

Economics 686 Lecture Notes

Gordon R. Fisher¹

Department of Economics

Concordia University

1455 de Maisonneuve West

Montreal, Quebec

H3G 1M8

May 20, 2001

¹Copyright Gordon R. Fisher, 1998. Not to be cited or quoted without the permission of the author.

Contents

1	Basic Matrix Algebra	5
1.1	Vectors in \mathbb{R}^n	5
1.2	Some Simple Matrices	9
1.3	Some Special Types of Matrices	13
1.4	Orthogonal Matrices and Eigenvalues	14
1.5	A Special Matrix Product	16
1.6	Scalar Functions of Square Matrices	16
1.7	The Rank of a Matrix	19
1.8	Determinants as the Product of Eigenvalues	20
1.9	Determinants as Areas and Volumes	21
1.10	Vectorization	25
1.11	An Important Partitioned Inverse	27
1.12	Useful Inverses and Projection Matrices	29
2	Vector Spaces	33
2.1	Introduction	33
2.2	Dimensionality, Basis and Co-ordinates	37
2.3	Euclidean Space	42
2.4	Subspaces	53

3	Projection and Least Squares	58
3.1	Projection Matrices	58
3.2	Means and Decompositions	64
3.3	Orthogonality and Least Squares	70
3.4	Restricted Least Squares	76
3.5	The Frisch-Waugh Theorem	81
3.6	Invariance	86
4	The Multivariate Normal and Related Distributions	90
4.1	Probability and Random Variables	90
4.2	Expectations	95
4.3	Special Covariance Structures	100
4.4	The Multivariate Normal Distribution	101
4.5	The Non-central Chi-square Distribution	107
4.6	Non-central F-distribution	110
4.7	Non-central t-distribution	110
4.8	Cochran's Theorem	110
5	Gauss-Markov Estimation and the Linear Hypothesis	112
5.1	The Gauss-Markov Theorem	112
5.2	The Linear Hypothesis	116

5.2.1	The Setting	116
5.2.2	Decomposition	118
5.2.3	Alternative Forms of the F-test	120
5.2.4	The Lagrange Multiplier Principle	124
5.3	Applications of the F-test	126
5.3.1	Testing the Di®erence Between Two Means	126
5.3.2	Testing a Regression Coe±cient	128
5.3.3	Testing a Linear Combination of Regression Coe±cients	129
5.3.4	Testing a Block of Regression Coe±cients	130
5.3.5	Durbin-Hausman Testing	131
5.3.6	Testing a Set of r Restrictions	135
5.3.7	Testing for Structural Change	137
5.3.8	Recursive Residuals	141
6	Limits, Continuity and Convergence	146
6.1	Introduction	146
6.2	Real Numbers	146
6.3	Continuity	148
6.4	The Order of a Sequence	148
6.5	Almost Sure Convergence	149
6.6	Convergence in Probability	150

6.7	Convergence in Distribution	151
7	Maximum-likelihood Estimation Procedures and Associated Tests of Significance	154
7.1	The General Problem	154
7.2	General Justification	155
7.3	Notation	157
7.4	Regularity	159
7.5	The Cramér-Rao Inequality	162
7.6	Maximum-likelihood Estimation in \mathbb{R}^k	164
7.7	Maximum-likelihood Estimation in \mathbb{R}^k	167
7.8	Associated Tests of Significance	172

$x \in \mathbb{R}^n$ and

$$x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

$x_i \in \mathbb{R}^n$; $i = 1, 2, \dots, n$. The column vector, x , may also be represented as a row by transposing it to yield x^T :

$$x^T = (x_1 \ x_2 \ \dots \ x_n)$$

Normally, if $x \in \mathbb{R}^n$, then x is to be thought of as a column or $(n \times 1)$ vector. If w and z are each elements of \mathbb{R}^n , they are said to be conformable, whereupon their sum, $w + z$, and their scalar product may be defined. If the elements of w and z are w_i and z_i , $i = 1, 2, \dots, n$, then

$$w + z = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{pmatrix} + \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_n \end{pmatrix} = \begin{pmatrix} w_1 + z_1 \\ w_2 + z_2 \\ \vdots \\ w_n + z_n \end{pmatrix}$$

The scalar (or inner) product is defined by

$$w^T z = \sum_i w_i z_i \in \mathbb{R}.$$

Of course, because $w^T z$ is a real number (referred to as a scalar, in general),

$$w^T z = z^T w.$$

The scalar product of two (conformable) vectors plays an important role because it is a means of measuring length. For $x \in \mathbb{R}^n$, the (Euclidean) length of x is written $\|x\|$ and is defined by

$$\|x\| = \sqrt{(x^T x)} = \sqrt{x_1^2 + \dots + x_n^2},$$

thus ensuring that length is always a positive number, or zero.

Consider now $k < n$ vectors $x_{:1}; x_{:2}; \dots; x_{:k}$, from \mathbb{R}^n . These may be arranged as n rows and k columns denoted by X

$$[x_{:1}; x_{:2}; \dots; x_{:k}] = \begin{matrix} & \begin{matrix} 2 & & & 3 \end{matrix} \\ \begin{matrix} 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 4 \end{matrix} & \begin{matrix} x_{11} & x_{12} & \dots & x_{1k} \\ x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nk} \end{matrix} & \begin{matrix} 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 5 \end{matrix} \end{matrix} = X.$$

X is an $(n \times k)$ matrix. X may be expressed as k columns of n elements x_{ij} in \mathbb{R}^n , $j = 1; 2; \dots; k$, or as n rows of k elements x_{ij} in \mathbb{R}^k , $i = 1; 2; \dots; n$:

$$X = \begin{matrix} & \begin{matrix} 2 & 3 \end{matrix} \\ \begin{matrix} 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 4 \end{matrix} & \begin{matrix} x_{1:} \\ x_{2:} \\ \vdots \\ x_{n:} \end{matrix} & \begin{matrix} 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 5 \end{matrix} \end{matrix}.$$

Combining the two alternative vector notations, $x_{:j}$ and $x_{i:}$, the nk individual real numbers comprising X are x_{ij} , $i = 1; 2; \dots; n$; $j = 1; 2; \dots; k$, the subscripts $i; j$ denoting position in the i 'th row and j 'th column.

Some discussion will be centered on the linear model

$$y_i = x_{i1} \beta_1 + x_{i2} \beta_2 + \dots + x_{ik} \beta_k + \epsilon_i \quad (1.1)$$

in which $i = 1; 2; \dots; n$, $n > k$ as above. Equation (1.1) may also be written as

$$y_i = \sum_j x_{ij} \beta_j + \epsilon_i \quad (1.2)$$

where $j = 1; 2; \dots; k$; or as

$$y = x_{:1} \beta_1 + x_{:2} \beta_2 + \dots + x_{:k} \beta_k + \epsilon \quad (1.3)$$

where y , the k vectors $x_{:j}$ and ϵ are all n -tuples of elements y_i , x_{ij} and ϵ_i respectively.

Finally, let $\beta \in \mathbb{R}^k$ be the $(k \times 1)$ vector of elements $\beta_1; \beta_2; \dots; \beta_k$, then equations (1:1), (1:2) and (1:3) may be written compactly as

$$y = X\beta + \epsilon, \quad (1.4)$$

X being an $(n \times k)$ matrix. In the usual interpretation, $y; x_{:1}; x_{:2}; \dots; x_{:k}$ are vectors of n observations on each of the $k + 1$ variables, β is a vector of k coefficients and ϵ is regarded as a vector of n unknown errors. Given $y; X$ and β , then $\epsilon = y - X\beta$. In practice, neither β nor ϵ is known and hence must be estimated in some way. To begin to study this problem, it is useful to understand some of the theory of linear spaces. As a first step, some of the algebra of matrices is reviewed.

1.2 Some Simple Matrices

Generally, vectors will be denoted by lower case letters and matrices by upper case letters: thus x will generally be a column vector and X a matrix. Consider the $(n \times m)$ matrix A of elements $a_{ij} \in \mathbb{R}; i = 1; 2; \dots; n; j = 1; 2; \dots; m$. It is common to write $A = [a_{ij}]$ to describe a matrix and its elements. If $n \neq m$, then A is said to be rectangular; if $n = m$, then A is said to be square. Matrices of elements from \mathbb{R} are said to be real matrices, meaning matrices comprising real elements. A matrix comprising zero elements is called the null matrix.

When A is square and $a_{ij} = \delta_{ij}$ where δ_{ij} is Kronecker's delta (i.e. $\delta_{ij} = 1$ for $i = j$ and $\delta_{ij} = 0$ for $i \neq j, i; j = 1; 2; \dots; n$) then

$$A = \begin{matrix} & \begin{matrix} 2 & & & & 3 \end{matrix} \\ \begin{matrix} 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 4 \end{matrix} & \begin{matrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{matrix} \end{matrix} = I.$$

I is the identity matrix. When I is $(n \times n)$, it is common to write the identity matrix as I_n so that its order, $(n \times n)$, is recognized explicitly in the notation. Similarly, an $(n \times m)$ matrix is said to be of order $(n \times m)$. If a matrix is simply of order n , then it is square with n rows and n columns.

As in the case of the matrix X in section 1.1, the notation

$$A = [a_{ij}] = \begin{matrix} & \begin{matrix} 2 & & & 3 \end{matrix} \\ \begin{matrix} 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 4 \end{matrix} & \begin{matrix} a_{11} & a_{21} & \dots & a_{n1} \\ a_{12} & a_{22} & \dots & a_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ a_{1m} & a_{2m} & \dots & a_{nm} \end{matrix} \end{matrix}$$

Thus, A is (m × n).

If A and B are of the same order, then A + B and A - B are defined as the matrices C and D respectively as follows

$$A + B = [a_{ij} + b_{ij}] = C$$

$$A - B = [a_{ij} - b_{ij}] = D.$$

Matrices that are not of the same order cannot be added together or subtracted from one another. If $\lambda \in \mathbb{R}$, i.e. λ is a scalar, then

$$\lambda A = A_\lambda = [\lambda a_{ij}].$$

This operation is called scalar multiplication.

The product of two matrices A and B, say AB, is only defined when the number of columns of A is equal to the number of rows of B. Thus if A is (n × k) and B is (k × m) then

$$AB = \sum_{i=1}^k a_{iq} b_{qj}$$

for $q = 1; 2; \dots; k$; $i = 1; 2; \dots; n$ and $j = 1; 2; \dots; m$. The product, AB , then has order $(n \times m)$. Writing down the orders explicitly:

$$\begin{matrix} A & B & C \\ (n \times k) & (k \times m) & (n \times m) \end{matrix} = \begin{matrix} \mathbf{h} \times \mathbf{i} \\ a_{iq} b_{qj} \end{matrix} :$$

Thus, the inner orders must be the same (k) and the outer orders ($n; m$) determine the order of the product. It follows that, if AB is defined, BA is not defined unless the outer orders are the same ($n = m$). When ($n = m$), AB is $(n \times n)$ and BA is $(k \times k)$, whereupon $AB \neq BA$.

Notice that, for vectors x and y in \mathbb{R}^n , while the scalar, or inner, product $x \cdot y$ is a scalar, the outer product xy^T is an $(n \times n)$ matrix.

Matrix multiplication and addition have the distributive property

$$(A + B)C = AC + BC$$

and both operations are associative

$$(A + B) + C = A + (B + C)$$

$$(AB)C = A(BC).$$

Since $AB \neq BA$ in general, matrix multiplication is not commutative. For any $(n \times m)$ matrix A ,

$$I_n A = A = A I_m.$$

A square matrix, A , is idempotent if $A = A^2$. An idempotent matrix is often called a projection matrix for reasons that will become obvious. If an idempotent matrix is also symmetric, $A = A^T = A^2$, then it is said to be an orthogonal projection matrix. Projection matrices play a crucial role in the statistical analysis of linear models like (1.4).

If x is a vector and A is square such that $x^T A x = x^T A^T x > 0 \quad \forall x \neq 0$, then A is said to be positive definite (pd); if $x^T A x = x^T A^T x \geq 0 \quad \forall x \neq 0$, then A is non-negative definite (nnd) or positive semi-definite. If $x^T A x = -x^T A^T x > 0 \quad \forall x \neq 0$, then A is negative definite. If $A = A^T = A^2$, then A is nnd unless it is pd, whereupon $A = I$.

