lambda}\mathbf{H}'$  where  $\mathbf{\Lambda}$  has the eigenvalues  $\lambda_j$  of  $\mathbf{A}$  on the diagonal and  $\mathbf{H}$  is orthonormal. Define  $\mathbf{C} = \mathbf{H}'\mathbf{B}\mathbf{H}$  which has non-negative diagonal elements  $C_{jj}$  since  $\mathbf{B}$  is positive semi-definite. Then

$$\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{\Lambda C}) = \sum_{j=1}^m \lambda_j C_{jj} \leq \max_j |\lambda_j| \sum_{j=1}^m C_{jj} = \|\mathbf{A}\|_2 \text{tr}(\mathbf{C})$$

where the inequality uses the fact that  $C_{jj} \geq 0$ . But note that

$$\text{tr}(\mathbf{C}) = \text{tr}(\mathbf{H}'\mathbf{B}\mathbf{H}) = \text{tr}(\mathbf{H}\mathbf{H}'\mathbf{B}) = \text{tr}(\mathbf{B})$$

since  $\mathbf{H}$  is orthonormal. Thus  $\text{tr}(\mathbf{AB}) \leq \|\mathbf{A}\|_2 \text{tr}(\mathbf{B})$  as stated. ■

**Proof of Quadratic Inequality:** In the Trace Inequality set  $\mathbf{B} = \mathbf{b}\mathbf{b}'$  and note  $\text{tr}(\mathbf{AB}) = \mathbf{b}'\mathbf{A}\mathbf{b}$  and  $\text{tr}(\mathbf{B}) = \mathbf{b}'\mathbf{b}$ . ■

**Proof of Strong Schwarz Matrix Inequality.** By the definition of the Frobenius norm, the property of the trace, the Trace Inequality (noting that both  $\mathbf{A}'\mathbf{A}$  and  $\mathbf{B}\mathbf{B}'$  are symmetric and positive semi-definite), and the Schwarz matrix inequality

$$\begin{aligned}\|\mathbf{AB}\|_F &= (\text{tr}(\mathbf{B}'\mathbf{A}'\mathbf{AB}))^{1/2} \\ &= (\text{tr}(\mathbf{A}'\mathbf{A}\mathbf{B}\mathbf{B}'))^{1/2} \\ &\leq (\|\mathbf{A}'\mathbf{A}\|_2 \text{tr}(\mathbf{B}\mathbf{B}'))^{1/2} \\ &= \|\mathbf{A}\|_2 \|\mathbf{B}\|_F.\end{aligned}$$

■

# Appendix B

## Probability

### B.1 Foundations

The set  $S$  of all possible outcomes of an experiment is called the **sample space** for the experiment. Take the simple example of tossing a coin. There are two outcomes, heads and tails, so we can write  $S = \{H, T\}$ . If two coins are tossed in sequence, we can write the four outcomes as  $S = \{HH, HT, TH, TT\}$ .

An **event**  $A$  is any collection of possible outcomes of an experiment. An event is a subset of  $S$ , including  $S$  itself and the null set  $\emptyset$ . Continuing the two coin example, one event is  $A = \{HH, HT\}$ , the event that the first coin is heads. We say that  $A$  and  $B$  are **disjoint** or **mutually exclusive** if  $A \cap B = \emptyset$ . For example, the sets  $\{HH, HT\}$  and  $\{TH\}$  are disjoint. Furthermore, if the sets  $A_1, A_2, \dots$  are pairwise disjoint and  $\cup_{i=1}^{\infty} A_i = S$ , then the collection  $A_1, A_2, \dots$  is called a **partition** of  $S$ .

The following are elementary set operations:

**Union:**  $A \cup B = \{x : x \in A \text{ or } x \in B\}$ .

**Intersection:**  $A \cap B = \{x : x \in A \text{ and } x \in B\}$ .

**Complement:**  $A^c = \{x : x \notin A\}$ .

The following are useful properties of set operations.

**Commutativity:**  $A \cup B = B \cup A$ ;  $A \cap B = B \cap A$ .

**Associativity:**  $A \cup (B \cap C) = (A \cup B) \cap C$ ;  $A \cap (B \cup C) = (A \cap B) \cup C$ .

**Distributive Laws:**  $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$ ;  $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$ .

$$(A \cup B) \cap (A \cup C).$$

**DeMorgan's Laws:**  $(A \cup B)^c = A^c \cap B^c$ ;  $(A \cap B)^c = A^c \cup B^c$ .

A **probability function** assigns probabilities (numbers between 0 and 1) to events  $A$  in  $\mathcal{S}$ . This is straightforward when  $\mathcal{S}$  is countable; when  $\mathcal{S}$  is uncountable we must be somewhat more careful. A set  $\mathcal{B}$  is called a **sigma algebra** (or Borel field) if  $\emptyset \in \mathcal{B}$ ,  $A \in \mathcal{B}$  implies  $A^c \in \mathcal{B}$ , and  $A_1, A_2, \dots \in \mathcal{B}$  implies  $\cup_{i=1}^{\infty} A_i \in \mathcal{B}$ . A simple example is  $\{\emptyset, \mathcal{S}\}$  which is known as the trivial sigma algebra. For any sample space  $\mathcal{S}$ , let  $\mathcal{B}$  be the smallest sigma algebra which contains all of the open sets in  $\mathcal{S}$ . When  $\mathcal{S}$  is countable,  $\mathcal{B}$  is simply the collection of all subsets of  $\mathcal{S}$ , including  $\emptyset$  and  $\mathcal{S}$ . When  $\mathcal{S}$  is the real line, then  $\mathcal{B}$  is the collection of all open and closed intervals. We call  $\mathcal{B}$  the sigma algebra associated with  $\mathcal{S}$ . We only define probabilities for events contained in  $\mathcal{B}$ .

We now can give the axiomatic definition of probability. Given  $\mathcal{S}$  and  $\mathcal{B}$ , a probability function  $\Pr$  satisfies  $\Pr(\mathcal{S}) = 1$ ,  $\Pr(A) \geq 0$  for all  $A \in \mathcal{B}$ , and if  $A_1, A_2, \dots \in \mathcal{B}$  are pairwise disjoint, then  $\Pr(\cup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} \Pr(A_i)$ .

Some important properties of the probability function include the following

- $\Pr(\emptyset) = 0$
- $\Pr(A) \leq 1$
- $\Pr(A^c) = 1 - \Pr(A)$
- $\Pr(B \cap A^c) = \Pr(B) - \Pr(A \cap B)$
- $\Pr(A \cup B) = \Pr(A) + \Pr(B) - \Pr(A \cap B)$
- If  $A \subset B$  then  $\Pr(A) \leq \Pr(B)$
- Bonferroni's Inequality:  $\Pr(A \cap B) \geq \Pr(A) + \Pr(B) - 1$
- Boole's Inequality:  $\Pr(A \cup B) \leq \Pr(A) + \Pr(B)$

For some elementary probability models, it is useful to have simple rules to count the number of objects in a set. These counting rules are facilitated by using the binomial coefficients which are defined for nonnegative integers  $n$  and  $r$ ,  $n \geq r$ , as

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

When counting the number of objects in a set, there are two important distinctions. Counting may be **with replacement** or **without replacement**. Counting may be **ordered** or **unordered**. For example, consider a lottery where you pick six numbers from the set 1, 2, ..., 49. This selection is without replacement if you are not allowed to select the same number twice, and is with replacement if this is allowed. Counting is ordered or not depending on whether the sequential order of the numbers is relevant to winning the lottery. Depending on these two distinctions, we have four expressions for the number of objects (possible arrangements) of size  $r$  from  $n$  objects.

	Without Replacement	With Replacement
Ordered	$\frac{n!}{(n-r)!}$	$n^r$
Unordered	$\binom{n}{r}$	$\binom{n+r-1}{r}$

In the lottery example, if counting is unordered and without replacement, the number of potential combinations is  $\binom{49}{6} = 13,983,816$ .

If  $\Pr(B) > 0$  the **conditional probability** of the event  $A$  given the event  $B$  is

$$\Pr(A | B) = \frac{\Pr(A \cap B)}{\Pr(B)}.$$

For any  $B$ , the conditional probability function is a valid probability function where  $S$  has been replaced by  $B$ . Rearranging the definition, we can write

$$\Pr(A \cap B) = \Pr(A | B) \Pr(B)$$

which is often quite useful. We can say that the occurrence of  $B$  has no information about the likelihood of event  $A$  when  $\Pr(A | B) = \Pr(A)$ , in which case we find

$$\Pr(A \cap B) = \Pr(A) \Pr(B) \tag{B.1}$$

We say that the events  $A$  and  $B$  are **statistically independent** when (B.1) holds. Furthermore, we say that the collection of events  $A_1, \dots, A_k$  are **mutually independent** when for any subset  $\{A_i : i \in I\}$ ,

$$\Pr\left(\bigcap_{i \in I} A_i\right) = \prod_{i \in I} \Pr(A_i).$$

**Theorem 1** (*Bayes' Rule*). For any set  $B$  and any partition  $A_1, A_2, \dots$  of the sample space, then for each  $i = 1, 2, \dots$

$$\Pr(A_i | B) = \frac{\Pr(B | A_i) \Pr(A_i)}{\sum_{j=1}^{\infty} \Pr(B | A_j) \Pr(A_j)}$$

## B.2 Random Variables

A **random variable**  $X$  is a function from a sample space  $S$  into the real line. This induces a new sample space – the real line – and a new probability function on the real line. Typically, we denote random variables by uppercase letters such as  $X$ , and use lower case letters such as  $x$  for potential values and realized values. (This is in contrast to the notation adopted for most of the textbook.) For a random variable  $X$  we define its **cumulative distribution function** (CDF) as

$$F(x) = \Pr(X \leq x). \tag{B.2}$$

Sometimes we write this as  $F_X(x)$  to denote that it is the CDF of  $X$ . A function  $F(x)$  is a CDF if and only if the following three properties hold:

1.  $\lim_{x \rightarrow -\infty} F(x) = 0$  and  $\lim_{x \rightarrow \infty} F(x) = 1$
2.  $F(x)$  is nondecreasing in  $x$
3.  $F(x)$  is right-continuous

We say that the random variable  $X$  is **discrete** if  $F(x)$  is a step function. In the latter case, the range of  $X$  consists of a countable set of real numbers  $\tau_1, \dots, \tau_r$ . The probability function for  $X$  takes the form

$$\Pr(X = \tau_j) = \pi_j, \quad j = 1, \dots, r \tag{B.3}$$

where  $0 \leq \pi_j \leq 1$  and  $\sum_{j=1}^r \pi_j = 1$ .

We say that the random variable  $X$  is **continuous** if  $F(x)$  is continuous in  $x$ . In this case  $\Pr(X = \tau) = 0$  for all  $\tau \in R$  so the representation (B.3) is unavailable. Instead, we represent the relative probabilities by the **probability density function** (PDF)

$$f(x) = \frac{d}{dx} F(x)$$

so that

$$F(x) = \int_{-\infty}^x f(u)du$$

and

$$\Pr(a \leq X \leq b) = \int_a^b f(u)du.$$

These expressions only make sense if  $F(x)$  is differentiable. While there are examples of continuous random variables which do not possess a PDF, these cases are unusual and are typically ignored.

A function  $f(x)$  is a PDF if and only if  $f(x) \geq 0$  for all  $x \in R$  and  $\int_{-\infty}^{\infty} f(x)dx = 1$ .

## B.3 Expectation

For any real function  $g$ , we define the **mean** or **expectation**  $\mathbb{E}g(X)$  as follows. If  $X$  is discrete,

$$\mathbb{E}(g(X)) = \sum_{j=1}^r g(\tau_j)\pi_j,$$

and if  $X$  is continuous

$$\mathbb{E}(g(X)) = \int_{-\infty}^{\infty} g(x)f(x)dx.$$

The latter is well defined and finite if

$$\int_{-\infty}^{\infty} |g(x)|f(x)dx < \infty. \tag{B.4}$$

If (B.4) does not hold, evaluate

$$I_1 = \int_{g(x)>0} g(x)f(x)dx$$
$$I_2 = - \int_{g(x)<0} g(x)f(x)dx$$

If  $I_1 = \infty$  and  $I_2 < \infty$  then we define  $\mathbb{E}(g(X)) = \infty$ . If  $I_1 < \infty$  and  $I_2 = \infty$  then we define  $\mathbb{E}(g(X)) = -\infty$ . If both  $I_1 = \infty$  and  $I_2 = \infty$  then  $\mathbb{E}(g(X))$  is undefined.

Since  $\mathbb{E}(a + bX) = a + b\mathbb{E}(X)$ , we say that expectation is a linear operator.

For  $m > 0$ , we define the  $m^{\text{th}}$  **moment** of  $X$  as  $\mathbb{E}(X^m)$  and the  $m^{\text{th}}$  **central moment** as  $\mathbb{E}((X - \mathbb{E}X)^m)$ .

Two special moments are the **mean**  $\mu = \mathbb{E}(X)$  and **variance**  $\sigma^2 = \mathbb{E}(X - \mu)^2 = \mathbb{E}(X^2) - \mu^2$ . We call  $\sigma = \sqrt{\sigma^2}$  the **standard deviation** of  $X$ . We can also write  $\sigma^2 = \text{var}(X)$ . For example, this allows the convenient expression  $\text{var}(a + bX) = b^2 \text{var}(X)$ .

The **moment generating function** (MGF) of  $X$  is

$$M(\lambda) = \mathbb{E}(\exp(\lambda X)).$$

The MGF does not necessarily exist. However, when it does and  $\mathbb{E}(|X|^m) < \infty$  then

$$\left. \frac{d^m}{d\lambda^m} M(\lambda) \right|_{\lambda=0} = \mathbb{E}(X^m)$$

which is why it is called the moment generating function.

More generally, the **characteristic function** (CF) of  $X$  is

$$C(\lambda) = \mathbb{E}(\exp(i\lambda X))$$

where  $i = \sqrt{-1}$  is the imaginary unit. The CF always exists, and when  $\mathbb{E}(|X|^m) < \infty$

$$\left. \frac{d^m}{d\lambda^m} C(\lambda) \right|_{\lambda=0} = i^m \mathbb{E}(X^m).$$

The  $L^p$  **norm**,  $p \geq 1$ , of the random variable  $X$  is

$$\|X\|_p = (\mathbb{E}(|X|^p))^{1/p}.$$

## B.4 Gamma Function

The gamma function is defined for  $\alpha > 0$  as

$$\Gamma(\alpha) = \int_0^{\infty} x^{\alpha-1} \exp(-x) dx. \tag{B.5}$$



By integration by parts you can show that it satisfies the property

$$\Gamma(1 + \alpha) = \Gamma(\alpha)\alpha.$$

Thus for positive integers  $n$ ,

$$\Gamma(n) = (n - 1)!$$

Hence the gamma function can be viewed as a continuous version of the factorial.

Special values include

$$\Gamma(1) = \int_0^{\infty} \exp(-x) dx = 1 \tag{B.6}$$

and

$$\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}. \tag{B.7}$$

The latter holds by making the change of variables  $x = u^2$  in (B.5) and applying (5.2).

A useful fact is

$$\int_0^{\infty} y^{a-1} \exp(-by) dy = b^{-a}\Gamma(a) \tag{B.8}$$

which can be found by applying change-of-variables to the definition (B.5).

Another is for  $\alpha \in \mathbb{R}$

$$\lim_{n \rightarrow \infty} \frac{\Gamma(n + \alpha)}{\Gamma(n) n^\alpha} = 1. \tag{B.9}$$

**Sterling's formula** is an expansion for the the logarithm of the gamma function (and hence for the factorial as well).

$$\log \Gamma(\alpha) = \frac{1}{2} \log(2\pi) + \left(\alpha - \frac{1}{2}\right) \log \alpha - z + \frac{1}{12\alpha} - \frac{1}{360\alpha^3} + \frac{1}{1260\alpha^5} + \dots$$

## B.5 Common Distributions

For reference, we now list some important discrete distribution function.

### Bernoulli

$$\Pr(X = x) = p^x(1 - p)^{1-x}, \quad x = 0, 1; \quad 0 \leq p \leq 1$$

$$\mathbb{E}(X) = p$$

$$\text{var}(X) = p(1 - p)$$

### Binomial

$$\Pr(X = x) = \binom{n}{x} p^x (1 - p)^{n-x}, \quad x = 0, 1, \dots, n; \quad 0 \leq p \leq 1$$

$$\mathbb{E}(X) = np$$

$$\text{var}(X) = np(1 - p)$$

### Geometric

$$\Pr(X = x) = p(1 - p)^{x-1}, \quad x = 1, 2, \dots; \quad 0 \leq p \leq 1$$

$$\mathbb{E}(X) = \frac{1}{p}$$

$$\text{var}(X) = \frac{1 - p}{p^2}$$

### Multinomial

$$\Pr(X_1 = x_1, X_2 = x_2, \dots, X_m = x_m) = \frac{n!}{x_1! x_2! \dots x_m!} p_1^{x_1} p_2^{x_2} \dots p_m^{x_m},$$

$$x_1 + \dots + x_m = n;$$

$$p_1 + \dots + p_m = 1$$

$$\mathbb{E}(X_i) = p_i$$

$$\text{var}(X_i) = np_i(1 - p_i)$$

$$\text{cov}(X_i, X_j) = -np_i p_j$$

## Negative Binomial

$$\Pr(X = x) = \frac{\Gamma(r+x)}{x!\Gamma(r)} p^r (1-p)^{x-1}, \quad x = 0, 1, 2, \dots; \quad 0 \leq p \leq 1$$

$$\mathbb{E}(X) = \frac{r(1-p)}{p}$$

$$\text{var}(X) = \frac{r(1-p)}{p^2}$$

## Poisson

$$\Pr(X = x) = \frac{\exp(-\lambda) \lambda^x}{x!}, \quad x = 0, 1, 2, \dots, \quad \lambda > 0$$

$$\mathbb{E}(X) = \lambda$$

$$\text{var}(X) = \lambda$$

We now list some important continuous distributions.

## Beta

$$f(x) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}, \quad 0 \leq x \leq 1; \quad \alpha > 0, \beta > 0$$

$$\mathbb{E}(X) = \frac{\alpha}{\alpha + \beta}$$

$$\text{var}(X) = \frac{\alpha\beta}{(\alpha + \beta + 1)(\alpha + \beta)^2}$$

## Cauchy

$$f(x) = \frac{1}{\pi(1+x^2)}, \quad -\infty < x < \infty$$

$$\mathbb{E}(X) \text{ not defined}$$

$$\text{var}(X) = \infty$$

## Exponential

$$f(x) = \frac{1}{\theta} \exp\left(-\frac{x}{\theta}\right), \quad 0 \leq x < \infty; \quad \theta > 0$$

$$\mathbb{E}X = \theta$$

$$\text{var}(X) = \theta^2$$

## Logistic

$$f(x) = \frac{\exp(-x)}{(1 + \exp(-x))^2}, \quad -\infty < x < \infty;$$

$$\mathbb{E}(X) = 0$$

$$\text{var}(X) = \frac{\pi^2}{3}$$

## Lognormal

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma x} \exp\left(-\frac{(\log x - \mu)^2}{2\sigma^2}\right), \quad 0 \leq x < \infty; \quad \sigma > 0$$

$$\mathbb{E}(X) = \exp(\mu + \sigma^2/2)$$

$$\text{var}(X) = \exp(2\mu + 2\sigma^2) - \exp(2\mu + \sigma^2)$$

## Pareto

$$f(x) = \frac{\beta\alpha^\beta}{x^{\beta+1}}, \quad \alpha \leq x < \infty, \quad \alpha > 0, \quad \beta > 0$$

$$\mathbb{E}(X) = \frac{\beta\alpha}{\beta - 1}, \quad \beta > 1$$

$$\text{var}(X) = \frac{\beta\alpha^2}{(\beta - 1)^2(\beta - 2)}, \quad \beta > 2$$

## Uniform

$$f(x) = \frac{1}{b - a}, \quad a \leq x \leq b$$

$$\mathbb{E}(X) = \frac{a + b}{2}$$

$$\text{var}(X) = \frac{(b - a)^2}{12}$$

## Weibull

$$f(x) = \frac{\gamma}{\beta} x^{\gamma-1} \exp\left(-\frac{x^\gamma}{\beta}\right), \quad 0 \leq x < \infty; \quad \gamma > 0, \beta > 0$$

$$\mathbb{E}(X) = \beta^{1/\gamma} \Gamma\left(1 + \frac{1}{\gamma}\right)$$

$$\text{var}(X) = \beta^{2/\gamma} \left( \Gamma\left(1 + \frac{2}{\gamma}\right) - \Gamma^2\left(1 + \frac{1}{\gamma}\right) \right)$$