1.4 Orthogonal Matrices and Eigenvalues

An orthogonal matrix, M , is square and invertible such that $M^{-1} = M^T$. If A is a real, square symmetric matrix, then there exists an orthogonal matrix, M , such that

$$M^T A M = \Lambda \tag{1.6}$$

in which Λ is diagonal having diagonal elements λ_i (say, $i = 1; 2; \dots; n$). If A is invertible, then all the λ_i will be non-zero. If the largest square sub-matrix of A that is invertible has order $r < n$, then r of the λ_i will be non-zero and the remaining $(n - r)$ will be zero. The non-zero elements λ_i are called eigenvalues or characteristic

roots or proper values. Notice that

$$M^T A M = \Lambda, \quad A M = M \Lambda. \quad (1.7)$$

Let the i -th column of M be m_i . Then, if λ_i is the i -th eigenvalue of Λ , (1:7) implies $A m_i = \lambda_i m_i$ and hence

$$(A - \lambda_i I) m_i = 0. \quad (1.8)$$

Equation (1:8) will hold for any column, m_i , of M choosing the appropriate λ_i from the diagonal elements of Λ . Thus (1:8) is quite general in the sense that for each λ_i , there will be a corresponding m_i . Indeed, the solution to m_i in (1:8) is called the eigenvector corresponding to the eigenvalue (or characteristic or proper value) λ_i . More on this later.

If A is pd, it is invertible and all eigenvalues in Λ will be positive. Taking the reciprocal of the square root of each and forming a corresponding diagonal matrix, say $\Lambda^{-1/2}$, it follows from (1:6) that

$$\Lambda^{-1/2} M^T A M \Lambda^{-1/2} = \Lambda^{-1/2} \Lambda \Lambda^{-1/2} = I.$$

Putting $M \Lambda^{-1/2} = Q$, this last equation may then be written as

$$Q^T A Q = I. \quad (1.9)$$

Equation (1:9) may be regarded as a definition of a pd matrix and, since M and Λ are invertible, Q is also invertible.

1.5 A Special Matrix Product

The Kronecker product of two matrices, A and B, is written $A \otimes B$ and is defined by

$$A \otimes B = [a_{ij} B].$$

If A is $(n \times m)$ and B is $(p \times q)$, then $A \otimes B$ is $(np \times mq)$. The following are properties of the Kronecker product:

$$(A \otimes B)(C \otimes D) = AC \otimes BD,$$

whereupon A and C and B and D must be conformable pairs,

$$A \otimes (B + C) = (A \otimes B) + (A \otimes C),$$

$$(A \otimes B)^{-1} = (A^{-1} \otimes B^{-1});$$

when A and B are invertible,

$$[A \otimes B]^{-1} = A^{-1} \otimes B^{-1}.$$

1.6 Scalar Functions of Square Matrices

The trace of a square matrix, A, is the sum of its diagonal elements. This is written:

$$\text{tr } A = \sum a_{ii}.$$

If A, B and C form a square matrix product, for example A is $(n \times p)$, B is $(p \times r)$ and C is $(r \times n)$ so that ABC is $(n \times n)$, then $\text{tr } ABC = \text{tr } BCA = \text{tr } CAB$. In regard

to a Kronecker product of square matrices, $\text{tr}(A \otimes B) = \text{tr} A \text{tr} B$. Quite clearly, the trace of a square matrix is a scalar.

The determinant of a square matrix may be calculated in a variety of equivalent ways. The classical expansion is based upon the (2×2) case. If

$$A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}$$

then the determinant of A, written $\det A$ or $|A|$, is given by

$$\det A = a_{11}a_{22} - a_{12}a_{21}.$$

Then if A is (3×3)

$$\det \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix} = a_{11}(a_{22}a_{33} - a_{23}a_{32}) - a_{12}(a_{21}a_{33} - a_{23}a_{31}) + a_{13}(a_{21}a_{32} - a_{22}a_{31})$$

which equals

$$= a_{11}\hat{A}_{11} - a_{12}\hat{A}_{12} + a_{13}\hat{A}_{13} \quad (1.10)$$

in which \hat{A}_{ij} equals the determinant obtained by eliminating row i and column j ($i, j = 1; 2; 3$). If now

$$A_{ij} = (-1)^{i+j} \hat{A}_{ij} \quad (1.11)$$

then (1:11) allows (1:10) to be written as:

$$\det A = \sum_j a_{1j} A_{1j}. \quad (1.12)$$

Equation (1:12) is a way of calculating the determinant of A by expansion of the first row. Equally, the same result would be obtained by expanding the second row or column and following the same rules. In the present (3 × 3) case,

$$\begin{aligned} \det A = & a_{11}(a_{22}a_{33} - a_{32}a_{23}) \\ & - a_{12}(a_{21}a_{33} - a_{31}a_{23}) \\ & + a_{13}(a_{21}a_{32} - a_{31}a_{22}) \end{aligned} \quad (1.13)$$

which reduced to (1:10). Generally, let $A = [a_{ij}]$ and $(-1)^{i+j} \hat{A}_{ij} = A_{ij}$ where \hat{A}_{ij} is the determinant formed by deleting the i'th row and the j'th column of A ($i, j = 1, 2, \dots, n$), then

$$\begin{aligned} \det A &= \sum_j a_{ij} A_{ij} \text{ for any fixed } i \\ &= \sum_i a_{ij} A_{ij} \text{ for any fixed } j. \end{aligned}$$

For any two conformable square matrices, A and B, $\det AB = \det A \times \det B$. If two columns (rows) of A are interchanged, $\det A$ merely changes sign. If a proportion of one row (column) is added to or subtracted from another, then the determinant remains unchanged. This last property is a reflection of the role played by linear dependence in determining the value of a determinant. Let the columns of A be

$a_{:1}; a_{:2}; \dots; a_{:m}$. Then these columns are said to be linearly independent if for scalars $\alpha_1; \alpha_2; \dots; \alpha_m$ the equation

$$\alpha_1 a_{:1} + \alpha_2 a_{:2} + \dots + \alpha_m a_{:m} = 0 \quad (1.14)$$

implies that $\alpha_1 = \alpha_2 = \dots = \alpha_m = 0$. If, on the contrary, equation (1.14) implies that not all α_i are zero, then $a_{:1}; a_{:2}; \dots; a_{:m}$ are said to be linearly dependent. Assuming linear dependence, say of the form $\alpha_1; \alpha_2; \alpha_3$ are non-zero, $\alpha_4; \alpha_5; \dots; \alpha_m$ are each zero, then

$$\begin{aligned} a_{:1} &= -\alpha_2 a_{:2} - \alpha_3 a_{:3} \\ &= \alpha_2 a_{:2} + \alpha_3 a_{:3} \end{aligned}$$

in an obvious notation. In this case, taking α_2 times column two from column one, and then α_3 times column three from the remainder yields:

$$\det A = \det \begin{pmatrix} 0 & a_{:2} & a_{:3} & \dots & a_{:m} \\ \vdots & \vdots & \vdots & \ddots & \vdots \end{pmatrix}$$

Expanding by the first column then yields $\det A = 0$. Thus a square matrix containing linearly dependent columns (rows) has determinant equal to zero.

1.7 The Rank of a Matrix

A square matrix whose determinant is zero is said to be singular. Non-singular matrices are invertible; singular matrices are not invertible. For any matrix, square or

rectangular, square matrices may be formed by eliminating rows and columns and then corresponding determinants may be calculated. If the largest non-vanishing determinant, formed of the elements of a matrix by eliminating rows and columns, is of order r , then the matrix is said to have rank r . If the matrix is $(n \times n)$ of rank r , then $r \leq n$ and the matrix has r linearly independent rows and columns, and the remaining $(n - r)$ rows and columns respectively are linearly dependent on them. If the matrix is rectangular of order $(n \times m)$ of rank r , then $r \leq \min\{n, m\}$. As an example, consider the $(n \times k)$ matrix X of equation (1:4) with $k < n$. This matrix can never have rank greater than k ; normally its rank is k , in which case X is said to be of full rank, i.e. its rank is as large as it can be.

1.8 Determinants as the Product of Eigenvalues

In regard to square, symmetric matrices like the matrix A in equation (1:6), notice that $\det A = \det^i M^T A M^C$ because M is an orthogonal matrix implying that $M^T = M^{-1}$ and

$$\begin{aligned} \det^i M^T A M^C &= \det M^T \det A \det M & (1.15) \\ &= \det^{-1} M \det A \det M. \end{aligned}$$

Now the determinant of any square matrix is the same as the determinant of its transpose; and the determinant of the inverse of a matrix is equal to the reciprocal of the determinant of the matrix itself. Thus (1:15) reduces to $\det A$. Since also

$\det(M^T A M) = \det A = \prod_{i=1}^n \lambda_i$, the product of the eigenvalues of A , the determinant of A is seen to be the product of its eigenvalues.

When A is square but not symmetric, the same rule (that $\det A$ is the product of its eigenvalues) holds, but for a different reason: there will exist a non-singular (but not necessarily real) matrix Q , such that

$$Q^{-1} A Q = T, \quad (1.16)$$

T being an upper triangular (not necessarily real) matrix whose diagonal elements are the eigenvalues, λ_i , of A . Since the determinant of an upper (or lower) triangular matrix is simply the product of its diagonal elements, $\det A = \det T = \prod_{i=1}^n \lambda_i$ as before.

1.9 Determinants as Areas and Volumes

An important application of determinants, which plays a vital role in multiple integration theory and hence in multivariate probability, is the evaluation of areas of parallelograms in \mathbb{R}^2 and of hypervolumes of n -dimensional parallelepipeds in \mathbb{R}^n . An example in \mathbb{R}^2 is given in Figure 1. Here, interest centres on the area of the parallelogram defined by the sum of the vectors z and w , that is the area $OABC$. The sum, $w + z$, is represented by \vec{OB} . The area of the parallelogram $OABC$ is the same as the area of the rectangle $OHGC$ defined by the sum of the vectors r and w . The vector r has been constructed so that r and w are orthogonal, or at right angles to one

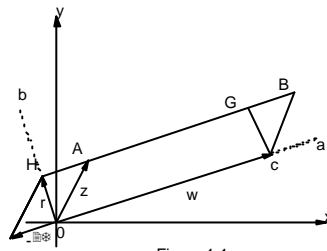


Figure 1.1

Figure 1:

another. When vectors are orthogonal, their scalar product is zero. Thus $r = z - \mu w$ for some positive fraction μ and hence

$$w \cdot (z - \mu w) = 0$$

implies

$$w \cdot z = w \cdot w \mu$$

or

$$\mu = \frac{w \cdot z}{w \cdot w}$$

Setting $\frac{w \cdot z}{w \cdot w} = P$, it follows that $P = P^T = P^2$, so P is an orthogonal projection matrix. Thus

$$r = (I_2 - P)z.$$

Recall, first, that the matrix formed by w and z , say A , has a determinant whose

value is unaffected by replacing z by $(z - w\mu) = (I_2 - P)z$:

$$\det A = \det [w : z] = \det [w : z - w\mu], \quad (1.17)$$

because a determinant is unaffected by taking away from any one column a fixed proportion of another. Recall also that while the columns of A are not orthogonal, w and $(I_2 - P)z$ are orthogonal by construction. Now the area of interest is $OHGC$ and this is obviously $\|w\| \| (I_2 - P)z \|$. Writing $B = [w : (I_2 - P)z]$,

$$|A| = \det [w : (I_2 - P)z] = |B|$$

and

$$\begin{aligned} |B| &= \det \begin{bmatrix} w & (I_2 - P)z \end{bmatrix} \\ &= \det \begin{bmatrix} w & 0 \\ 0 & z^T (I_2 - P)z \end{bmatrix} \end{aligned}$$

because $(I_2 - P) = (I_2 - P)^T = (I_2 - P)^2$. Thus

$$\begin{aligned} |A| &= \sqrt{w^T w : z^T (I_2 - P)z} \\ &= \|w\| \|(I_2 - P)z\| \\ &= \|w\| \|r\| \end{aligned}$$

Thus the area of $OABC$ equals the area of $OHGC$ equals $\|w\| \|r\|$, the determinant of the matrix comprising the co-ordinates of the vectors w and z .