## Gamma

$$f(x) = \frac{1}{\Gamma(\alpha)\theta^\alpha} x^{\alpha-1} \exp\left(-\frac{x}{\theta}\right), \quad 0 \leq x < \infty; \quad \alpha > 0, \theta > 0$$

$$\mathbb{E}(X) = \alpha\theta$$

$$\text{var}(X) = \alpha\theta^2$$

## Chi-Square

$$f(x) = \frac{1}{\Gamma(r/2)2^{r/2}} x^{r/2-1} \exp\left(-\frac{x}{2}\right), \quad 0 \leq x < \infty; \quad r > 0$$

$$\mathbb{E}(X) = r$$

$$\text{var}(X) = 2r$$

## Normal

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), \quad -\infty < x < \infty; \quad -\infty < \mu < \infty, \sigma^2 > 0$$

$$\mathbb{E}(X) = \mu$$

$$\text{var}(X) = \sigma^2$$

## Student t

$$f(x) = \frac{\Gamma\left(\frac{r+1}{2}\right)}{\sqrt{r\pi}\Gamma\left(\frac{r}{2}\right)} \left(1 + \frac{x^2}{r}\right)^{-\left(\frac{r+1}{2}\right)}, \quad -\infty < x < \infty; \quad r > 0$$

$$\mathbb{E}(X) = 0 \text{ if } r > 1$$

$$\text{var}(X) = \frac{r}{r-2} \text{ if } r > 2$$

## B.6 Multivariate Random Variables

A pair of bivariate random variables  $(X, Y)$  is a function from the sample space into  $\mathbb{R}^2$ . The joint CDF of  $(X, Y)$  is

$$F(x, y) = \Pr(X \leq x, Y \leq y).$$

If  $F$  is continuous, the joint probability density function is

$$f(x, y) = \frac{\partial^2}{\partial x \partial y} F(x, y).$$

For a Borel measurable set  $A \in \mathbb{R}^2$ ,

$$\Pr((X, Y) \in A) = \int \int_A f(x, y) dx dy$$

For any measurable function  $g(x, y)$ ,

$$\mathbb{E}g(X, Y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) f(x, y) dx dy.$$

The **marginal distribution** of  $X$  is

$$\begin{aligned} F_X(x) &= \Pr(X \leq x) \\ &= \lim_{y \rightarrow \infty} F(x, y) \\ &= \int_{-\infty}^x \int_{-\infty}^{\infty} f(u, y) dy du \end{aligned}$$

so the **marginal density** of  $X$  is

$$f_X(x) = \frac{d}{dx} F_X(x) = \int_{-\infty}^{\infty} f(x, y) dy.$$

Similarly, the marginal density of  $Y$  is

$$f_Y(y) = \int_{-\infty}^{\infty} f(x, y) dx.$$

The random variables  $X$  and  $Y$  are defined to be **independent** if  $f(x, y) = f_X(x)f_Y(y)$ . Furthermore,  $X$  and  $Y$  are independent if and only if there exist functions  $g(x)$  and  $h(y)$  such that  $f(x, y) = g(x)h(y)$ .

If  $X$  and  $Y$  are independent, then

$$\begin{aligned}
 \mathbb{E}(g(X)h(Y)) &= \int \int g(x)h(y)f(y, x)dydx \\
 &= \int \int g(x)h(y)f_Y(y)f_X(x)dydx \\
 &= \int g(x)f_X(x)dx \int h(y)f_Y(y)dy \\
 &= \mathbb{E}(g(X))\mathbb{E}(h(Y)).
 \end{aligned} \tag{B.10}$$

if the expectations exist. For example, if  $X$  and  $Y$  are independent then

$$\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y).$$

Another implication of (B.10) is that if  $X$  and  $Y$  are independent and  $Z = X + Y$ , then

$$\begin{aligned}
 M_Z(\lambda) &= \mathbb{E}(\exp(\lambda(X + Y))) \\
 &= \mathbb{E}((\exp(\lambda X)\exp(\lambda Y))) \\
 &= \mathbb{E}(\exp(\lambda'X))\mathbb{E}(\exp(\lambda'Y)) \\
 &= M_X(\lambda)M_Y(\lambda).
 \end{aligned} \tag{B.11}$$

The covariance between  $X$  and  $Y$  is

$$\text{cov}(X, Y) = \sigma_{XY} = \mathbb{E}((X - \mathbb{E}X)(Y - \mathbb{E}Y)) = \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y).$$

The correlation between  $X$  and  $Y$  is

$$\text{corr}(X, Y) = \rho_{XY} = \frac{\sigma_{XY}}{\sigma_X\sigma_Y}.$$

The Cauchy-Schwarz Inequality implies that

$$|\rho_{XY}| \leq 1. \tag{B.12}$$

The correlation is a measure of linear dependence, free of units of measurement.

If  $X$  and  $Y$  are independent, then  $\sigma_{XY} = 0$  and  $\rho_{XY} = 0$ . The reverse, however, is not true. For example, if  $\mathbb{E}(X) = 0$  and  $\mathbb{E}(X^3) = 0$ , then  $\text{cov}(X, X^2) = 0$ .

A useful fact is that

$$\text{var}(X + Y) = \text{var}(X) + \text{var}(Y) + 2 \text{cov}(X, Y).$$

An implication is that if  $X$  and  $Y$  are independent, then

$$\text{var}(X + Y) = \text{var}(X) + \text{var}(Y),$$

the variance of the sum is the sum of the variances.

A  $k \times 1$  random vector  $\mathbf{X} = (X_1, \dots, X_k)'$  is a function from  $S$  to  $\mathbb{R}^k$ . Let  $\mathbf{x} = (x_1, \dots, x_k)'$  denote a vector in  $\mathbb{R}^k$ . (In this Appendix, we use bold to denote vectors. Bold capitals  $\mathbf{X}$  are random vectors and bold lower case  $\mathbf{x}$  are nonrandom vectors. Again, this is in distinction to the notation used in the bulk of the text) The vector  $\mathbf{X}$  has the distribution and density functions

$$F(\mathbf{x}) = \Pr(\mathbf{X} \leq \mathbf{x})$$

$$f(\mathbf{x}) = \frac{\partial^k}{\partial x_1 \cdots \partial x_k} F(\mathbf{x}).$$

For a measurable function  $g: \mathbb{R}^k \rightarrow \mathbb{R}^s$ , we define the expectation

$$\mathbb{E}g(\mathbf{X}) = \int_{\mathbb{R}^k} g(\mathbf{x})f(\mathbf{x})d\mathbf{x}$$

where the symbol  $d\mathbf{x}$  denotes  $dx_1 \cdots dx_k$ . In particular, we have the  $k \times 1$  multivariate mean

$$\boldsymbol{\mu} = \mathbb{E}(\mathbf{X})$$

and  $k \times k$  covariance matrix

$$\begin{aligned} \boldsymbol{\Sigma} &= \mathbb{E}((\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})') \\ &= \mathbb{E}(\mathbf{X}\mathbf{X}') - \boldsymbol{\mu}\boldsymbol{\mu}' \end{aligned}$$

If the elements of  $\mathbf{X}$  are mutually independent, then  $\boldsymbol{\Sigma}$  is a diagonal matrix and

$$\text{var}\left(\sum_{i=1}^k \mathbf{X}_i\right) = \sum_{i=1}^k \text{var}(\mathbf{X}_i)$$



## B.7 Conditional Distributions and Expectation

The **conditional density** of  $Y$  given  $\mathbf{X} = \mathbf{x}$  is defined as

$$f_{Y|\mathbf{X}}(y | \mathbf{x}) = \frac{f(\mathbf{x}, y)}{f_{\mathbf{X}}(\mathbf{x})}$$

if  $f_{\mathbf{X}}(\mathbf{x}) > 0$ . One way to derive this expression from the definition of conditional probability is

$$\begin{aligned} f_{Y|\mathbf{X}}(y | \mathbf{x}) &= \frac{\partial}{\partial y} \lim_{\varepsilon \rightarrow 0} \Pr(Y \leq y | \mathbf{x} \leq \mathbf{X} \leq \mathbf{x} + \varepsilon) \\ &= \frac{\partial}{\partial y} \lim_{\varepsilon \rightarrow 0} \frac{\Pr(\{Y \leq y\} \cap \{\mathbf{x} \leq \mathbf{X} \leq \mathbf{x} + \varepsilon\})}{\Pr(\mathbf{x} \leq \mathbf{X} \leq \mathbf{x} + \varepsilon)} \\ &= \frac{\partial}{\partial y} \lim_{\varepsilon \rightarrow 0} \frac{F(\mathbf{x} + \varepsilon, y) - F(\mathbf{x}, y)}{F_{\mathbf{X}}(\mathbf{x} + \varepsilon) - F_{\mathbf{X}}(\mathbf{x})} \\ &= \frac{\partial}{\partial y} \lim_{\varepsilon \rightarrow 0} \frac{\frac{\partial}{\partial x} F(\mathbf{x} + \varepsilon, y)}{f_{\mathbf{X}}(\mathbf{x} + \varepsilon)} \\ &= \frac{\frac{\partial^2}{\partial x \partial y} F(\mathbf{x}, y)}{f_{\mathbf{X}}(\mathbf{x})} \\ &= \frac{f(\mathbf{x}, y)}{f_{\mathbf{X}}(\mathbf{x})}. \end{aligned}$$

The **conditional mean** or **conditional expectation** is the function

$$m(\mathbf{x}) = \mathbb{E}(Y | \mathbf{X} = \mathbf{x}) = \int_{-\infty}^{\infty} y f_{Y|\mathbf{X}}(y | \mathbf{x}) dy.$$

The conditional mean  $m(\mathbf{x})$  is a function, meaning that when  $\mathbf{X}$  equals  $\mathbf{x}$ , then the expected value of  $Y$  is  $m(\mathbf{x})$ .

Similarly, we define the conditional variance of  $Y$  given  $\mathbf{X} = \mathbf{x}$  as

$$\begin{aligned} \sigma^2(\mathbf{x}) &= \text{var}(Y | \mathbf{X} = \mathbf{x}) \\ &= \mathbb{E}\left((Y - m(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}\right) \\ &= \mathbb{E}(Y^2 | \mathbf{X} = \mathbf{x}) - m(\mathbf{x})^2. \end{aligned}$$

Evaluated at  $\mathbf{x} = \mathbf{X}$ , the conditional mean  $m(\mathbf{X})$  and conditional variance  $\sigma^2(\mathbf{X})$  are random variables, functions of  $\mathbf{X}$ . We write this as  $\mathbb{E}(Y | \mathbf{X}) = m(\mathbf{X})$  and  $\text{var}(Y | \mathbf{X}) = \sigma^2(\mathbf{X})$ . For example, if  $\mathbb{E}(Y | \mathbf{X} = \mathbf{x}) = \alpha + \beta'\mathbf{x}$ , then  $\mathbb{E}(Y | \mathbf{X}) = \alpha + \beta'\mathbf{X}$ , a transformation of  $\mathbf{X}$ .

The following are important facts about conditional expectations.

**Simple Law of Iterated Expectations:**

$$\mathbb{E}(\mathbb{E}(Y | \mathbf{X})) = \mathbb{E}(Y) \tag{B.13}$$

**Proof:**

$$\begin{aligned} \mathbb{E}(\mathbb{E}(Y | \mathbf{X})) &= \mathbb{E}(m(\mathbf{X})) \\ &= \int_{-\infty}^{\infty} m(\mathbf{x})f_{\mathbf{X}}(\mathbf{x})d\mathbf{x} \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} yf_{Y|\mathbf{X}}(y | \mathbf{x})f_{\mathbf{X}}(\mathbf{x})dyd\mathbf{x} \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} yf(y, \mathbf{x})dyd\mathbf{x} \\ &= \mathbb{E}(Y). \end{aligned}$$

**Law of Iterated Expectations:**

$$\mathbb{E}(\mathbb{E}(Y | \mathbf{X}, \mathbf{Z}) | \mathbf{X}) = \mathbb{E}(Y | \mathbf{X}) \tag{B.14}$$

**Conditioning Theorem.** For any function  $g(\mathbf{x})$ ,

$$\mathbb{E}(g(\mathbf{X})Y | \mathbf{X}) = g(\mathbf{X})\mathbb{E}(Y | \mathbf{X}) \tag{B.15}$$

**Proof:** Let

$$\begin{aligned} h(\mathbf{x}) &= \mathbb{E}(g(\mathbf{X})Y | \mathbf{X} = \mathbf{x}) \\ &= \int_{-\infty}^{\infty} g(\mathbf{x})yf_{Y|\mathbf{X}}(y | \mathbf{x})dy \\ &= g(\mathbf{x}) \int_{-\infty}^{\infty} yf_{Y|\mathbf{X}}(y | \mathbf{x})dy \\ &= g(\mathbf{x})m(\mathbf{x}) \end{aligned}$$

where  $m(\mathbf{x}) = \mathbb{E}(Y | \mathbf{X} = \mathbf{x})$ . Thus  $h(\mathbf{X}) = g(\mathbf{X})m(\mathbf{X})$ , which is the same as  $\mathbb{E}(g(\mathbf{X})Y | \mathbf{X}) = g(\mathbf{X})\mathbb{E}(Y | \mathbf{X})$ .

## B.8 Transformations

Suppose that  $\mathbf{X} \in \mathbb{R}^k$  with continuous distribution function  $F_{\mathbf{X}}(\mathbf{x})$  and density  $f_{\mathbf{X}}(\mathbf{x})$ . Let  $\mathbf{Y} = \mathbf{g}(\mathbf{X})$  where  $\mathbf{g}(\mathbf{x}) : \mathbb{R}^k \rightarrow \mathbb{R}^k$  is one-to-one, differentiable, and invertible. Let  $\mathbf{h}(\mathbf{y})$  denote the inverse of  $\mathbf{g}(\mathbf{x})$ . The **Jacobian** is

$$J(\mathbf{y}) = \det \left( \frac{\partial}{\partial \mathbf{y}'} \mathbf{h}(\mathbf{y}) \right).$$

Consider the univariate case  $k = 1$ . If  $g(x)$  is an increasing function, then  $g(X) \leq Y$  if and only if  $X \leq h(Y)$ , so the distribution function of  $Y$  is

$$\begin{aligned} F_Y(y) &= \Pr(g(X) \leq y) \\ &= \Pr(X \leq h(Y)) \\ &= F_X(h(Y)). \end{aligned}$$

Taking the derivative, the density of  $Y$  is

$$f_Y(y) = \frac{d}{dy} F_Y(y) = f_X(h(Y)) \frac{d}{dy} h(y).$$

If  $g(x)$  is a decreasing function, then  $g(X) \leq Y$  if and only if  $X \geq h(Y)$ , so

$$\begin{aligned} F_Y(y) &= \Pr(g(X) \leq y) \\ &= 1 - \Pr(X \geq h(Y)) \\ &= 1 - F_X(h(Y)) \end{aligned}$$

and the density of  $Y$  is

$$f_Y(y) = -f_X(h(Y)) \frac{d}{dy} h(y).$$

We can write these two cases jointly as

$$f_Y(y) = f_X(h(Y)) |J(y)|. \tag{B.16}$$

This is known as the **change-of-variables** formula. This same formula (B.16) holds for  $k > 1$ , but its justification requires deeper results from analysis.

As one example, take the case  $X \sim U[0, 1]$  and  $Y = -\log(X)$ . Here,  $g(x) = -\log(x)$  and  $h(y) = \exp(-y)$  so the Jacobian is  $J(y) = -\exp(-y)$ . As the range of  $X$  is  $[0, 1]$ , that for  $Y$  is  $[0, \infty)$ . Since  $f_X(x) = 1$  for  $0 \leq x \leq 1$  (B.16) shows that

$$f_Y(y) = \exp(-y), \quad 0 \leq y \leq \infty,$$

an exponential density.

## B.9 Inequalities

**Jensen's Inequality.** If  $g(\cdot) : \mathbb{R}^m \rightarrow \mathbb{R}$  is convex, then for any random vector  $\mathbf{x}$  for which  $\mathbb{E} \|\mathbf{x}\| < \infty$  and  $\mathbb{E} |g(\mathbf{x})| < \infty$ ,

$$g(\mathbb{E}(\mathbf{x})) \leq \mathbb{E}(g(\mathbf{x})). \quad (\text{B.17})$$

If  $g(\cdot)$  concave, then the inequality is reversed.

**Conditional Jensen's Inequality.** If  $g(\cdot) : \mathbb{R}^m \rightarrow \mathbb{R}$  is convex, then for any random vectors  $(\mathbf{y}, \mathbf{x})$  for which  $\mathbb{E} \|\mathbf{y}\| < \infty$  and  $\mathbb{E} \|g(\mathbf{y})\| < \infty$ ,

$$g(\mathbb{E}(\mathbf{y} \mid \mathbf{x})) \leq \mathbb{E}(g(\mathbf{y}) \mid \mathbf{x}). \quad (\text{B.18})$$

If  $g(\cdot)$  concave, then the inequality is reversed.

**Conditional Expectation Inequality.** For any  $r \geq 1$  such that  $\mathbb{E} |y|^r < \infty$ , then

$$\mathbb{E} (|\mathbb{E}(y \mid \mathbf{x})|^r) \leq \mathbb{E} (|y|^r) < \infty. \quad (\text{B.19})$$

**Expectation Inequality.** For any random matrix  $\mathbf{Y}$  for which  $\mathbb{E} \|\mathbf{Y}\| < \infty$ ,

$$\|\mathbb{E}(\mathbf{Y})\| \leq \mathbb{E} \|\mathbf{Y}\|. \quad (\text{B.20})$$

**Hölder's Inequality.** If  $p > 1$  and  $q > 1$  and  $\frac{1}{p} + \frac{1}{q} = 1$ , then for any random  $m \times n$  matrices  $\mathbf{X}$  and  $\mathbf{Y}$ ,

$$\mathbb{E} \|\mathbf{X}'\mathbf{Y}\| \leq (\mathbb{E} (\|\mathbf{X}\|^p))^{1/p} (\mathbb{E} (\|\mathbf{Y}\|^q))^{1/q}. \quad (\text{B.21})$$

**Cauchy-Schwarz Inequality.** For any random  $m \times n$  matrices  $\mathbf{X}$  and  $\mathbf{Y}$ ,

$$\mathbb{E} \|\mathbf{X}'\mathbf{Y}\| \leq \left( \mathbb{E} (\|\mathbf{X}\|^2) \right)^{1/2} \left( \mathbb{E} (\|\mathbf{Y}\|^2) \right)^{1/2}. \quad (\text{B.22})$$

**Matrix Cauchy-Schwarz Inequality.** Tripathi (1999). For any random  $\mathbf{x} \in \mathbb{R}^m$  and  $\mathbf{y} \in \mathbb{R}^\ell$ ,

$$\mathbb{E} (\mathbf{y}\mathbf{x}') (\mathbb{E} (\mathbf{x}\mathbf{x}'))^{-1} \mathbb{E} (\mathbf{x}\mathbf{y}') \leq \mathbb{E} (\mathbf{y}\mathbf{y}') \quad (\text{B.23})$$

**Minkowski's Inequality.** For any random  $m \times n$  matrices  $\mathbf{X}$  and  $\mathbf{Y}$ ,

$$(\mathbb{E} (\|\mathbf{X} + \mathbf{Y}\|^p))^{1/p} \leq (\mathbb{E} (\|\mathbf{X}\|^p))^{1/p} + (\mathbb{E} (\|\mathbf{Y}\|^p))^{1/p} \quad (\text{B.24})$$

**Liapunov's Inequality.** For any random  $m \times n$  matrix  $\mathbf{X}$  and  $1 \leq r \leq p$ ,

$$(\mathbb{E} (\|\mathbf{X}\|^r))^{1/r} \leq (\mathbb{E} (\|\mathbf{X}\|^p))^{1/p} \quad (\text{B.25})$$

**Markov's Inequality (standard form).** For any random vector  $\mathbf{x}$  and non-negative function  $g(\mathbf{x}) \geq 0$ ,

$$\Pr(g(\mathbf{x}) > \alpha) \leq \alpha^{-1} \mathbb{E}(g(\mathbf{x})). \quad (\text{B.26})$$

**Markov's Inequality (strong form).** For any random vector  $\mathbf{x}$  and non-negative function  $g(\mathbf{x}) \geq 0$ ,

$$\Pr(g(\mathbf{x}) > \alpha) \leq \alpha^{-1} \mathbb{E}(g(\mathbf{x}) \mathbf{1}(g(\mathbf{x}) > \alpha)). \quad (\text{B.27})$$

**Chebyshev's Inequality.** For any random variable  $x$ ,

$$\Pr(|x - \mathbb{E}x| > \alpha) \leq \frac{\text{var}(x)}{\alpha^2}. \quad (\text{B.28})$$

**Proof of Jensen's Inequality (B.17).** Since  $g(\mathbf{u})$  is convex, at any point  $\mathbf{u}$  there is a nonempty set of subderivatives (linear surfaces touching  $g(\mathbf{u})$  at  $\mathbf{u}$  but lying below  $g(\mathbf{u})$  for all  $\mathbf{u}$ ). Let  $a + \mathbf{b}'\mathbf{u}$  be a subderivative of  $g(\mathbf{u})$  at  $\mathbf{u} = \mathbb{E}(\mathbf{x})$ . Then for all  $\mathbf{u}$ ,  $g(\mathbf{u}) \geq a + \mathbf{b}'\mathbf{u}$  yet  $g(\mathbb{E}(\mathbf{x})) = a + \mathbf{b}'\mathbb{E}(\mathbf{x})$ . Applying expectations,  $\mathbb{E}(g(\mathbf{x})) \geq a + \mathbf{b}'\mathbb{E}(\mathbf{x}) = g(\mathbb{E}(\mathbf{x}))$ , as stated. ■