The same result could have been obtained by applying equation (1:16), or (1:6) if appropriate. This would be equivalent to forming the vector r and then rotating the axes $(x; y)$ about O to become the axes $(a; b)$. When this is done, w and r respectively have co-ordinate vectors of the form

$$\begin{matrix} 2 & 3 \\ 6 & 7 \\ 4 & 5 \\ 0 & \end{matrix} a_1 \quad \text{and} \quad \begin{matrix} 2 & 3 \\ 6 & 7 \\ 4 & 5 \\ b_2 & \end{matrix}$$

in the new co-ordinate system. That these vectors would have these co-ordinates in the system $(a; b)$ is obvious, since w lies on the a -axis and r lies on the b -axis.

As an example, let the matrix A be given by:

$$A = \begin{pmatrix} 2 & 3 \\ 6 & 7 \\ 4 & 5 \end{pmatrix} \quad |A| = 14.$$

The matrix A is represented in figure 1 by the two vectors w and z :

$$w = \begin{pmatrix} 2 \\ 6 \\ 4 \end{pmatrix} \quad z = \begin{pmatrix} 3 \\ 7 \\ 5 \end{pmatrix}$$

yielding A as specified below. The vector r may be determined from $r = z + \mu w$ and

$w \cdot r = 0$. Hence $w \cdot r = 6(2 + 6\mu) + 2(3 + 2\mu) = 0$. Solving for μ yields:

$$r = \begin{pmatrix} 0.7 \\ 2.1 \\ 5 \end{pmatrix}$$

It then follows that $\|w\|^2 \|r\|^2 = (w \cdot w)(r \cdot r) = (36 + 4)(0.7^2 + 2.1^2) = 196$. Thus $\|w\| \|r\| = \sqrt{196} = 14$, as already determined for $|A|$. Taking the alternative route,

$r = (I_2 - P)z$ in which

$$P = \frac{w(w^T w)^{-1} w^T}{3}$$

$$= \begin{bmatrix} 6 & 0.9 & 0.3 \\ 0 & 7 & 5 \\ 4 & 0.3 & 0.1 \end{bmatrix}$$

Thus

$$(I_2 - P)z = \begin{bmatrix} 2 & 3 & 2 & 3 \\ 6 & 0.1 & 0.3 & 7 & 6 & 2 & 7 \\ 0 & 7 & 5 & 5 & 4 & 7 \\ 2 & 0.3 & 0.9 & 3 & 3 \end{bmatrix} z$$

$$= \begin{bmatrix} 6 & 0.7 & 7 \\ 0 & 7 & 5 \\ 4 & 2 & 1 \end{bmatrix} z$$

as before.

1.10 Vectorization

It is occasionally useful to transform a matrix into a vector. As will become evident in chapter two, the axioms of a vector space, which define what may be regarded as a vector, apply also to matrices. Therefore it is quite legitimate to regard a matrix as a vector, in the sense that a matrix is an element of a vector space.

If A is an $(n \times m)$ matrix of m columns $a_j \in \mathbb{R}^n$ and n rows $a_i \in \mathbb{R}^m$ for

$i = 1; 2; \dots; n$ and $j = 1; 2; \dots; m$, then

$$\text{vec}A = \begin{matrix} 2 & 3 \\ \begin{matrix} a_{11} \\ a_{12} \\ \vdots \\ a_{1m} \end{matrix} \\ \begin{matrix} a_{21} \\ a_{22} \\ \vdots \\ a_{2m} \end{matrix} \\ \vdots \\ \begin{matrix} a_{n1} \\ a_{n2} \\ \vdots \\ a_{nm} \end{matrix} \end{matrix} \in \mathbb{R}^{nm}$$

and

$$\text{vec}A^> = \begin{matrix} 2 & 3 \\ \begin{matrix} a_{11}^> \\ a_{12}^> \\ \vdots \\ a_{1m}^> \end{matrix} \\ \begin{matrix} a_{21}^> \\ a_{22}^> \\ \vdots \\ a_{2m}^> \end{matrix} \\ \vdots \\ \begin{matrix} a_{n1}^> \\ a_{n2}^> \\ \vdots \\ a_{nm}^> \end{matrix} \end{matrix} \in \mathbb{R}^{mn}$$

Of course, $(\text{vec}A)$ is $(1 \times nm)$ which cannot be equal to $\text{vec}A^>$ which is $(mn \times 1)$.

Also

$$\mathbf{i}_{\text{vec}A^>} = \begin{bmatrix} a_{11}^> \\ a_{12}^> \\ \vdots \\ a_{1m}^> \\ a_{21}^> \\ a_{22}^> \\ \vdots \\ a_{2m}^> \\ \vdots \\ a_{n1}^> \\ a_{n2}^> \\ \vdots \\ a_{nm}^> \end{bmatrix}$$

whereas

$$(\text{vec}A)^> = \begin{bmatrix} a_{11} \\ a_{12} \\ \vdots \\ a_{1m} \\ a_{21} \\ a_{22} \\ \vdots \\ a_{2m} \\ \vdots \\ a_{n1} \\ a_{n2} \\ \vdots \\ a_{nm} \end{bmatrix}$$

An important result for later application is the vectorization of a real matrix product. If A, B and C are respectively $(n \times m)$, $(m \times p)$ and $(p \times k)$ real matrices, their matrix product is of order $(n \times k)$ and hence $\text{vec}(ABC)$ must be an nk -tuple.

This nk -tuple is obtained from

$$\text{vec}(ABC) = \begin{bmatrix} \mathbf{f} \\ \mathbf{C}^{\mathbf{p}} \end{bmatrix} - A^{\mathbf{m}} \text{vec}B.$$

The matrix $\begin{bmatrix} \mathbf{f} \\ \mathbf{C}^{\mathbf{p}} \end{bmatrix} - A$ is of order $(kn \times pm)$ and $\text{vec}B$ is $(mp \times 1)$, thus $\text{vec}(ABC)$ must be $(kn \times 1)$ or an nk -tuple. In the case of a matrix product AB , this may be written $AB = AB I_p = I_n AB$ whereupon, applying the rule,

$$\begin{aligned} \text{vec}(AB I_p) &= [I_p - A] \text{vec}B \\ \text{vec}(I_n AB) &= \begin{bmatrix} \mathbf{f} \\ \mathbf{B}^{\mathbf{p}} \end{bmatrix} - I_n^{\mathbf{m}} \text{vec}A \end{aligned}$$

it must clearly be the case that

$$\begin{aligned} \text{vec}(AB) &= [I_p - A] \text{vec}B \\ &= \begin{bmatrix} \mathbf{f} \\ \mathbf{B}^{\mathbf{p}} \end{bmatrix} - I_n^{\mathbf{m}} \text{vec}A. \end{aligned}$$

For the matrix product $ABCD$,

$$\begin{aligned} \text{vec}(ABCD) &= \begin{bmatrix} \mathbf{i} \\ \mathbf{D}^{\mathbf{q}} \end{bmatrix} - AB^{\mathbf{c}} \text{vec}C \\ &= \begin{bmatrix} \mathbf{i} \\ \mathbf{D}^{\mathbf{q}} \end{bmatrix} - A^{\mathbf{c}} \text{vec}(BC). \end{aligned}$$

1.11 An Important Partitioned Inverse

Let A and C be symmetric, non-singular matrices of order $(m \times m)$ and $(k \times k)$, and let B be an $(m \times k)$ matrix of rank k . Using these matrices, let the $(mk \times mk)$

symmetric matrix D be defined by

$$D = \begin{bmatrix} 2 & 3 \\ 6 & 7 \\ 4 & 5 \end{bmatrix} \begin{matrix} A & B \\ B^T & C \end{matrix} \quad (1.18)$$

It is taken that all the inverses in expressions to be developed below exist.

Before proceeding, it should be made explicit that matrices that are partitioned, like D, may be treated as if the sub-matrices are elements. For example, if

$$E = \begin{bmatrix} 2 & 3 \\ 6 & 7 \\ 4 & 5 \end{bmatrix} \begin{matrix} P & Q \\ R & S \end{matrix}$$

then, provided the products exist

$$DE = \begin{bmatrix} 2 & 3 \\ 6 & 7 \\ 4 & 5 \end{bmatrix} \begin{matrix} A & B \\ B^T & C \end{matrix} \begin{bmatrix} 2 & 3 \\ 6 & 7 \\ 4 & 5 \end{bmatrix} \begin{matrix} P & Q \\ R & S \end{matrix} = \begin{bmatrix} 2 & 3 \\ 6 & 7 \\ 4 & 5 \end{bmatrix} \begin{matrix} AP + BR & AQ + BS \\ B^T P + CR & B^T Q + CS \end{matrix}$$

Similarly

$$E^T = \begin{bmatrix} 2 & 3 \\ 6 & 7 \\ 4 & 5 \end{bmatrix} \begin{matrix} P^T & R^T \\ Q^T & S^T \end{matrix}$$

Returning to D of (1:18), a formula for D^{-1} is

$$D^{-1} = \begin{bmatrix} 2 & 3 \\ 6 & 7 \\ 4 & 5 \end{bmatrix} \begin{matrix} A_i & BC_i^{-1} B^T C_i^{-1} \\ C_i & B^T A_i^{-1} B^T C_i^{-1} \end{matrix} \quad (1.19)$$

In (1:19) it should be noticed that, since D is symmetric, so is D^{-1} and hence

$$A_i & BC_i^{-1} B^T C_i^{-1} = C_i & B^T A_i^{-1} B^T C_i^{-1}$$

or

$${}^i A_j B C_i^{-1} B^{-1} \mathbb{C}_i^{-1} B C_i^{-1} = A_i^{-1} B {}^i C_j B^{-1} A_i^{-1} B \mathbb{C}_i^{-1}. \quad (1.20)$$

That D_i^{-1} in (1:19) is indeed the inverse of D is easily checked by post-multiplying D_i^{-1} by D . This yields

$$\begin{aligned} & {}^i A_j B C_i^{-1} B^{-1} \mathbb{C}_i^{-1} A_j {}^i A_j B C_i^{-1} B^{-1} \mathbb{C}_i^{-1} B C_i^{-1} B^{-1} \\ &= {}^i A_j B C_i^{-1} B^{-1} \mathbb{C}_i^{-1} {}^i A_j B C_i^{-1} B^{-1} \mathbb{C}_i^{-1} \\ &= I_m \end{aligned} \quad (1.21)$$

$${}^i A_j B C_i^{-1} B^{-1} \mathbb{C}_i^{-1} B_j {}^i A_j B C_i^{-1} B^{-1} \mathbb{C}_i^{-1} B C_i^{-1} C = 0 \quad (1.22)$$

$${}_j {}^i C_j B^{-1} A_i^{-1} B \mathbb{C}_i^{-1} B^{-1} A_i^{-1} A + {}^i C_j B^{-1} A_i^{-1} B \mathbb{C}_i^{-1} B^{-1} = 0 \quad (1.23)$$

$$\begin{aligned} & {}_j {}^i C_j B^{-1} A_i^{-1} B \mathbb{C}_i^{-1} B^{-1} A_i^{-1} B + {}^i C_j B^{-1} A_i^{-1} B \mathbb{C}_i^{-1} C \\ &= {}^i C_j B^{-1} A_i^{-1} B \mathbb{C}_i^{-1} {}^i C_j B^{-1} A_i^{-1} B \\ &= I_k \end{aligned} \quad (1.24)$$

1.12 Useful Inverses and Projection Matrices

Re-writing equation (1:21) using equation (1:20) yields

$${}^i A_j B C_i^{-1} B^{-1} \mathbb{C}_i^{-1} A_j A_i^{-1} B {}^i C_j B^{-1} A_i^{-1} B \mathbb{C}_i^{-1} B^{-1} = I_m:$$

Post-multiplying throughout by A^{-1} , there results

$$A^{-1}B^{-1}C^{-1}B^{-1} = A^{-1} + A^{-1}B^{-1}C^{-1}B^{-1}A^{-1}. \quad (1.25)$$

Equation (1.25) finds application in the random-coefficients model associated with the linear model (1.4), i.e. model (1.4) augmented by the relation

$$y = A\theta + \hat{A} \quad (1.26)$$

in which A is a $(k \times m) (m \times k)$ known matrix of rank m , θ is an m -tuple of coefficients and \hat{A} is a $(k \times 1)$ vector of errors, independent of y in (1.4). Associated with \hat{A} is an $(m \times n)$ pd matrix Σ . Setting

$$\begin{bmatrix} y \\ X \end{bmatrix} = \begin{bmatrix} A & B \\ C & D \end{bmatrix} \begin{bmatrix} \theta \\ \alpha \end{bmatrix} + \begin{bmatrix} \hat{A} \\ \epsilon \end{bmatrix}$$

and applying (1.25),

$$\begin{bmatrix} I & X \Sigma^{-1} X' \\ X \Sigma^{-1} X' & X \Sigma^{-1} X' + D^{-1} \end{bmatrix}^{-1} \begin{bmatrix} y \\ X \alpha \end{bmatrix} = \begin{bmatrix} I & X \Sigma^{-1} X' \\ X \Sigma^{-1} X' & X \Sigma^{-1} X' + D^{-1} \end{bmatrix}^{-1} \begin{bmatrix} y \\ X \alpha \end{bmatrix} + \begin{bmatrix} I & X \Sigma^{-1} X' \\ X \Sigma^{-1} X' & X \Sigma^{-1} X' + D^{-1} \end{bmatrix}^{-1} \begin{bmatrix} \hat{A} \\ \epsilon \end{bmatrix}. \quad (1.27)$$

The matrix on the left-hand side of (1.27) is the inverse of the matrix associated with the composite error, $y + X\hat{A}$, obtained by substituting (1.26) into (1.4).