**Proof of Conditional Jensen's Inequality.** The same as the proof of (B.17), but using conditional expectations. The conditional expectations exist since  $\mathbb{E}\|\mathbf{y}\| < \infty$  and  $\mathbb{E}\|g(\mathbf{y})\| < \infty$ . ■

**Proof of Conditional Expectation Inequality.** As the function  $|u|^r$  is convex for  $r \geq 1$ , the Conditional Jensen's inequality implies

$$|\mathbb{E}(y \mid \mathbf{x})|^r \leq \mathbb{E}(|y|^r \mid \mathbf{x}).$$

Taking unconditional expectations and the law of iterated expectations, we obtain

$$\mathbb{E}(|\mathbb{E}(y \mid \mathbf{x})|^r) \leq \mathbb{E}(\mathbb{E}(|y|^r \mid \mathbf{x})) = \mathbb{E}(|y|^r) < \infty$$

as required. ■

**Proof of Expectation Inequality.** By the Triangle inequality, for  $\lambda \in [0, 1]$ ,

$$\|\lambda \mathbf{U}_1 + (1 - \lambda) \mathbf{U}_2\| \leq \lambda \|\mathbf{U}_1\| + (1 - \lambda) \|\mathbf{U}_2\|$$

which shows that the matrix norm  $g(\mathbf{U}) = \|\mathbf{U}\|$  is convex. Applying Jensen's Inequality (B.17) we find (B.20). ■

**Proof of Hölder's Inequality.** Since  $\frac{1}{p} + \frac{1}{q} = 1$  an application of Jensen's Inequality (A.8) shows that for any real  $a$  and  $b$

$$\exp\left[\frac{1}{p}a + \frac{1}{q}b\right] \leq \frac{1}{p}\exp(a) + \frac{1}{q}\exp(b).$$

Setting  $u = \exp(a)$  and  $v = \exp(b)$  this implies

$$u^{1/p}v^{1/q} \leq \frac{u}{p} + \frac{v}{q}$$

and this inequality holds for any  $u > 0$  and  $v > 0$ .

Set  $u = \|\mathbf{X}\|^p / \mathbb{E}(\|\mathbf{X}\|^p)$  and  $v = \|\mathbf{Y}\|^q / \mathbb{E}(\|\mathbf{Y}\|^q)$ . Note that  $\mathbb{E}(u) = \mathbb{E}(v) = 1$ . By the matrix Schwarz Inequality (A.20),  $\|\mathbf{X}'\mathbf{Y}\| \leq \|\mathbf{X}\| \|\mathbf{Y}\|$ . Thus

$$\begin{aligned} \frac{\mathbb{E} \|\mathbf{X}'\mathbf{Y}\|}{(\mathbb{E}(\|\mathbf{X}\|^p))^{1/p} (\mathbb{E}(\|\mathbf{Y}\|^q))^{1/q}} &\leq \frac{\mathbb{E}(\|\mathbf{X}\| \|\mathbf{Y}\|)}{(\mathbb{E}(\|\mathbf{X}\|^p))^{1/p} (\mathbb{E}(\|\mathbf{Y}\|^q))^{1/q}} \\ &= \mathbb{E}\left(u^{1/p} v^{1/q}\right) \\ &\leq \mathbb{E}\left(\frac{u}{p} + \frac{v}{q}\right) \\ &= \frac{1}{p} + \frac{1}{q} \\ &= 1, \end{aligned}$$

which is (B.21).  $\blacksquare$

**Proof of Cauchy-Schwarz Inequality.** Special case of Hölder's with  $p = q = 2$ .

**Proof of Matrix Cauchy-Schwarz Inequality.** Define  $e = \mathbf{y} - (\mathbb{E}(\mathbf{y}\mathbf{x}')) (\mathbb{E}(\mathbf{x}\mathbf{x}'))^{-1} \mathbf{x}$ . Note that  $\mathbb{E}(ee') \geq 0$  is positive semi-definite. We can calculate that

$$\mathbb{E}(ee') = \mathbb{E}(\mathbf{y}\mathbf{y}') - (\mathbb{E}(\mathbf{y}\mathbf{x}')) (\mathbb{E}(\mathbf{x}\mathbf{x}'))^{-1} \mathbb{E}(\mathbf{x}\mathbf{y}').$$

Since the left-hand-side is positive semi-definite, so is the right-hand-side, which means  $\mathbb{E}(\mathbf{y}\mathbf{y}') \geq (\mathbb{E}(\mathbf{y}\mathbf{x}')) (\mathbb{E}(\mathbf{x}\mathbf{x}'))^{-1} \mathbb{E}(\mathbf{x}\mathbf{y}')$  as stated.  $\blacksquare$

**Proof of Liapunov's Inequality.** The function  $g(u) = u^{p/r}$  is convex for  $u > 0$  since  $p \geq r$ . Set  $u = \|\mathbf{X}\|^r$ . By Jensen's inequality,  $g(\mathbb{E}(u)) \leq \mathbb{E}(g(u))$  or

$$(\mathbb{E}(\|\mathbf{X}\|^r))^{p/r} \leq \mathbb{E}\left((\|\mathbf{X}\|^r)^{p/r}\right) = \mathbb{E}(\|\mathbf{X}\|^p).$$

Raising both sides to the power  $1/p$  yields  $(\mathbb{E}(\|\mathbf{X}\|^r))^{1/r} \leq (\mathbb{E}(\|\mathbf{X}\|^p))^{1/p}$  as claimed.  $\blacksquare$

**Proof of Minkowski's Inequality.** Note that by rewriting, using the triangle inequality (A.21), and then Hölder's Inequality to the two expectations

$$\begin{aligned}
\mathbb{E}(\|\mathbf{X} + \mathbf{Y}\|^p) &= \mathbb{E}\left(\|\mathbf{X} + \mathbf{Y}\| \|\mathbf{X} + \mathbf{Y}\|^{p-1}\right) \\
&\leq \mathbb{E}\left(\|\mathbf{X}\| \|\mathbf{X} + \mathbf{Y}\|^{p-1}\right) + \mathbb{E}\left(\|\mathbf{Y}\| \|\mathbf{X} + \mathbf{Y}\|^{p-1}\right) \\
&\leq (\mathbb{E}(\|\mathbf{X}\|^p))^{1/p} \mathbb{E}\left(\left(\|\mathbf{X} + \mathbf{Y}\|^{q(p-1)}\right)^{1/q}\right) \\
&\quad + (\mathbb{E}(\|\mathbf{Y}\|^p))^{1/p} \mathbb{E}\left(\left(\|\mathbf{X} + \mathbf{Y}\|^{q(p-1)}\right)^{1/q}\right) \\
&= \left((\mathbb{E}(\|\mathbf{X}\|^p))^{1/p} + (\mathbb{E}(\|\mathbf{Y}\|^p))^{1/p}\right) \mathbb{E}\left(\left(\|\mathbf{X} + \mathbf{Y}\|^p\right)^{(p-1)/p}\right)
\end{aligned}$$

where the second equality picks  $q$  to satisfy  $1/p + 1/q = 1$ , and the final equality uses this fact to make the substitution  $q = p/(p-1)$  and then collects terms. Dividing both sides by  $\mathbb{E}\left(\left(\|\mathbf{X} + \mathbf{Y}\|^p\right)^{(p-1)/p}\right)$ , we obtain (B.24). ■

**Proof of Markov's Inequality.** Let  $F$  denote the distribution function of  $\mathbf{x}$ . Then

$$\begin{aligned}
\Pr(g(\mathbf{x}) \geq \alpha) &= \int_{\{g(\mathbf{u}) \geq \alpha\}} dF(\mathbf{u}) \\
&\leq \int_{\{g(\mathbf{u}) \geq \alpha\}} \frac{g(\mathbf{u})}{\alpha} dF(\mathbf{u}) \\
&= \alpha^{-1} \int 1(g(\mathbf{u}) > \alpha) g(\mathbf{u}) dF(\mathbf{u}) \\
&= \alpha^{-1} \mathbb{E}(g(\mathbf{x}) 1(g(\mathbf{x}) > \alpha))
\end{aligned}$$

the inequality using the region of integration  $\{g(\mathbf{u}) > \alpha\}$ . This establishes the strong form (B.27). Since  $1(g(\mathbf{x}) > \alpha) \leq 1$ , the final expression is less than  $\alpha^{-1} \mathbb{E}(g(\mathbf{x}))$ , establishing the standard form (B.26). ■

**Proof of Chebyshev's Inequality.** Define  $y = (x - \mathbb{E}x)^2$  and note that  $\mathbb{E}(y) = \text{var}(x)$ . The events  $\{|x - \mathbb{E}x| > \alpha\}$  and  $\{y > \alpha^2\}$  are equal, so by an application Markov's inequality we find

$$\Pr(|x - \mathbb{E}x| > \alpha) = \Pr(y > \alpha^2) \leq \alpha^{-2} \mathbb{E}(y) = \alpha^{-2} \text{var}(x)$$

as stated. ■



# Appendix C

## Numerical Optimization

Many econometric estimators are defined by an optimization problem of the form

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmin}} Q(\boldsymbol{\theta}) \quad (\text{C.1})$$

where the parameter is  $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^m$  and the criterion function is  $Q(\boldsymbol{\theta}) : \Theta \rightarrow \mathbb{R}$ . For example NLLS, GLS, MLE and GMM estimators take this form. In most cases,  $Q(\boldsymbol{\theta})$  can be computed for given  $\boldsymbol{\theta}$ , but  $\hat{\boldsymbol{\theta}}$  is not available in closed form. In this case, numerical methods are required to obtain  $\hat{\boldsymbol{\theta}}$ .

### C.1 Grid Search

Many optimization problems are either one dimensional ( $m = 1$ ) or involve one-dimensional optimization as a sub-problem (for example, a line search). In this context grid search may be employed.

**Grid Search.** Let  $\Theta = [a, b]$  be an interval. Pick some  $\varepsilon > 0$  and set  $G = (b - a)/\varepsilon$  to be the number of gridpoints. Construct an equally spaced grid on the region  $[a, b]$  with  $G$  gridpoints, which is  $\{\boldsymbol{\theta}(j) = a + j(b - a)/G : j = 0, \dots, G\}$ . At each point evaluate the criterion function and find the gridpoint which yields the smallest value of the criterion, which is  $\boldsymbol{\theta}(\hat{j})$  where  $\hat{j} = \operatorname{argmin}_{0 \leq j \leq G} Q(\boldsymbol{\theta}(j))$ . This value  $\boldsymbol{\theta}(\hat{j})$  is the gridpoint estimate of  $\boldsymbol{\theta}$ . If the grid is sufficiently fine to capture small oscillations in  $Q(\boldsymbol{\theta})$ , the approximation error is bounded by  $\varepsilon$ , that is,  $|\boldsymbol{\theta}(\hat{j}) - \hat{\boldsymbol{\theta}}| \leq \varepsilon$ . Plots of  $Q(\boldsymbol{\theta}(j))$  against  $\boldsymbol{\theta}(j)$

can help diagnose errors in grid selection. This method is quite robust but potentially costly.