Two applications of equation (1.19) arise in least squares estimation theory applied to model (1.4). First, it is often useful to have available a particular form of the pd

With

$$F_1 = X_1 (X_1' M_2 X_1)^{-1} X_1' M_2$$

$$F_2 = X_2 (X_2' M_1 X_2)^{-1} X_2' M_1$$

$$P = X (X' X)^{-1} X'$$

equation (1:28) becomes

$$P = F_1 + F_2. \tag{1.29}$$

Here, $P = P' = P^2$, $F_1 = F_1^2$ and $F_2 = F_2^2$ but neither F_1 nor F_2 is symmetric. While P is referred to as an orthogonal projection, F_1 and F_2 are referred to as oblique projections. In model (1:4)

$$\begin{aligned} y &= X\beta + \epsilon \\ &= X_1\beta_1 + X_2\beta_2 + \epsilon \end{aligned}$$

where $\beta = \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}$. Py is the least squares estimator of $X\beta$, $F_1 y$ and $F_2 y$ are the corresponding least squares estimators of $X_1\beta_1$ and $X_2\beta_2$ respectively.

2 Vector Spaces

2.1 Introduction

The aim of chapter one was to review some of the theory of matrices from the classical viewpoint of manipulative algebra. Along with a number of specific definitions and properties of matrices, three points were noted:

1. A vector in \mathbb{R}^n may be regarded as a special case of an $(n \times m)$ matrix simply by setting $m = 1$.
2. Any $(n \times m)$ matrix, A , defined over the field \mathbb{R} , may be transformed into a $(nm \times 1)$ vector by the operation of vectorization. Thus, $\text{vec}A$ is an $(nm \times 1)$ vector.
3. Let M be the set of all $(n \times m)$ matrices defined over \mathbb{R} and let $\alpha \in \mathbb{R}$. Then for any pair, A and B in M , $A + B \in M$ and $\alpha A \in M$.

Point 1 suggests that matrices are more general objects than vectors, but this seems to be contradicted by point 2, since any matrix may be written as a vector. As will become apparent, matrices and vectors are merely special examples of a class of objects which obey the property recognized in point 3. This property is fundamental to the definition of a linear or vector space.

It is the purpose of this chapter to introduce the theory of linear or vector spaces and to explain, in reasonable detail, a particular type of vector space with which there will be especial concern. This is a Euclidean space and such a space is the setting for the linear model outlined in equations (1:1) to (1:4). A Euclidean space is a vector space endowed with a means of measuring the length of vectors and the angles between them. The means of measuring lengths and angles is the scalar product.

The first step is to characterize point 3 above in the form of a definition.

Definition 1 A set V of elements $x; y; z; \dots$ is said to form a vector space over the field F if:

- (a) For every two elements x and y in V , there is associated an element z in V called the sum of x and y , written $x + y = z$.
- (b) For every x in V and every α in F , there is an element αx in V called the product of x and α .

These operators must obey the following axioms.

- (1) Vector operations:

Commutativity: $x + y = y + x$

Associativity: $(x + y) + z = x + (y + z)$

Zero: V contains an element, 0 ,
such that $x + 0 = x$ $\forall x$ in V .

Negativity: For every x in V , there exists an
element $(-x)$ such that $x + (-x) = 0$.

(2) Scalar Multiplication:

of a vector: $1x = x$

$$(\alpha \beta)x = \alpha(\beta x) \quad \alpha, \beta \in F;$$

of a sum of vectors $(\alpha + \beta)x = \alpha x + \beta x$

or scalars $\alpha(x + y) = \alpha x + \alpha y$. ■

The theory of vector spaces is much more general than is suggested by the examples of \mathbb{R}^n . Two examples illustrate the point.

Example 1 The set of all polynomials of degree not exceeding some natural number, n , constitutes a vector space, whereas the set of polynomials of degree exactly n does not. To see this, the following example suffices:

$$i(t^3 + t) + i(t^3 + t) = 2t$$

This is not a polynomial of degree $n = 3$, but it is a polynomial of degree not exceeding $n = 3$. ■

Example 2 The set of all real matrices of order $(n \times m)$, $A; B; C; \dots$ constitutes a vector space since $A + B = C$ lies in the same set and so does λA where $\lambda \in \mathbb{R}$. For example

$$\begin{matrix} 2 & 3 & 2 & 3 & 2 & 3 \\ \begin{matrix} 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 4 \end{matrix} & \begin{matrix} 1 & 4 \\ 2 & 1 \\ 3 & 3 \end{matrix} & \begin{matrix} 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 5 \end{matrix} & + & \begin{matrix} 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 4 \end{matrix} & \begin{matrix} 6 & 1 \\ 2 & 1 \\ 6 & 1 \end{matrix} & \begin{matrix} 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 5 \end{matrix} & = & \begin{matrix} 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 4 \end{matrix} & \begin{matrix} 7 & 5 \\ 4 & 2 \\ 9 & 4 \end{matrix} & \begin{matrix} 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 5 \end{matrix} \end{matrix};$$

and

$$\begin{matrix} 2 & 3 & 2 & 3 \\ \begin{matrix} 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 4 \end{matrix} & \begin{matrix} 1 & 4 \\ 2 & 1 \\ 3 & 3 \end{matrix} & \begin{matrix} 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 5 \end{matrix} & = & \begin{matrix} 6 \\ 6 \\ 6 \\ 6 \\ 6 \\ 4 \end{matrix} & \begin{matrix} 2 & 8 \\ 4 & 2 \\ 6 & 6 \end{matrix} & \begin{matrix} 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 5 \end{matrix} \end{matrix}.$$

The answer in each case is a (3×2) matrix of real numbers. This example merely illustrates point 3 in the opening paragraph of this sections. ■

The elements of a vector space are generally called vectors. There is some confusion caused by this terminology in examples 1 and 2. Thus a matrix is an element of a vector space, namely the set of all real matrices of the same order. This implies that a matrix is a vector; and, in a formal sense, it is. Yet there will be no serious confusion so long as it is recognized that a matrix will be called a matrix and that such a matrix is an element of a vector space. As will become apparent below, vector spaces may be distinguished by their dimensionality. In the case of the vectors from \mathbb{R}^n , \mathbb{R}^n is said to be an n-dimensional vector space. In the case of $(n \times m)$ matrices of

real elements, these are elements of \mathbb{R}^{nm} , notwithstanding the fact that such matrices may be thought of as m vectors drawn from \mathbb{R}^n . Of course, if A is an $(n \times m)$ matrix, then $\text{vec}A \in \mathbb{R}^{nm}$. In respect of example 1, a terminology may be used corresponding to that of example 2. Thus polynomials of degree not exceeding n may be recognized as polynomials while also appreciating that such polynomials are elements of a vector space.

Vectors defined over the field \mathbb{R} are called real vectors; the spaces in which these vectors lie are called real vector spaces. In econometrics, real vector spaces are of great importance, especially real vector spaces endowed with a scalar product. Before introducing a scalar product, three important features of general vector spaces are introduced. These are dimensionality, basis and co-ordinates.

2.2 Dimensionality, Basis and Co-ordinates

If $x; y; z; \dots; w$ are vectors in a vector space and $\alpha; \beta; \gamma; \dots; \mu$ is a corresponding set of scalars, then

$$\alpha x + \beta y + \gamma z + \dots + \mu w$$

is said to be a linear combination of the vectors $x; y; z; \dots; w$ with coefficients $\alpha; \beta; \gamma; \dots; \mu$.

Definition 2 Let V be a vector space. The vectors $x; y; z; \dots; w$ in V are said to be linearly dependent if there exists a corresponding set of scalars $\alpha; \beta; \gamma; \dots; \mu$, not

all zero, such that the linear combination

$$\alpha x + \beta y + \gamma z + \dots + \mu w = 0. \quad (2.1)$$

The vectors $x; y; z; \dots; w$ are said to be linearly independent if the equality

$$\alpha x + \beta y + \gamma z + \dots + \mu w = 0$$

implies that $\alpha = \beta = \gamma = \dots = \mu = 0$. ■

Let the vectors $x; y; z; \dots; w$ be linearly dependent. Then there exists scalars such that equation (2:1) exists with at least one coefficient non-zero. Let $\alpha \neq 0$. Then equation (2:1) may be re-written

$$x = -\frac{\beta}{\alpha}y - \frac{\gamma}{\alpha}z - \dots - \frac{\mu}{\alpha}w. \quad (2.2)$$

Putting $-\frac{\beta}{\alpha} = \lambda; -\frac{\gamma}{\alpha} = 1; \dots; -\frac{\mu}{\alpha} = \lambda$, then

$$x = \lambda y + z + \dots + \lambda w \quad (2.3)$$

and x may be expressed as a linear combination of the remaining vectors. Going the other way round, if equation (2:3) holds, then equation (2:1) will also hold by setting $\alpha = \lambda; \beta = \lambda; \gamma = 1; \dots; \mu = \lambda$. Thus if a set of vectors is linearly dependent, then one vector can be expressed as a linear combination of the others and if one vector may be expressed as a linear combination of the others, then those vectors are linearly dependent.

Definition 3 A vector space V is said to be n -dimensional if it contains n linearly independent vectors and any $(n + 1)$ vectors in V are linearly dependent. ■

The idea of a basis follows directly from the idea of dimensionality.

Definition 4 Any set of n linearly independent vectors in an n -dimensional space V is called a basis of V . ■

From the definition of a basis comes the idea of co-ordinates. Co-ordinates are a product of the following theorem.

Theorem 1 Every vector in an n -dimensional vector space V can be represented as a linear combination of basis vectors.

Proof. Let $e_1; e_2; \dots; e_n$ be a basis of V . For any arbitrary x in V , $x; e_1; e_2; \dots; e_n$ represents a linearly dependent set of vectors. Therefore there will exist scalars $\alpha_0; \alpha_1; \dots; \alpha_n$, not all zero, such that

$$\alpha_0 x + \alpha_1 e_1 + \alpha_2 e_2 + \dots + \alpha_n e_n = 0.$$

Now $\alpha_0 \neq 0$ for if it were zero, then the basis vectors would be linearly dependent which they cannot be. Thus $\alpha_0 \neq 0$ and

$$x = \alpha_0^{-1} \alpha_1 e_1 + \alpha_0^{-1} \alpha_2 e_2 + \dots + \alpha_0^{-1} \alpha_n e_n.$$

This proves that x may be expressed as a linear combination of the basis vectors. To prove uniqueness, let there be two linear combinations of the same basis vectors which equate with x :

$$x = \alpha_1 e_1 + \alpha_2 e_2 + \dots + \alpha_n e_n$$

and

$$x = \beta_1 e_1 + \beta_2 e_2 + \dots + \beta_n e_n.$$

Subtracting, it follows that

$$0 = (\alpha_1 - \beta_1) e_1 + (\alpha_2 - \beta_2) e_2 + \dots + (\alpha_n - \beta_n) e_n.$$

Since $e_1; e_2; \dots; e_n$ are linearly independent vectors, it follows that $(\alpha_1 - \beta_1) = (\alpha_2 - \beta_2) = \dots = (\alpha_n - \beta_n) = 0$, implying that $\alpha_1 = \beta_1; \alpha_2 = \beta_2; \dots; \alpha_n = \beta_n$. Thus, the representation is unique. ■

In theorem 1, it may be taken that, relative to the basis $e_1; e_2; \dots; e_n$, x is determined by the unique set of coefficients $\alpha_1; \alpha_2; \dots; \alpha_n$. These coefficients are called the co-ordinates of x relative to the basis $e_1; e_2; \dots; e_n$. This basis may be defined as an $(n \times n)$ matrix since each e_i is an $(n \times 1)$ vector. Let this matrix, called a basis matrix, be defined as

$$[e_1; e_2; \dots; e_n] = E,$$

whereupon, writing $B = \begin{bmatrix} b_{11} & b_{12} & \dots & b_{1n} \\ b_{21} & b_{22} & \dots & b_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ b_{n1} & b_{n2} & \dots & b_{nn} \end{bmatrix}$, since B is another $(n \times n)$ matrix and

$$G = EB: \quad (2.6)$$

It follows from the representation of x that

$$x = E\xi = G\mu = EB\mu \quad (2.7)$$

and immediately $\xi = B\mu$: B must be invertible because equation (2.6) indicates that the rank of G (which must be n since it has n linearly independent columns, by definition) is equal to the rank of EB . But E also has rank n , since it is a basis matrix. Hence, B also has rank n . Therefore, $\mu = B^{-1}\xi$. If the matrix B (whose columns represent the co-ordinates of the basis matrix G relative to the basis matrix E) is known, then $G = EB$ and the co-ordinates, μ_i , of x relative to the basis G are also known from the co-ordinates, ξ_i , of x relative to the basis E via $\mu = B^{-1}\xi$.

Finally note that if $x = E\xi$ and $y = E\mu$, then $x + y = E\xi + E\mu = E(\xi + \mu)$, i.e. given the same basis, addition of vectors corresponds to the addition of their co-ordinates.

2.3 Euclidean Space

A Euclidean space is a real vector space endowed with a scalar product, i.e. a mapping $\mathbb{R}^n \rightarrow \mathbb{R}$. The purpose of a scalar (or inner) product is to provide the means to measure length and angle in the space. Vectors have length and direction and these ideas may be given numerical expression once the scalar product is defined.

The natural scalar product for any two vectors x and y in \mathbb{R}^n is

$$x \cdot y = \sum_i x_i y_i \quad (2.8)$$

for all $i = 1; 2; \dots; n$, the x_i and y_i being the co-ordinates of x and y relative to the basis $[e_1; e_2; \dots; e_n] = I_n$. Notice the following properties of $x \cdot y$ for every pair of vectors in \mathbb{R}^n :

- i) $x \cdot y = y \cdot x$;
- ii) $(\lambda x) \cdot y = \lambda (x \cdot y)$; $\lambda \in \mathbb{R}$;
- iii) $(x + z) \cdot y = x \cdot y + z \cdot y$;
- iv) $x \cdot x \geq 0$ with equality only if $x = 0$.

The scalar product $x \cdot y$ is just a special case of $x \cdot Ay$ where $A = I_n$. A more general scalar product would be

$$x \cdot Ay = \sum_i \sum_j x_i a_{ij} y_j$$

for $i, j = 1; 2; \dots; n$, subject to the following conditions corresponding to (i)-(iv) immediately above for every pair of vectors in \mathbb{R}^n :

- i) $x \cdot Ay = y \cdot Ax$;
- ii) $(\lambda x) \cdot Ay = \lambda (x \cdot Ay)$; $\lambda \in \mathbb{R}$;
- iii) $(x + z) \cdot Ay = x \cdot Ay + z \cdot Ay$;
- iv) $x \cdot Ax \geq 0$ with equality only if $x = 0$.

Condition (i) implies that $y^T A^T x = y^T Ax$ and hence that A is a symmetric matrix. Condition (iv) implies that A is a pd matrix. This last condition is required to measure length as $\sqrt{x^T Ax}$ and so $x^T Ax$ must be a pd quadratic form for a real square root to exist. In practice, A may be chosen as a matter of convenience.

Because there are many possibilities for defining a scalar product, a more general notation is required. Usually, this takes the form $(x; y)$ or $hx; yi$; sometimes $[x; y]$. To indicate a particular product, for example the one that takes the form $x^T Ax$ above, it would be usual to define this relative to $x^T y$ in the following notation: $(; ;)$ is the natural scalar and $h; ; i = (; ; A;)$. A general definition of a scalar product is now set out as

Definition 5 If for every pair of vectors $x; y$ in \mathbb{R}^n there is associated a real number $(x; y)$ such that

- i) $(x; y) = (y; x)$
- ii) $(\alpha x; y) = \alpha (x; y); \alpha \in \mathbb{R}$
- iii) $(x_1 + x_2; y) = (x_1; y) + (x_2; y)$
- iv) $(x; x) \geq 0$ with equality only when $x = 0$,

then $(x; y)$ is called the scalar product defined on \mathbb{R}^n and \mathbb{R}^n is then referred to as an n -dimensional Euclidean space.

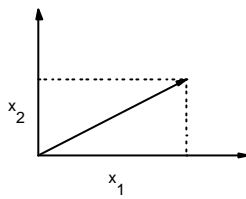


Figure 2.1

Figure 2:

Definition 6 Let x be a vector in Euclidean space in which the scalar product $(; ;)$ is defined. By the length of x is meant the positive number $\|x\| = \sqrt{(x; x)}$, unless $x = 0$, whereupon $\|x\| = 0$.

Example 3 In 2-dimensional Euclidean space, the vector x has co-ordinates x_1 and x_2 as indicated in Figure 2.1. The scalar product defined on this space is $(x; y) = x^T y$.

Thus

$$x = \begin{pmatrix} 6 \\ 4 \end{pmatrix} \begin{matrix} x_1 \\ x_2 \end{matrix}$$

and

$$\|x\|^2 = x^T x = x_1^2 + x_2^2.$$

It follows that $\|x\| = \sqrt{(x_1^2 + x_2^2)}$, which is the same length as according to Pythagoras' Theorem. ■

Let x and y be two non-zero vectors in \mathbb{R}^n on which the scalar product $(; ;)$ is

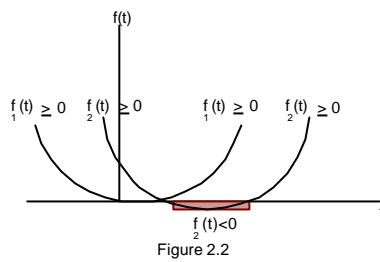


Figure 3:

defined. Let t be a real number. Axiom 4 of definition 5 requires

$$(x \dot{;} ty; x \dot{;} ty) \dot{;} 0$$

with equality only if $x = ty$. Consequently

$$(x; x) \dot{;} t(x; y) \dot{;} t(y; x) + t^2(y; y) \dot{;} 0.$$

From axiom 1 of definition 5, this last expression becomes

$$(y; y)t^2 \dot{;} 2(x; y)t + (x; x) = f(t) \dot{;} 0. \tag{2.9}$$

Thus the value of t must be such that $f(t)$, a quadratic form in t , is always positive or zero, and never negative; thus $f(t)$ may be like $f_1(t)$ or $f_3(t)$ in figure 2, but not like $f_2(t)$. The quadratic $f(t)$, therefore, can never cross the t -axis (diagram 2.2) i.e. $f(t)$ can never have two distinct real roots. If $f(t)$ cannot have two distinct real roots, then in equation (2:9) it must have either two coincident roots or imaginary roots. Hence

$$4(x; y)^2 \dot{;} 4(y; y)(x; x) \tag{2.10}$$

i.e. in $ax^2 + bx + c$; $b^2 = 4ac$ when coincident roots emerge, or $b^2 < 4ac$ and hence the roots are imaginary. Thus from equation (2:10)

$$(x; y)^2 \leq (x; x)(y; y). \quad (2.11)$$

Inequality (2:11) is known as the Cauchy-Schwartz Inequality. It implies

$$| \cos \mu | \leq \frac{(x; y)}{\|x\| \|y\|} \leq 1. \quad (2.12)$$

For this reason, it is possible to state

Definition 7 By the angle μ between two non-zero vectors in n -dimensional Euclidean space is meant

$$\mu = \cos^{-1} \left(\frac{(x; y)}{\|x\| \|y\|} \right). \quad \blacksquare$$

The bounds expressed by equation (2:12) are the bounds of the cosine function and this is the justification for using Definition 6. It is helpful, nevertheless, to explore the geometry in \mathbb{R}^2 that lies behind Definition 6.

Example 4 Consider x and y in \mathbb{R}^n on which the natural scalar product $x \cdot y = (x; y)$ is defined. The notation is described in Figure 2.3.

The basis in the space is

$$[e_1 \ e_2] = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = I_2.$$

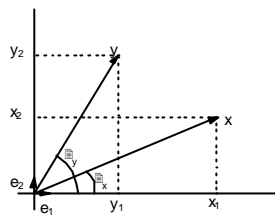


Figure 2.3

Figure 4:

Relative to the basis

$$x = \begin{pmatrix} 6 \\ 4 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 6 \\ 4 \end{pmatrix}$$

and

$$y = \begin{pmatrix} 7 \\ 5 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 7 \\ 5 \end{pmatrix}$$

The lengths of x and y are quite natural in view of Pythagoras' Theorem

$$\|x\|^2 = x_1^2 + x_2^2, \quad \|x\| = \sqrt{x_1^2 + x_2^2},$$

$$\|y\|^2 = y_1^2 + y_2^2, \quad \|y\| = \sqrt{y_1^2 + y_2^2}.$$

Moreover

$$\sin \mu_x = \frac{x_2}{\|x\|}; \cos \mu_x = \frac{x_1}{\|x\|};$$

$$\sin \mu_y = \frac{y_2}{\|y\|}; \cos \mu_y = \frac{y_1}{\|y\|}.$$

Now interest centres on the angle between x and y , that is $\mu = \mu_y - \mu_x$ and $\cos(\mu_y - \mu_x)$

$\mu_x) = \cos \mu_y \cos \mu_x + \sin \mu_y \sin \mu_x$: Therefore

$$\begin{aligned} \cos \mu &= \frac{x_1 y_1 + x_2 y_2}{\|x\| \|y\|} \\ &= \frac{x \cdot y}{\|x\| \|y\|}. \end{aligned}$$

It follows that the Cauchy-Schwartz inequality is vindicated in this case and that

$$x \cdot y = \|x\| \|y\| \cos \mu.$$

When $\mu = \frac{\pi}{2} = 90^\circ$, $\cos \mu = 0$ and hence $x \cdot y = 0$; x and y are said to be orthogonal.

If $\mu = 0$ $\cos \mu = 1$ and if $\mu = \pi = 180^\circ$, $\cos \mu = -1$: In these two cases x and y are co-linear. For example, when $\mu = 0$

$$\begin{aligned} x \cdot y &= \|x\| \|y\| \\ &= \|x\|^2 \frac{\|y\|}{\|x\|} \\ &= x \cdot x \frac{\|y\|}{\|x\|} \end{aligned}$$

Setting $\lambda = \frac{\|y\|}{\|x\|}$, $\lambda \in \mathbb{R}$ and $x \cdot y = x \cdot (\lambda x) = \lambda (x \cdot x) = \lambda \|x\|^2$: If $\mu = \pi$; then $\lambda < 0$ and y would be opposite in direction to x . ■

Definition 8 If for any pair of vectors x and y in Euclidean space

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

then x and y are said to be orthogonal.

Definition 8 implies that $(x; y) = 0$ and is in fact a straightforward statement of Pythagoras' Theorem. Notice also that the statement holds for any scalar product.