**Two-Step Grid Search.** The gridsearch method can be refined by a two-step execution. For an error bound of  $\varepsilon$  pick  $G$  so that  $G^2 = (b - a)/\varepsilon$ . For the first step define an equally spaced grid on the region  $[a, b]$  with  $G$  gridpoints, which is  $\{\boldsymbol{\theta}(j) = a + j(b - a)/G : j = 0, \dots, G\}$ . At each point evaluate the criterion function and let  $\hat{j} = \operatorname{argmin}_{0 \leq j \leq G} Q(\boldsymbol{\theta}(j))$ . For the second step define an equally spaced grid on  $[\boldsymbol{\theta}(\hat{j} - 1), \boldsymbol{\theta}(\hat{j} + 1)]$  with  $G$  gridpoints, which is  $\{\boldsymbol{\theta}'(k) = \boldsymbol{\theta}(\hat{j} - 1) + 2k(b - a)/G^2 : k = 0, \dots, G\}$ . Let  $\hat{k} = \operatorname{argmin}_{0 \leq k \leq G} Q(\boldsymbol{\theta}'(k))$ . The estimate of  $\hat{\boldsymbol{\theta}}$  is  $\boldsymbol{\theta}(\hat{k})$ . The advantage of the two-step method over a one-step grid search is that the number of function evaluations has been reduced from  $(b - a)/\varepsilon$  to  $2\sqrt{(b - a)/\varepsilon}$  which can be substantial. The disadvantage is that if the function  $Q(\boldsymbol{\theta})$  is irregular, the first-step grid may not bracket  $\hat{\boldsymbol{\theta}}$  which thus would be missed.

## C.2 Gradient Methods

Gradient Methods are iterative methods which produce a sequence  $\boldsymbol{\theta}_i : i = 1, 2, \dots$  which are designed to converge to  $\hat{\boldsymbol{\theta}}$ . All require the choice of a starting value  $\boldsymbol{\theta}_1$ , and all require the computation of the **gradient** of  $Q(\boldsymbol{\theta})$

$$\mathbf{g}(\boldsymbol{\theta}) = \frac{\partial}{\partial \boldsymbol{\theta}} Q(\boldsymbol{\theta})$$

and some require the **Hessian**

$$\mathcal{H}(\boldsymbol{\theta}) = \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} Q(\boldsymbol{\theta}).$$

If the functions  $\mathbf{g}(\boldsymbol{\theta})$  and  $\mathcal{H}(\boldsymbol{\theta})$  are not analytically available, they can be calculated numerically. Take the  $j'$ th element of  $\mathbf{g}(\boldsymbol{\theta})$ . Let  $\delta_j$  be the  $j'$ th unit vector (zeros everywhere except for a one in the  $j'$ th row). Then for  $\varepsilon$  small

$$g_j(\boldsymbol{\theta}) \simeq \frac{Q(\boldsymbol{\theta} + \delta_j \varepsilon) - Q(\boldsymbol{\theta})}{\varepsilon}.$$

Similarly,

$$g_{jk}(\boldsymbol{\theta}) \simeq \frac{Q(\boldsymbol{\theta} + \delta_j \varepsilon + \delta_k \varepsilon) - Q(\boldsymbol{\theta} + \delta_k \varepsilon) - Q(\boldsymbol{\theta} + \delta_j \varepsilon) + Q(\boldsymbol{\theta})}{\varepsilon^2}$$

In many cases, numerical derivatives can work well but can be computationally costly relative to analytic derivatives. In some cases, however, numerical derivatives can be quite unstable.

Most gradient methods are a variant of **Newton's method** which is based on a quadratic approximation. By a Taylor's expansion for  $\boldsymbol{\theta}$  close to  $\hat{\boldsymbol{\theta}}$

$$0 = \mathbf{g}(\hat{\boldsymbol{\theta}}) \simeq \mathbf{g}(\boldsymbol{\theta}) + \mathcal{H}(\boldsymbol{\theta}) (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})$$

which implies

$$\hat{\boldsymbol{\theta}} = \boldsymbol{\theta} - \mathcal{H}(\boldsymbol{\theta})^{-1} \mathbf{g}(\boldsymbol{\theta}).$$

This suggests the iteration rule

$$\hat{\boldsymbol{\theta}}_{i+1} = \boldsymbol{\theta}_i - \mathcal{H}(\boldsymbol{\theta}_i)^{-1} \mathbf{g}(\boldsymbol{\theta}_i).$$

where

One problem with Newton's method is that it will send the iterations in the wrong direction if  $\mathcal{H}(\boldsymbol{\theta}_i)$  is not positive definite. One modification to prevent this possibility is quadratic hill-climbing which sets

$$\hat{\boldsymbol{\theta}}_{i+1} = \boldsymbol{\theta}_i - (\mathcal{H}(\boldsymbol{\theta}_i) + \alpha_i \mathbf{I}_m)^{-1} \mathbf{g}(\boldsymbol{\theta}_i).$$

where  $\alpha_i$  is set just above the smallest eigenvalue of  $\mathcal{H}(\boldsymbol{\theta}_i)$  if  $\mathcal{H}(\boldsymbol{\theta}_i)$  is not positive definite.

Another productive modification is to add a scalar **steplength**  $\lambda_i$ . In this case the iteration rule takes the form

$$\boldsymbol{\theta}_{i+1} = \boldsymbol{\theta}_i - \mathbf{D}_i \mathbf{g}_i \lambda_i \tag{C.2}$$

where  $\mathbf{g}_i = \mathbf{g}(\boldsymbol{\theta}_i)$  and  $\mathbf{D}_i = \mathcal{H}(\boldsymbol{\theta}_i)^{-1}$  for Newton's method and  $\mathbf{D}_i = (\mathcal{H}(\boldsymbol{\theta}_i) + \alpha_i \mathbf{I}_m)^{-1}$  for quadratic hill-climbing.

Allowing the steplength to be a free parameter allows for a line search, a one-dimensional optimization. To pick  $\lambda_i$  write the criterion function as a function of  $\lambda$

$$Q(\lambda) = Q(\boldsymbol{\theta}_i + \mathbf{D}_i \mathbf{g}_i \lambda)$$

a one-dimensional optimization problem. There are two common methods to perform a line search. A **quadratic approximation** evaluates the first

and second derivatives of  $Q(\lambda)$  with respect to  $\lambda$ , and picks  $\lambda_i$  as the value minimizing this approximation. The **half-step** method considers the sequence  $\lambda = 1, 1/2, 1/4, 1/8, \dots$ . Each value in the sequence is considered and the criterion  $Q(\boldsymbol{\theta}_i + \mathbf{D}_i \mathbf{g}_i \lambda)$  evaluated. If the criterion has improved over  $Q(\boldsymbol{\theta}_i)$ , use this value, otherwise move to the next element in the sequence.

Newton's method does not perform well if  $Q(\boldsymbol{\theta})$  is irregular, and it can be quite computationally costly if  $\mathbf{H}(\boldsymbol{\theta})$  is not analytically available. These problems have motivated alternative choices for the weight matrix  $\mathbf{D}_i$ . These methods are called **Quasi-Newton** methods. Two popular methods are due to Davidson-Fletcher-Powell (DFP) and Broyden-Fletcher-Goldfarb-Shanno (BFGS).

Let

$$\begin{aligned}\Delta \mathbf{g}_i &= \mathbf{g}_i - \mathbf{g}_{i-1} \\ \Delta \boldsymbol{\theta}_i &= \boldsymbol{\theta}_i - \boldsymbol{\theta}_{i-1}\end{aligned}$$

and . The DFP method sets

$$\mathbf{D}_i = \mathbf{D}_{i-1} + \frac{\Delta \boldsymbol{\theta}_i \Delta \boldsymbol{\theta}_i'}{\Delta \boldsymbol{\theta}_i' \Delta \mathbf{g}_i} + \frac{\mathbf{D}_{i-1} \Delta \mathbf{g}_i \Delta \mathbf{g}_i' \mathbf{D}_{i-1}}{\Delta \mathbf{g}_i' \mathbf{D}_{i-1} \Delta \mathbf{g}_i}.$$

The BFGS methods sets

$$\mathbf{D}_i = \mathbf{D}_{i-1} + \frac{\Delta \boldsymbol{\theta}_i \Delta \boldsymbol{\theta}_i'}{\Delta \boldsymbol{\theta}_i' \Delta \mathbf{g}_i} - \frac{\Delta \boldsymbol{\theta}_i \Delta \boldsymbol{\theta}_i'}{(\Delta \boldsymbol{\theta}_i' \Delta \mathbf{g}_i)^2} \Delta \mathbf{g}_i' \mathbf{D}_{i-1} \Delta \mathbf{g}_i + \frac{\Delta \boldsymbol{\theta}_i \Delta \mathbf{g}_i' \mathbf{D}_{i-1}}{\Delta \boldsymbol{\theta}_i' \Delta \mathbf{g}_i} + \frac{\mathbf{D}_{i-1} \Delta \mathbf{g}_i \Delta \boldsymbol{\theta}_i'}{\Delta \boldsymbol{\theta}_i' \Delta \mathbf{g}_i}.$$

For any of the gradient methods, the iterations continue until the sequence has converged in some sense. This can be defined by examining whether  $|\boldsymbol{\theta}_i - \boldsymbol{\theta}_{i-1}|$ ,  $|Q(\boldsymbol{\theta}_i) - Q(\boldsymbol{\theta}_{i-1})|$  or  $|g(\boldsymbol{\theta}_i)|$  has become small.

### C.3 Derivative-Free Methods

All gradient methods can be quite poor in locating the global minimum when  $Q(\boldsymbol{\theta})$  has several local minima. Furthermore, the methods are not well defined when  $Q(\boldsymbol{\theta})$  is non-differentiable. In these cases, alternative optimization methods are required. One example is the **simplex method** of Nelder-Mead (1965).

A more recent innovation is the method of **simulated annealing (SA)**. For a review see Goffe, Ferrier, and Rodgers (1994). The SA method is a sophisticated random search. Like the gradient methods, it relies on an iterative sequence. At each iteration, a random variable is drawn and added to the current value of the parameter. If the resulting criterion is decreased, this new value is accepted. If the criterion is increased, it may still be accepted depending on the extent of the increase and another randomization. The latter property is needed to keep the algorithm from selecting a local minimum. As the iterations continue, the variance of the random innovations is shrunk. The SA algorithm stops when a large number of iterations is unable to improve the criterion. The SA method has been found to be successful at locating global minima. The downside is that it can take considerable computer time to execute.