If for example $h; \cdot; i = (\cdot; A; \cdot)$, then

$$\|x\|^2 = \overline{(x; Ax)}.$$

Moreover, if $x; Ay = 0$ then x and y are orthogonal relative to the scalar product $h; \cdot; i$. Thus relative to $h; \cdot; i$

$$\begin{aligned} \|x + y\|^2 &= (x + y); A(x + y) \\ &= x; Ax + x; Ay + y; Ax + y; Ay \\ &= x; Ax + 2x; Ay + y; Ay \\ &= \|x\|^2 + \|y\|^2 + 2x; Ay. \end{aligned}$$

Thus $\|x + y\|^2 = \|x\|^2 + \|y\|^2$ only when $x; Ay = 0$, i.e.

$$h; x; y; i = (x; Ay) = 0.$$

Definition 9 In n -dimensional Euclidean space, the vectors $x_1; x_2; \dots; x_n$ are said to be pairwise orthogonal if

$$\|x_1 + x_2 + \dots + x_n\|^2 = \|x_1\|^2 + \|x_2\|^2 + \dots + \|x_n\|^2.$$

If the vectors $x_1; x_2; \dots; x_n$ constitute a set of non-zero, pairwise orthogonal vectors in n -dimensional Euclidean space, these vectors constitute a basis of the space, called an orthogonal basis. ■

Notice that any vector x in Euclidean space may be normalized to have length one, by the simple device of scalar multiplication by $\frac{1}{\|x\|}$. Thus if $z = \frac{x}{\|x\|}$ then

$$\begin{aligned} \|z\|^2 &= (z; z) \\ &= \frac{(x; x)}{\|x\|^2} \\ &= 1. \end{aligned}$$

Definition 10 The non-zero vectors $e_1; e_2; \dots; e_n$ of an n -dimensional Euclidean space are said to form an orthonormal basis if they are pairwise orthogonal and if, in addition, each has unit length; that is if $(e_i; e_j) = \delta_{ij}$ where δ_{ij} is Kronecker's delta $i; j = 1; 2; \dots; n$ (i.e. $\delta_{ij} = 0$ for $i \neq j$ and $\delta_{ij} = 1$ for $i = j$). ■

Before moving on to discuss orthogonal bases, one theorem and one definition will be introduced.

Theorem 2 In Euclidean space, let x and y be two vectors. Then $\|x+y\|^2 = \|x\|^2 + \|y\|^2 + 2(x; y)$.

Proof.

$$\begin{aligned} \|x+y\|^2 &= (x+y; x+y) \\ &= (x; x) + 2(x; y) + (y; y). \end{aligned}$$

Now $(x; y) = \|x\|\|y\|\cos\theta$: Hence

$$\begin{aligned} \|x+y\|^2 &= (x+y; x+y) \\ &= (x; x) + 2\|x\|\|y\|\cos\theta + (y; y). \end{aligned}$$

But

$$\begin{aligned} (x+y; x+y) + 2\|x\|\|y\| + (y; y) &= \|x\|^2 + 2\|x\|\|y\| + \|y\|^2 \\ &= (\|x\| + \|y\|)^2. \end{aligned}$$

Thus

$$\|x+y\| = \|x\| + \|y\|. \blacksquare$$

Definition 11 The distance between the vectors x and y in Euclidean space is defined as $\|x-y\|$. \blacksquare

Every n -dimensional Euclidean vector space contains orthogonal bases, and hence orthonormal bases. Let e_i ($i = 1; 2; \dots; n$) be an $(n \times 1)$ vector with zeros everywhere except in the i 'th position which is unity. Then $e_1; e_2; \dots; e_n$ constitutes an orthonormal basis:

$$[e_1; e_2; \dots; e_n] = I_n.$$

Moreover, $\sum_{i=1}^n e_i \cdot e_i$ which has one in every position; e is called the equiangular line in \mathbb{R}^n . Every n -tuple x in \mathbb{R}^n may be written $x = I_n x$, whereupon x_i is seen to be the i 'th co-ordinate of x relative to the orthonormal basis matrix I_n . Finally, consider

$$\begin{aligned} (e_i; x) &= e_i^T x \\ &= x_i. \end{aligned}$$

Also $\langle e_i, e_i \rangle = 1$. Hence

$$\langle e_i, e_i \rangle = \begin{pmatrix} 0 & \dots & 0 & 1 & 0 & \dots & 0 \\ 0 & \dots & 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & 0 & 0 & \dots & 0 \\ 0 & \dots & 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & 0 & 0 & \dots & 0 \\ 0 & \dots & 0 & 0 & 0 & \dots & 0 \end{pmatrix} x = \begin{pmatrix} x_i \\ \vdots \\ x_i \\ \vdots \\ 0 \\ \vdots \\ 0 \end{pmatrix}.$$

$\langle e_i, e_i \rangle = 1$ is an example of an orthogonal projection. Before turning our attention to these, it is important to introduce the concept of a subspace.

2.4 Subspaces

A subset of vectors in a vector space is called a subspace if it forms a vector space according to Definition 1. Consider the basis matrix I_n of \mathbb{R}^n . The first k columns of I_n form an $(n \times k)$ matrix E_k and the remaining $(n - k)$ columns an $n \times (n - k)$ matrix E_{n-k} . Clearly E_k has rank k and E_{n-k} has rank $(n - k)$. E_k may be regarded as a basis of a k -dimensional subspace and E_{n-k} as a basis for an $(n - k)$ -dimensional subspace. Let the first subspace be denoted L and the second as L^\perp . It is clear that any vector in L is orthogonal to any vector in L^\perp because

$$E_k^T E_{n-k} = 0.$$

This implies that for all $\mu \in \mathbb{R}^k$ and all $v \in \mathbb{R}^{n-k}$:

$$\begin{aligned} (E_k \mu)' E_{n-k} v &= \mu' E_k' E_{n-k} v \\ &= 0. \end{aligned}$$

Thus any vector in L , written $E_k \mu$, is orthogonal to all vectors in L^\perp , written $E_{n-k} v$, indeed, the superscript \perp denotes orthogonal complement since $L + L^\perp = \mathbb{R}^n$, by definition, and both L and L^\perp intersect, by construction, only at the origin $= \{0\}$. In this particular case, there is only intersection at the origin; it is then usual to write $L \oplus L^\perp = \mathbb{R}^n$ where \oplus denotes direct sum; however, it is not necessary for the subspaces to be orthogonal. Thus two subspaces L and M might not be orthogonal yet $L \setminus M = \{0\}$; in which case we could write $L \oplus M = V$, say, and M would be the complement of L in V .

Definition 12 Let L be a subspace of an n -dimensional Euclidean space \mathbb{R}^n . A vector in \mathbb{R}^n is orthogonal to L if it is orthogonal to every vector in L . ■

Two orthogonal subspaces of great interest in econometrics may now be introduced. Consider the linear model introduced earlier in equation (1:4)

$$y = X\beta + \epsilon,$$

where X has k linearly independent columns and n rows. The first subspace in \mathbb{R}^n of interest is the range (or span) of X which is the set of all vectors that are linear

combinations of the columns of X :

$$R[X] = \{x \in \mathbb{R}^n : x = X_{:,j}, j \in \{1, \dots, k\}\}.$$

Clearly, $R[X]$ is a subspace. With $x_1 = X_{:,1}$ and $x_2 = X_{:,2}$, $x_1 + x_2 = X_{:,1+2} = X_{:,j_1+j_2}$. Hence, $X_{:,j} \in R[X]$. Also if α is a scalar and $x \in R[X]$, $\alpha x = X_{:,j} \alpha = X_{:,j} \mu$; $\mu \in \mathbb{R}^k$, hence $\alpha x \in R[X]$. Since X has rank k , it is always possible to form k linearly independent linear combinations of its columns such that

$$XM = W$$

where M is a $(k \times k)$, non-singular matrix. Then the rank of W = the rank of X and $R[W] = \{w \in \mathbb{R}^n : w = W\mu; \mu \in \mathbb{R}^k\}$. But $W\mu = XM\mu = X_{:,j}$; say, and hence $R[X] = R[W]$. Clearly the dimension of $R[X]$ is k .

The other subspace of interest is the null-space or kernel of $X^>$:

$$N(X^>) = \{z \in \mathbb{R}^n : X^>z = 0\}.$$

Clearly, $N(X^>)$ is a subspace, for if z_1 and $z_2 \in N(X^>)$ then $X^>(z_1 + z_2) = X^>z_1 + X^>z_2 = 0 + 0 = 0$. $z_1 + z_2 = z$, whereupon $z \in N(X^>)$. Moreover, if $X^>z = 0$, then $X^>(z) = \sum_j X^>z_j = 0$, $\sum_j z_j = 0$, $\sum_j z_j \in R[X]$.

Theorem 3 For the $(n \times k)$ real matrix X of rank k , $N(X^>) = R[X]^\perp$.

Proof. Let $z \in N(X^>)$, then $X^>z = 0$ which implies, for all $j \in \{1, \dots, k\}$, $\sum_j X^>z_j = 0$ or $(X_{:,j})^>z = 0$ and $z \in R[X]^\perp$. Let $z \in R[X]^\perp$: Then $X^>z = 0$ and $z \in N(X^>)$. ■

The dimension of $R[X] = \text{rank } X = k$. The dimension of $N[X]$ = the number of columns of X ; $\text{rank } X = (n - k)$. The dimension of $N[A]$, where A is any matrix, is called the nullity of A . Moreover, in view of Theorem 3, it is clearly the case that, if A has order $(n \times k)$ and rank p ,

$$\dim R[A] + \dim N[A] = p + (n - p) = n$$

and this is equivalent to saying

$$\text{rank } A + \text{nullity of } A = n.$$

Of course, the nullity of $X = k$; $\text{rank } X = k$; $k = 0$ i.e. $N[X]$ contains only the zero vector. i.e. $Xz = 0$ has no solution save $z = 0$.

Every vector subspace may be represented as the range or the null-space of a matrix. Let V_1 and V_2 be subspaces of R^n of dimension k_1 and k_2 . Let X_1 and X_2 be $(n \times k_1)$ and $(n \times k_2)$ matrices of rank k_1 and k_2 . Then V_1 could be defined as $R[X_1]$ and V_2 as $R[X_2]$.

Theorem 4 If V_1 and V_2 are subspaces of R^n , then $(V_1 + V_2)^\perp = V_1^\perp \cap V_2^\perp$.

Proof. Let $V_1 = R[X_1]$ with X_1 an $(n \times k_1)$ matrix of full rank. Let $V_2 = R[X_2]$ with X_2 an $(n \times k_2)$ matrix of rank k_2 . Then

$$\begin{aligned}
 (V_1 + V_2)^\perp &= [R[X_1] + R[X_2]]^\perp \\
 &= R[X_1 : X_2]^\perp \\
 &= N \begin{bmatrix} X_1^\top \\ X_2^\top \end{bmatrix} \\
 &= \{ \mu \in \mathbb{R}^n : \begin{bmatrix} X_1^\top \\ X_2^\top \end{bmatrix} \mu = 0 \} \\
 &= V_1^\perp \cap V_2^\perp
 \end{aligned}$$

since $N[X_1^\top] = V_1^\perp$ and $N[X_2^\top] = V_2^\perp$. ■

Theorem 5 $(V_1 \cap V_2)^\perp = V_1^\perp + V_2^\perp$.

Proof. Let $V_1 = N[X_1^\top]$ and $V_2 = N[X_2^\top]$. Where X_1 and X_2 are appropriate $(n \times k_1)$ and $(n \times k_2)$ matrices of rank k_1 and k_2 respectively. Then

$$N[X_1^\top] \cap N[X_2^\top] = N \begin{bmatrix} X_1^\top \\ X_2^\top \end{bmatrix}$$

It follows that

$$\begin{aligned}
 N[X_1^\top] \cap N[X_2^\top] &= R[X_1 : X_2]^\perp \\
 &= R[X_1] + R[X_2] \\
 &= V_1 + V_2. \blacksquare
 \end{aligned}$$

3 Projection and Least Squares

3.1 Projection Matrices

The idea of orthogonal projection arises in its simplest form when there are two arbitrary vectors in the plane \mathbb{R}^2 , separated by an acute angle μ , as indicated in Figure 3.1. A perpendicular is dropped

from the vector y onto the line on which the vector x lies, yielding another vector which will be labelled $m = Py$. Here P is the matrix that transforms y into m , and the vector difference $(y - m) = (I_2 - P)y$ is labelled r . The problem is to determine the matrix P . Clearly, since y , x , m and r all lie in the plane \mathbb{R}^2 , P must be a (2×2) matrix.

Notice first that, since x and m are co-linear, m must be obtainable from x by scalar multiplication. If this is written $m = \alpha x$, then α is a positive fraction because m is in the same direction as x but shorter. Indeed, α is nothing more than the proportion of the length of x by which the vector x must be reduced to have the

length of m . Thus

$$s = \frac{kmk}{kxk},$$

where length is determined by the natural scalar product which, for x and y , will be written $(x; y) = x^>y$: It follows that

$$m = Py = x_s = x \frac{\mu kmk}{kxk}. \quad (3.1)$$

Now it was established in chapter 2 that $(x; y) = kxkkyk \cos \mu$ and $\cos \mu = \frac{kmk}{kyk}$.

Thus

$$x^>y = kxkkyk \frac{kmk}{kyk} = kxkkmk,$$

or

$$kmk = \frac{x^>y}{kxk}. \quad (3.2)$$

Substituting (3:2) into (3:1),

$$m = Py = x \frac{1/2 x^>y}{kxk^2}. \quad (3.3)$$

But $kxk^2 = x^>x$, whereupon

$$m = Py = x(x^>x)^{-1} x^>y, \quad (3.4)$$

implying that $P = x x^>x^{-1} x^>$, a (2×2) matrix as required. Notice that P is symmetric and idempotent, $P = P^> = P^2$; this may be checked by applying the rules

for transposition and matrix multiplication. This symmetry and idempotence of P is no coincidence, but rather a consequence of the orthogonality of m and $(y - m) = r$, for $m^T r = 0$ implies

$$(Py)^T (I_2 - P)y = 0$$

or $y^T P^T (I_2 - P)y = 0$ for all y , which cannot be unless $P^T = P^T P$. Upon transposing P^T in the last expression,

$$P = P^T = P^2.$$

This is quite a general rule: if the line x is replaced by the k -dimensional linear subspace L , while y in \mathbb{R}^2 is replaced by y in \mathbb{R}^n , $k < n$, the required $(n \times n)$ projection matrix, P , that takes y in \mathbb{R}^n orthogonally into Py in L , will have the properties of symmetry $P = P^T$ and idempotence ($P = P^2$). These properties are a consequence of the orthogonality condition arising from the natural scalar product.

Consider again the linear model (1:4) of chapter 1 and let $R[X] = L$, L being a subspace of \mathbb{R}^n on which the natural scalar product is defined as before. The vector y is an element of \mathbb{R}^n and the vector X^{-1} is set equal to 1 . Then (1:4) becomes

$$y = 1 + \epsilon \tag{3.5}$$

where $y \in \mathbb{R}^n$ and $1 \in L \subset \mathbb{R}^n$: The coefficient vector ϵ in equation (1:4) is unknown and so the location of 1 in L is unknown; 1 must therefore be estimated and this

may be done by setting $\hat{y} = P y$ where P is the orthogonal projection of y on L , that is, $\hat{y} = P y$ in which P is the orthogonal projection matrix from \mathbb{R}^n onto L . This leaves \hat{y} to be estimated as $\hat{y} = y - (I_n - P) y$, \hat{y} being orthogonal to L , i.e. $\hat{y} \perp L$. Since orthogonality is defined according to $(\cdot; \cdot)$ and L has basis matrix X , then $\hat{y} = P y = X \hat{\beta}$, say; and $\hat{y} = y - X \hat{\beta} = (I_n - P) y$ is required to be orthogonal to every vector x in L , that is, setting $x = X s$; for all $s \in \mathbb{R}^k$,

$$(x; y - \hat{y}) = (X s; y - P y) = s^T X^T (y - P y) = 0.$$

Thus

$$s^T X^T (y - P y) = 0$$

for all $s \in \mathbb{R}^k$ and since the last equality must hold identically in s ,

$$X^T (y - P y) = 0$$

and

$$X^T P y = X^T X \hat{\beta} = X^T y.$$

This reveals that $P = X (X^T X)^{-1} X^T$ and again $P = P^2 = P^T$.

Another property of P that is important is that it is unique for L . For suppose that $W = X M$ is another basis matrix of L , M being non-singular of order k . By corresponding arguments to those used above, the orthogonal projection matrix from

\mathbb{R}^n to L relative to the basis matrix W is $W^{-1}W^T W^{-1}W^T$. But

$$W^{-1}W^T W^{-1}W^T = X M^{-1} M^T X^T X M^{-1} M^T X^T.$$

By the rule for inversion of a matrix product and the non-singularity of M and $X^T X^{-1} X^T$,

$$\begin{aligned} X M^{-1} M^T X^T X M^{-1} M^T X^T &= X M M^{-1} X^T X^{-1} X^T M^{-1} M^T X^T \\ &= X^T X^T X^{-1} X^T \\ &= P. \end{aligned}$$

Thus the matrix P is invariant to the basis used to represent L . P is, of course, $n \times n$ and, not surprisingly since P depends only on X , $R[P] = R[X]$. Moreover, the rank of $P = \text{tr}P = \text{tr}X^T X^T X^{-1} X^T = \text{tr}X^T X^T X^{-1} X^T X = k$.

Finally note that for all $x \in L$,

$$(y - Px) = (y - Py) + (Py - Px). \quad (3.6)$$

Now $(Py - Px) \in L$ and $(y - Py)$, by definition, is orthogonal to every vector in L , including $(Py - Px)$. Thus $(y - Py)^T (Py - Px) = (Py - Px)^T (y - Py) = 0$ and

$$\|y - Px\|^2 = \|y - Py\|^2 + \|Py - Px\|^2 \quad (3.7)$$

which is a representation of Pythagoras' Theorem. Since x is any vector in L and Py is unique (because y is given and P is unique), x may be chosen so as to make $\|y - Px\|^2$ as small as possible. Equation (3:7) indicates that

$$\min_{x \in L} \|y - Px\|^2 = \|y - Py\|^2 \quad (3.8)$$

i.e. that by choosing $x = \hat{x} = Py$, $(y - x)^T (y - x)$ is minimized for given y and all x in L . This is the principle of least squares: the sum of squares $(y - x)^T (y - x)$ is made least by setting $x = Py = \hat{x}$. An alternative way of saying the same thing is that the shortest distance between y and the hyperplane L is $\|(I_n - P)y\|$.

For any $x \in L$, x may be represented uniquely by $x = X\alpha$ for some $\alpha \in \mathbb{R}^k$; because X is a basis of L . Thus if P is the orthogonal projection from \mathbb{R}^n onto L , then

$$Px = X(X^T X)^{-1} X^T x = X\alpha = x.$$

Since x is any vector in L , then $Px = x \forall x \in L$. In the orthogonal complement L^\perp of L in \mathbb{R}^n , recall that $L^\perp = \mathcal{R}[X]^\perp = \mathcal{N}(X^T)$, of dimension $(n - k)$. Let $z \in L^\perp$. Then $Pz = X(X^T X)^{-1} X^T z$ and since $z \in \mathcal{N}(X^T)$, $X^T z = 0$, a condition that holds for every $z \in L^\perp$. Thus, given that P is the orthogonal projection onto L , then $Px = x$ for every $x \in L$ and $Pz = 0$ for every $z \in L^\perp$. Indeed, for any vector y in \mathbb{R}^n , there exists a unique decomposition

$$y = x + z \tag{3.9}$$

with $x \in L$ and $z \in L^\perp$. Moreover

$$Py = Px + Pz = x + 0 = x$$

and

$$(I_n - P)y = (I_n - P)x + (I_n - P)z = 0 + z = z.$$

Since P is the unique orthogonal projection matrix onto L , then $(I_n - P)$ is the unique orthogonal projection matrix onto L^\perp :

Suppose now that L_0 is a subspace that lies entirely in L : $L_0 \subseteq L \subseteq \mathbb{R}^n$. L_0 may be formulated as L subject to a set of r linear restrictions, $r < k$, like $A^T x = 0$, $x \in L$, A being a given $(n \times r)$ matrix of rank r . Then $L_0 = \{x \in L : A^T x = 0\}$. Alternatively, $X_0 = XM$, where M is a known $[k \times (k - r)]$ matrix of rank $(k - r)$, and X_0 is a basis of L_0 . The unique orthogonal projection matrix from \mathbb{R}^n onto L_0 , is then given by $P_0 = X_0 (X_0^T X_0)^{-1} X_0^T = XM (M^T X^T X M)^{-1} M^T X^T$. The last expression cannot be simplified because M is not invertible. Let x_0 be any vector in L_0 . Since X_0 is the $[n \times (k - r)]$ basis for L_0 , x_0 has unique representation $x_0 = X_0 s$ for some $s \in \mathbb{R}^{k-r}$. Thus $x_0 = X_0 s = XM s$. Notice that $P_0 x = x_0$ and that $P x_0 = P X M s = X M s = X_0 s = x_0$. It follows that $P P_0 = P X M (M^T X^T X M)^{-1} M^T X^T = P_0$. Moreover $P_0 P = P_0$. Indeed, if P and Q are orthogonal projection matrices and $PQ = Q$, $QP = Q$ then $\mathcal{R}[Q] \subseteq \mathcal{R}[P]$.

3.2 Means and Decompositions

The equiangular vector in \mathbb{R}^n , which was introduced in chapter 2, plays an important role in statistics as the following special case will reveal. Let x be a vector of n observations x_i ($i = 1, 2, \dots, n$); x will be regarded as an element in n -dimensional Euclidean space \mathbb{R}^n on which the natural scalar product is defined. The equiangular

vector in \mathbb{R}^n is defined by $e = [1; 1; \dots; 1]$. The orthogonal projection from \mathbb{R}^n onto the equiangular line $E = \mathbb{R}[e]$ is

$$p = \frac{1}{n} e e^T$$

Notice that $e^T x = x^T e = \sum_{i=1}^n x_i$ and that $e^T e = n$. It follows that $\frac{1}{n} e e^T x = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$, the mean of the observations x_1, x_2, \dots, x_n . Thus

$$p x = \begin{bmatrix} \bar{x} \\ \bar{x} \\ \vdots \\ \bar{x} \end{bmatrix} = e \bar{x}; \quad (I_n - p) x = \begin{bmatrix} x_1 - \bar{x} \\ x_2 - \bar{x} \\ \vdots \\ x_n - \bar{x} \end{bmatrix}$$

and it follows immediately that $x^T (I_n - p) x = \sum_{i=1}^n (x_i - \bar{x})^2$ because $(I_n - p)$ is symmetric and idempotent.

Suppose now that the n observations x_1, x_2, \dots, x_n represent a randomly selected sample from a population which has mean μ and variance σ^2 . Then

$$(x_i - \mu) = p(x_i - \mu) + (I_n - p)(x_i - \mu) \tag{3.10}$$

is the unique decomposition of the vector $(x_i - \mu) \in \mathbb{R}^n$ into $p(x_i - \mu)$ on E and $(I_n - p)(x_i - \mu)$ on E^\perp . Notice carefully that the dimension of E ($= \text{rank } p$) is $\dim E = 1$ while $\dim E^\perp = \text{rank } (I_n - p) = (n - 1)$; also, since $e \in E$, $(I_n - p)e = 0$. The decomposition (3:10) thus corresponds to the decomposition (3:9), where $N = \frac{1}{n} e e^T =$

$E^?$. Now the right-hand side (RHS) of (3:10) comprises two orthogonal vectors and hence their scalar product is zero. Consequently,

$$(x_i - \bar{x})^2 = (x_i - \bar{x})^2 + n(\bar{x} - \bar{x})^2 \quad (3.11)$$

and, in more familiar statistical notation, (3:11) may be written

$$\sum (x_i - \bar{x})^2 = n \sum (\bar{x} - \bar{x})^2 + \sum (x_i - \bar{x})^2. \quad (3.12)$$

The decomposition (3:11) or (3:12) is important in the development of the F- and t-distributions, on the assumption that the population from which x is drawn is normal.