# Bibliography

- [1] Abadir, Karim M. and Jan R. Magnus (2005): *Matrix Algebra*, Cambridge University Press.
- [2] Aitken, A.C. (1935): “On least squares and linear combinations of observations,” *Proceedings of the Royal Statistical Society*, 55, 42-48.
- [3] Akaike, H. (1973): “Information theory and an extension of the maximum likelihood principle.” In B. Petroc and F. Csake, eds., *Second International Symposium on Information Theory*.
- [4] Anderson, T.W. and H. Rubin (1949): “Estimation of the parameters of a single equation in a complete system of stochastic equations,” *The Annals of Mathematical Statistics*, 20, 46-63.
- [5] Andrews, Donald W. K. (1988): “Laws of large numbers for dependent non-identically distributed random variables,” *Econometric Theory*, 4, 458-467.
- [6] Andrews, Donald W. K. (1991), “Asymptotic normality of series estimators for nonparametric and semiparametric regression models,” *Econometrica*, 59, 307-345.
- [7] Andrews, Donald W. K. (1993), “Tests for parameter instability and structural change with unknown change point,” *Econometrica*, 61, 821-8516.
- [8] Andrews, Donald W. K. and Moshe Buchinsky: (2000): “A three-step method for choosing the number of bootstrap replications,” *Econometrica*, 68, 23-51.

- [9] Andrews, Donald W. K. and Werner Ploberger (1994): “Optimal tests when a nuisance parameter is present only under the alternative,” *Econometrica*, 62, 1383-1414.
- [10] Ash, Robert B. (1972): *Real Analysis and Probability*, Academic Press.
- [11] Basman, R. L. (1957): “A generalized classical method of linear estimation of coefficients in a structural equation,” *Econometrica*, 25, 77-83.
- [12] Bekker, P.A. (1994): “Alternative approximations to the distributions of instrumental variable estimators,” *Econometrica*, 62, 657-681.
- [13] Billingsley, Patrick (1968): *Convergence of Probability Measures*. New York: Wiley.
- [14] Billingsley, Patrick (1995): *Probability and Measure*, 3rd Edition, New York: Wiley.
- [15] Bose, A. (1988): “Edgeworth correction by bootstrap in autoregressions,” *Annals of Statistics*, 16, 1709-1722.
- [16] Box, George E. P. and Dennis R. Cox, (1964). “An analysis of transformations,” *Journal of the Royal Statistical Society, Series B*, 26, 211-252.
- [17] Breusch, T.S. and A.R. Pagan (1979): “The Lagrange multiplier test and its application to model specification in econometrics,” *Review of Economic Studies*, 47, 239-253.
- [18] Brown, B. W. and Whitney K. Newey (2002): “GMM, efficient bootstrapping, and improved inference ,” *Journal of Business and Economic Statistics*.
- [19] Card, David (1995): “Using geographic variation in college proximity to estimate the return to schooling,” in *Aspects of Labor Market Behavior: Essays in Honour of John Vanderkamp*, L.N. Christofides, E.K. Grant, and R. Swidinsky, editors. Toronto: University of Toronto Press.
- [20] Carlstein, E. (1986): “The use of subseries methods for estimating the variance of a general statistic from a stationary time series,” *Annals of Statistics*, 14, 1171-1179.

- [21] Casella, George and Roger L. Berger (2002): *Statistical Inference*, 2nd Edition, Duxbury Press.
- [22] Chamberlain, Gary (1987): "Asymptotic efficiency in estimation with conditional moment restrictions," *Journal of Econometrics*, 34, 305-334.
- [23] Choi, In and Peter C.B. Phillips (1992): "Asymptotic and finite sample distribution theory for IV estimators and tests in partially identified structural equations," *Journal of Econometrics*, 51, 113-150.
- [24] Chow, G.C. (1960): "Tests of equality between sets of coefficients in two linear regressions," *Econometrica*, 28, 591-603.
- [25] Cragg, John (1992): "Quasi-Aitken Estimation for Heterskedasticity of Unknown Form" *Journal of Econometrics*, 54, 179-201.
- [26] Davidson, James (1994): *Stochastic Limit Theory: An Introduction for Econometricians*. Oxford: Oxford University Press.
- [27] Davison, A.C. and D.V. Hinkley (1997): *Bootstrap Methods and their Application*. Cambridge University Press.
- [28] Dickey, D.A. and W.A. Fuller (1979): "Distribution of the estimators for autoregressive time series with a unit root," *Journal of the American Statistical Association*, 74, 427-431.
- [29] Donald Stephen G. and Whitney K. Newey (2001): "Choosing the number of instruments," *Econometrica*, 69, 1161-1191.
- [30] Dufour, J.M. (1997): "Some impossibility theorems in econometrics with applications to structural and dynamic models," *Econometrica*, 65, 1365-1387.
- [31] Efron, Bradley (1979): "Bootstrap methods: Another look at the jack-knife," *Annals of Statistics*, 7, 1-26.
- [32] Efron, Bradley (1982): *The Jackknife, the Bootstrap, and Other Resampling Plans*. Society for Industrial and Applied Mathematics.



- [33] Efron, Bradley and R.J. Tibshirani (1993): *An Introduction to the Bootstrap*, New York: Chapman-Hall.
- [34] Eicker, F. (1963): "Asymptotic normality and consistency of the least squares estimators for families of linear regressions," *Annals of Mathematical Statistics*, 34, 447-456.
- [35] Engle, Robert F. and Clive W. J. Granger (1987): "Co-integration and error correction: Representation, estimation and testing," *Econometrica*, 55, 251-276.
- [36] Frisch, Ragnar (1933): "Editorial," *Econometrica*, 1, 1-4.
- [37] Frisch, Ragnar and F. Waugh (1933): "Partial time regressions as compared with individual trends," *Econometrica*, 1, 387-401.
- [38] Gallant, A. Ronald and D.W. Nychka (1987): "Seminonparametric maximum likelihood estimation," *Econometrica*, 55, 363-390.
- [39] Gallant, A. Ronald and Halbert White (1988): *A Unified Theory of Estimation and Inference for Nonlinear Dynamic Models*. New York: Basil Blackwell.
- [40] Galton, Francis (1886): "Regression Towards Mediocrity in Hereditary Stature," *The Journal of the Anthropological Institute of Great Britain and Ireland*, 15, 246-263.
- [41] Goldberger, Arthur S. (1964): *Econometric Theory*, Wiley.
- [42] Goldberger, Arthur S. (1968): *Topics in Regression Analysis*, Macmillan
- [43] Goldberger, Arthur S. (1991): *A Course in Econometrics*. Cambridge: Harvard University Press.
- [44] Goffe, W.L., G.D. Ferrier and J. Rogers (1994): "Global optimization of statistical functions with simulated annealing," *Journal of Econometrics*, 60, 65-99.
- [45] Gosset, William S. (a.k.a. "Student") (1908): "The probable error of a mean," *Biometrika*, 6, 1-25.

- [46] Gauss, K.F. (1809): "Theoria motus corporum coelestium," in *Werke*, Vol. VII, 240-254.
- [47] Granger, Clive W. J. (1969): "Investigating causal relations by econometric models and cross-spectral methods," *Econometrica*, 37, 424-438.
- [48] Granger, Clive W. J. (1981): "Some properties of time series data and their use in econometric specification," *Journal of Econometrics*, 16, 121-130.
- [49] Granger, Clive W. J. and Timo Teräsvirta (1993): *Modelling Nonlinear Economic Relationships*, Oxford University Press, Oxford.
- [50] Gregory, A. and M. Veall (1985): "On formulating Wald tests of nonlinear restrictions," *Econometrica*, 53, 1465-1468,
- [51] Haavelmo, T. (1944): "The probability approach in econometrics," *Econometrica*, supplement, 12.
- [52] Hall, A. R. (2000): "Covariance matrix estimation and the power of the overidentifying restrictions test," *Econometrica*, 68, 1517-1527,
- [53] Hall, P. (1992): *The Bootstrap and Edgeworth Expansion*, New York: Springer-Verlag.
- [54] Hall, P. (1994): "Methodology and theory for the bootstrap," *Handbook of Econometrics, Vol. IV*, eds. R.F. Engle and D.L. McFadden. New York: Elsevier Science.
- [55] Hall, P. and J.L. Horowitz (1996): "Bootstrap critical values for tests based on Generalized-Method-of-Moments estimation," *Econometrica*, 64, 891-916.
- [56] Hahn, J. (1996): "A note on bootstrapping generalized method of moments estimators," *Econometric Theory*, 12, 187-197.
- [57] Hamilton, James D. (1994) *Time Series Analysis*.
- [58] Hansen, Bruce E. (1992): "Efficient estimation and testing of cointegrating vectors in the presence of deterministic trends," *Journal of Econometrics*, 53, 87-121.

- [59] Hansen, Bruce E. (1996): "Inference when a nuisance parameter is not identified under the null hypothesis," *Econometrica*, 64, 413-430.
- [60] Hansen, Bruce E. (2006): "Edgeworth expansions for the Wald and GMM statistics for nonlinear restrictions," *Econometric Theory and Practice: Frontiers of Analysis and Applied Research*, edited by Dean Corbae, Steven N. Durlauf and Bruce E. Hansen. Cambridge University Press.
- [61] Hansen, Lars Peter (1982): "Large sample properties of generalized method of moments estimators," *Econometrica*, 50, 1029-1054.
- [62] Hansen, Lars Peter, John Heaton, and A. Yaron (1996): "Finite sample properties of some alternative GMM estimators," *Journal of Business and Economic Statistics*, 14, 262-280.
- [63] Hausman, J.A. (1978): "Specification tests in econometrics," *Econometrica*, 46, 1251-1271.
- [64] Heckman, J. (1979): "Sample selection bias as a specification error," *Econometrica*, 47, 153-161.
- [65] Horn, S.D., R.A. Horn, and D.B. Duncan. (1975) "Estimating heteroscedastic variances in linear model," *Journal of the American Statistical Association*, 70, 380-385.
- [66] Horowitz, Joel (2001): "The Bootstrap," *Handbook of Econometrics*, Vol. 5, J.J. Heckman and E.E. Leamer, eds., Elsevier Science, 3159-3228.
- [67] Imbens, G.W. (1997): "One step estimators for over-identified generalized method of moments models," *Review of Economic Studies*, 64, 359-383.
- [68] Imbens, G.W., R.H. Spady and P. Johnson (1998): "Information theoretic approaches to inference in moment condition models," *Econometrica*, 66, 333-357.
- [69] Jarque, C.M. and A.K. Bera (1980): "Efficient tests for normality, homoskedasticity and serial independence of regression residuals," *Economic Letters*, 6, 255-259.

- [70] Johansen, S. (1988): “Statistical analysis of cointegrating vectors,” *Journal of Economic Dynamics and Control*, 12, 231-254.
- [71] Johansen, S. (1991): “Estimation and hypothesis testing of cointegration vectors in the presence of linear trend,” *Econometrica*, 59, 1551-1580.
- [72] Johansen, S. (1995): *Likelihood-Based Inference in Cointegrated Vector Auto-Regressive Models*, Oxford University Press.
- [73] Johansen, S. and K. Juselius (1992): “Testing structural hypotheses in a multivariate cointegration analysis of the PPP and the UIP for the UK,” *Journal of Econometrics*, 53, 211-244.
- [74] Kitamura, Y. (2001): “Asymptotic optimality and empirical likelihood for testing moment restrictions,” *Econometrica*, 69, 1661-1672.
- [75] Kitamura, Y. and M. Stutzer (1997): “An information-theoretic alternative to generalized method of moments,” *Econometrica*, 65, 861-874..
- [76] Koenker, Roger (2005): *Quantile Regression*. Cambridge University Press.
- [77] Kunsch, H.R. (1989): “The jackknife and the bootstrap for general stationary observations,” *Annals of Statistics*, 17, 1217-1241.
- [78] Kwiatkowski, D., P.C.B. Phillips, P. Schmidt, and Y. Shin (1992): “Testing the null hypothesis of stationarity against the alternative of a unit root: How sure are we that economic time series have a unit root?” *Journal of Econometrics*, 54, 159-178.
- [79] Lafontaine, F. and K.J. White (1986): “Obtaining any Wald statistic you want,” *Economics Letters*, 21, 35-40.
- [80] Lehmann, E.L. and George Casella (1998): *Theory of Point Estimation*, 2nd Edition, Springer.
- [81] Lehmann, E.L. and Joseph P. Romano (2005): *Testing Statistical Hypotheses*, 3rd Edition, Springer.

- [82] Lindeberg, Jarl Waldemar, (1922): “Eine neue Herleitung des Exponentialgesetzes in der Wahrscheinlichkeitsrechnung,” *Mathematische Zeitschrift*, 15, 211-225.
- [83] Li, Qi and Jeffrey Racine (2007) *Nonparametric Econometrics*.
- [84] Lovell, M.C. (1963): “Seasonal adjustment of economic time series,” *Journal of the American Statistical Association*, 58, 993-1010.
- [85] MacKinnon, James G. (1990): “Critical values for cointegration,” in Engle, R.F. and C.W. Granger (eds.) *Long-Run Economic Relationships: Readings in Cointegration*, Oxford, Oxford University Press.
- [86] MacKinnon, James G. and Halbert White (1985): “Some heteroskedasticity-consistent covariance matrix estimators with improved finite sample properties,” *Journal of Econometrics*, 29, 305-325.
- [87] Magnus, J. R., and H. Neudecker (1988): *Matrix Differential Calculus with Applications in Statistics and Econometrics*, New York: John Wiley and Sons.
- [88] Mann, H.B. and A. Wald (1943). “On stochastic limit and order relationships,” *The Annals of Mathematical Statistics* 14, 217–226.
- [89] Muirhead, R.J. (1982): *Aspects of Multivariate Statistical Theory*. New York: Wiley.
- [90] Nelder, J. and R. Mead (1965): “A simplex method for function minimization,” *Computer Journal*, 7, 308-313.
- [91] Nerlove, Marc (1963): “Returns to Scale in Electricity Supply,” Chapter 7 of *Measurement in Economics* (C. Christ, et al, eds.). Stanford: Stanford University Press, 167-198.
- [92] Newey, Whitney K. (1990): “Semiparametric efficiency bounds,” *Journal of Applied Econometrics*, 5, 99-135.
- [93] Newey, Whitney K. (1997): “Convergence rates and asymptotic normality for series estimators,” *Journal of Econometrics*, 79, 147-168.

- [94] Newey, Whitney K. and Daniel L. McFadden (1994): “Large Sample Estimation and Hypothesis Testing,” in Robert Engle and Daniel McFadden, (eds.) *Handbook of Econometrics*, vol. IV, 2111-2245, North Holland: Amsterdam.
- [95] Newey, Whitney K. and Kenneth D. West (1987): “Hypothesis testing with efficient method of moments estimation,” *International Economic Review*, 28, 777-787.
- [96] Owen, Art B. (1988): “Empirical likelihood ratio confidence intervals for a single functional,” *Biometrika*, 75, 237-249.
- [97] Owen, Art B. (2001): *Empirical Likelihood*. New York: Chapman & Hall.
- [98] Park, Joon Y. and Peter C. B. Phillips (1988): “On the formulation of Wald tests of nonlinear restrictions,” *Econometrica*, 56, 1065-1083,
- [99] Phillips, Peter C.B. (1989): “Partially identified econometric models,” *Econometric Theory*, 5, 181-240.
- [100] Phillips, Peter C.B. and Sam Ouliaris (1990): “Asymptotic properties of residual based tests for cointegration,” *Econometrica*, 58, 165-193.
- [101] Politis, D.N. and J.P. Romano (1996): “The stationary bootstrap,” *Journal of the American Statistical Association*, 89, 1303-1313.
- [102] Potscher, B.M. (1991): “Effects of model selection on inference,” *Econometric Theory*, 7, 163-185.
- [103] Qin, J. and J. Lawless (1994): “Empirical likelihood and general estimating equations,” *The Annals of Statistics*, 22, 300-325.
- [104] Ramsey, J. B. (1969): “Tests for specification errors in classical linear least-squares regression analysis,” *Journal of the Royal Statistical Society*, Series B, 31, 350-371.
- [105] Rudin, W. (1987): *Real and Complex Analysis*, 3rd edition. New York: McGraw-Hill.

- [106] Runge, Carl (1901): “Über empirische Funktionen und die Interpolation zwischen äquidistanten Ordinaten,” *Zeitschrift für Mathematik und Physik*, 46, 224-243.
- [107] Said, S.E. and D.A. Dickey (1984): “Testing for unit roots in autoregressive-moving average models of unknown order,” *Biometrika*, 71, 599-608.
- [108] Secrist, Horace (1933): *The Triumph of Mediocrity in Business*. Evanston: Northwestern University.
- [109] Shao, J. and D. Tu (1995): *The Jackknife and Bootstrap*. NY: Springer.
- [110] Sargan, J.D. (1958): “The estimation of economic relationships using instrumental variables,” *Econometrica*, 26, 393-415.
- [111] Shao, Jun (2003): *Mathematical Statistics*, 2nd edition, Springer.
- [112] Sheather, S.J. and M.C. Jones (1991): “A reliable data-based bandwidth selection method for kernel density estimation,” *Journal of the Royal Statistical Society, Series B*, 53, 683-690.
- [113] Shin, Y. (1994): “A residual-based test of the null of cointegration against the alternative of no cointegration,” *Econometric Theory*, 10, 91-115.
- [114] Silverman, B.W. (1986): *Density Estimation for Statistics and Data Analysis*. London: Chapman and Hall.
- [115] Sims, C.A. (1972): “Money, income and causality,” *American Economic Review*, 62, 540-552.
- [116] Sims, C.A. (1980): “Macroeconomics and reality,” *Econometrica*, 48, 1-48.
- [117] Staiger, D. and James H. Stock (1997): “Instrumental variables regression with weak instruments,” *Econometrica*, 65, 557-586.
- [118] Stock, James H. (1987): “Asymptotic properties of least squares estimators of cointegrating vectors,” *Econometrica*, 55, 1035-1056.

- [119] Stock, James H. (1991): “Confidence intervals for the largest autoregressive root in U.S. macroeconomic time series,” *Journal of Monetary Economics*, 28, 435-460.
- [120] Stock, James H. and Jonathan H. Wright (2000): “GMM with weak identification,” *Econometrica*, 68, 1055-1096.
- [121] Stock, James H. and Mark W. Watson (2010): *Introduction to Econometrics*, 3rd edition, Addison-Wesley.
- [122] Stone, Marshall H. (1937): “Applications of the Theory of Boolean Rings to General Topology,” *Transactions of the American Mathematical Society*, 41, 375-481.
- [123] Stone, Marshall H. (1948): “The Generalized Weierstrass Approximation Theorem,” *Mathematics Magazine*, 21, 167-184.
- [124] Theil, Henri. (1953): “Repeated least squares applied to complete equation systems,” The Hague, Central Planning Bureau, mimeo.
- [125] Theil, Henri (1961): *Economic Forecasts and Policy*. Amsterdam: North Holland.
- [126] Theil, Henri. (1971): *Principles of Econometrics*, New York: Wiley.
- [127] Tobin, James (1958): “Estimation of relationships for limited dependent variables,” *Econometrica*, 26, 24-36.
- [128] Tripathi, Gautam (1999): “A matrix extension of the Cauchy-Schwarz inequality,” *Economics Letters*, 63, 1-3.
- [129] van der Vaart, A.W. (1998): *Asymptotic Statistics*, Cambridge University Press.
- [130] Wald, A. (1943): “Tests of statistical hypotheses concerning several parameters when the number of observations is large,” *Transactions of the American Mathematical Society*, 54, 426-482.
- [131] Wang, J. and E. Zivot (1998): “Inference on structural parameters in instrumental variables regression with weak instruments,” *Econometrica*, 66, 1389-1404.



- [132] Weierstrass, K. (1885): “Über die analytische Darstellbarkeit sogenannter willkürlicher Functionen einer reellen Veränderlichen,” *Sitzungsberichte der Königlich Preussischen Akademie der Wissenschaften zu Berlin*, 1885.
- [133] White, Halbert (1980): “A heteroskedasticity-consistent covariance matrix estimator and a direct test for heteroskedasticity,” *Econometrica*, 48, 817-838.
- [134] White, Halbert (1984): *Asymptotic Theory for Econometricians*, Academic Press.
- [135] Wooldridge, Jeffrey M. (2010) *Econometric Analysis of Cross Section and Panel Data*, 2nd edition, MIT Press.
- [136] Wooldridge, Jeffrey M. (2009) *Introductory Econometrics: A Modern Approach*, 4th edition, Southwestern
- [137] Zellner, Arnold. (1962): “An efficient method of estimating seemingly unrelated regressions, and tests for aggregation bias,” *Journal of the American Statistical Association*, 57, 348-368.
- [138] Zhang, Fuzhen and Qingling Zhang (2006): “Eigenvalue inequalities for matrix product,” *IEEE Transactions on Automatic Control*, 51, 1506-1509.)