Note first that the rank of the left-hand side (LHS) of (3:11), and hence of (3:12), is n , the first term on the RHS has rank 1 and the second term has rank $(n - 1)$. Thus the ranks on the RHS sum to the rank of the LHS. Secondly, if the population from which x is drawn is normal, then each

$$\frac{x_i - \bar{x}}{\sigma} \gg \text{NID}(0; 1); \quad i = 1; 2; \dots; n.$$

Hence, dividing (3:12) throughout by σ^2 , the LHS is distributed as the sum of squares of n independent $N(0; 1)$ variates, that is, as a central $\chi^2(n)$ distribution. Thirdly, as a consequence of the sum of ranks on the RHS being equal to the rank of the LHS and the known distribution of $\sum (x_i - \bar{x})^2$, $\frac{n(\bar{x} - \bar{x})^2}{\sigma^2}$ and $\frac{\sum (x_i - \bar{x})^2}{\sigma^2}$ are independently distributed as $\chi^2(1)$ and $\chi^2(n - 1)$. Hence their ratio adjusted by rank (that is

degrees of freedom) is

$$\frac{\frac{1}{n_i} \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2}{(n_i - 1) s^2} \sim F(1; n_i - 1).$$

Hence, eliminating s^2 , taking the square root and re-arranging

$$\frac{(\bar{x}_i - \bar{x})}{\frac{s}{\sqrt{n_i}}} \sim t(n_i - 1) \quad (3.13)$$

where $s^2 = \frac{1}{n_i - 1} \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2$. The expression (3.13) should be familiar from elementary statistics.

Turning back now to the model (1.4) of chapter 1, let the first column of X be e , so that the model may be written

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \dots + \beta_k x_{ik} + \epsilon_i \quad (3.14)$$

for $i = 1; 2; 3; \dots; n$. In this case

$$X = [e; x_{:2}; x_{:3}; \dots; x_{:k}]. \quad (3.15)$$

Notice that for $L = R[X]$, $e \in L$ and hence for the orthogonal projection matrix P onto L ,

$$Pe = e; Pp = p.$$

Similarly, if $X_1 = [e; x_{:2}; x_{:3}; \dots; x_{:j}]$ and $X_2 = [x_{:j+1}; x_{:j+2}; \dots; x_{:k}]$ then $PX_1 = X_1$ and $PX_2 = X_2$. With this background, consider, $y'(I - P)y = \sum_{i=1}^n (y_i - \hat{y}_i)^2$ where

\bar{y} is the mean of the elements of y . A decomposition, corresponding to (3:11), that applies to (3:14) is

$$y' (I_n - p) y = y' (P - p) y + y' (I_n - P) y. \quad (3.16)$$

Such a decomposition is said to be 'adjusted for the mean' since, instead of decomposing $y'y$ into $y'Py$ and $y' (I_n - P) y$, the projection p onto the equiangular line is inserted into the LHS and the first term on the RHS and $py = e\bar{y}$; this is precisely the adjustment for the mean.

The following nomenclature is used:

$$y' (I_n - p) y = S_T^2, \text{ or the total sum of squares;}$$

$$y' (P - p) y = S_E^2, \text{ or the explained sum of squares;}$$

$$y' (I_n - P) y = S_U^2, \text{ or the unexplained sum of squares.}$$

S_T^2 has rank $(n - 1)$, S_E^2 has rank $(k - 1)$ and S_U^2 has rank $(n - k)$ because $\text{tr} (I_n - p) = (n - 1)$, $\text{tr} (P - p) = (k - 1)$ and $\text{tr} (I_n - P) = (n - k)$. Clearly $(n - 1) = (k - 1) + (n - k)$ and, re-arranging (3:16),

$$S_U^2 = S_T^2 - 1 - \frac{S_E^2}{S_T^2}.$$

If $\frac{S_E^2}{S_T^2} = R^2$ then

$$S_U^2 = S_T^2 (1 - R^2). \quad (3.17)$$

Equation (3:17), in particular R^2 , has a nice interpretation arising from the Cauchy-Schwartz inequality. This comes about as follows.

$$\begin{aligned} \frac{S_E^2}{S_T^2} &= \frac{y^{\circ} (P_i - p)y \cdot y^{\circ} (P_i - p)y}{y^{\circ} (I_{n_i} - p)y \cdot y^{\circ} (P_i - p)y} \\ &= \frac{y^{\circ} (P_i - p)y^2}{k(I_{n_i} - p)yk^2 \cdot k(P_i - p)yk^2}. \end{aligned}$$

Now $y^{\circ} (P_i - p)y$ is equal to $y^{\circ} (P_i - p)(I_{n_i} - p)y^a$ because

$$(P_i - p)(I_{n_i} - p) = (P_i - p - p + p) = (P_i - p).$$

Hence

$$R^2 = \frac{y^{\circ} (P_i - p)(I_{n_i} - p)y^a}{k(I_{n_i} - p)yk^2 \cdot k(P_i - p)yk^2}$$

It follows from the Cauchy-Schwartz inequality that

$$0 \leq \frac{y^{\circ} (P_i - p)(I_{n_i} - p)y^a}{k(I_{n_i} - p)yk^2 \cdot k(P_i - p)yk^2} \leq 1. \quad (3.18)$$

Let μ be the angle between $(P_i - p)y$ and $(I_{n_i} - p)y$. Then

$$R^2 = \cos^2 \mu \quad (3.19)$$

and from equation (3:18)

$$1 - R^2 = \sin^2 \mu. \quad (3.20)$$

From (3:17) and (3:18)), $(1 - R^2)$ is seen as the proportion of S_T^2 that is unexplained (S_U^2) and hence R^2 is the proportion of S_T^2 that is explained. Since $(P_i - p)y =$

$\hat{X} = e\hat{y}$, writing \hat{y} for X^\wedge ,

$$y^>(P - p)y = \sum (\hat{y}_i - y)^2$$

which is seen as that part of S_T^2 which the least squares line $X^\wedge = \hat{y} = e\hat{y}$, adjusted by the mean of y ($e\hat{y}$), explains. Finally, from (3:18),

$$R^2 = \frac{\sum (\hat{y}_i - y)(y_i - y)}{\sum (\hat{y}_i - y)^2 \sum (y_i - y)^2}$$

the square of the product moment correlation coefficient between y_i and \hat{y}_i from elementary statistics. Of course, $py = e\hat{y}$ and $py = p[Py + (I_n - P)y]$; but $E \frac{1}{2} L$ and so $pP = p$, $p(I - P) = 0$. Thus $py = e\hat{y} = pPy = pX^\wedge$. It follows that the mean of $y_1; y_2; \dots; y_n$ is \hat{y} and the mean of $\hat{y}_1; \hat{y}_2; \dots; \hat{y}_n$ is also \hat{y} .

3.3 Orthogonality and Least Squares

There are various different ways to formulate the method of least squares using the natural scalar product. The one that emphasizes minimizing the sum of squared errors, from which the term least squares comes, may be written

$$\arg \min_{x \in L} \|y - x\|^2 = \hat{x} = Py. \tag{3.21}$$

In equation (3:21), the expression $\|y - x\|^2$ is squared length defined by the natural scalar product $(\cdot; \cdot)$. A corresponding expression emphasizes the orthogonality of $(y - \hat{x})$ with L . Since Py is unique for given y , then \hat{x} is unique, being defined by the condition

$$(y - \hat{x}; x) = 0 \quad \forall x \in L. \tag{3.22}$$

Expression (3:22) is important and it plays an especially important role in the discussion of invariance below (section 3.6).

Least squares which defines orthogonality relative to the natural scalar product (\cdot, \cdot) is usually referred to as ordinary least squares or OLS. A more general expression for the scalar product is $h(\cdot, \cdot) = (\cdot, A\cdot)$ for some $(n \times n)$, nnd matrix A . To be a proper scalar product, A must be pd but, under certain conditions, A may be nnd whereupon $h(\cdot, \cdot)$ is said to be a quasi-scalar product. If in fact $A (\in I_n)$ is pd, then least square relative to $h(\cdot, \cdot)$ is referred to as generalized least squares (GLS); where A is nnd then least squares relative to $h(\cdot, \cdot)$ belongs to the class of generalized instrumental variables estimators (GIVE). GIVE covers the following methods: instrumental variables (IV), two-stage least squares (2SLS), three-stage least squares (3SLS), generalized 2SLS and generalized 3SLS (G2SLS and G3SLS) and the generalized method of moments (GMM). At this stage, it is not convenient or helpful to enter into a discussion of these methods. Rather, attention will directed toward the general case of least squares relative to the quasi-scalar product $h(\cdot, \cdot)$.

Model (1:4) will take the form of (3:14) and X the form of (3:15). The matrix A may be pd or nnd and its rank will be K , $k \cdot K \cdot n$, so that X^TAX is invertible; alternatively $L^T \setminus R[A] = \hat{A}$. In this setting, the least squares criterion will be (3:22), relative to the scalar product $h(\cdot, \cdot)$. Thus $\hat{\beta} = X^{-1}$ is selected at $\hat{\beta} = X^{\hat{A}}$ such that

$$h(y - \hat{\beta}; x) = (y - \hat{\beta}; Ax) = 0 \quad (3.23)$$

Writing $x = X\beta$ for any $\beta \in \mathbb{R}^k$, then

$$\beta^T X^T A y - \beta^T X^T A X \beta = 0 \quad (3.24)$$

identically in β . Thus since $X^T A X$ is invertible,

$$\hat{\beta} = (X^T A X)^{-1} X^T A y. \quad (3.25)$$

If $F = (X^T A X)^{-1} X^T A = F$, then $F \in \mathbb{R}^{n \times n}$ but $F = F^2$. For this reason F is often referred to as an oblique projection matrix. More precisely, F is the projection matrix on L , orthogonal relative to $h(\cdot, \cdot)$. In the same language, $P = X(X^T X)^{-1} X^T$ is the projection matrix on L , orthogonal relative to the natural scalar product (\cdot, \cdot) .

Recall that in section 3.3, relative to the scalar product (\cdot, \cdot) , $p = e^i e^j e^i e^j$ and

$$\begin{aligned} S_T^2 &= y^T (I_n - p) y \\ &= (y; (I_n - p) y) \\ &= (y - py; y - py) \\ &= k(I_n - p) y k^2 \end{aligned}$$

because $(I_n - p)$ is symmetric and idempotent. Similarly,

$$\begin{aligned} S_E^2 &= y^T (P - p) y = k(P - p) y k^2, \\ S_U^2 &= y^T (I_n - P) y = k(I_n - P) y k^2. \end{aligned}$$

What happens to these quantities when the scalar product is altered from $(; ;)$ to $h; ; i = (; ; A;)$?

² It is easily seen that p has to be replaced by $f = e^i e^j A^k e^l e^m A$ following the form of F using the equiangular vector e . But neither F nor f is symmetric (though each is idempotent). On the other hand, for two vectors z and w , $h z; ; w i = (z; ; A w)$. Thus, for example,

$$\begin{aligned} h y; ; (I_n - f) y i &= (y; ; A (I_n - f) y) \\ &= y; ; A - A e^i e^j A e^k e^l e^m A y \end{aligned}$$

Now

$$\begin{aligned} h (I_n - f) y; ; (I_n - f) y i &= (y; ; e^i e^j A e^k e^l e^m A y; ; \\ &\quad A y; ; A e^i e^j A e^k e^l e^m A y) \\ &= y^> A y; ; y^> A e^i e^j A e^k e^l e^m A y \\ &\quad - y^> A e^i e^j A e^k e^l e^m A y \\ &\quad + y^> A e^i e^j A e^k e^l e^m A e^i e^j A e^k e^l e^m A y \\ &= y^> A y; ; y^> A f y; ; y^> A f y + y^> A f y, \end{aligned}$$

because $A f = f^> A$, that is $A f$ is symmetric; and $f^> A f = A f$. Hence

$$\begin{aligned} h y; ; (I_n - f) y i &= h (I_n - f) y; ; (I_n - f) y i \\ &= k (I_n - f) y k_{ii}^2 \end{aligned}$$

where k_{hi}^2 denotes length relative to the scalar product $h; ; i$.

It follows from this discussion that $(I_n - F)$, $(F - f)$ and $(I_n - f)$ are symmetric relative to the scalar product $h; ; i$; using $(F - f)$, then

$$A(F - f) = (F - f)^T A = (F - f)^T A(F - f)$$

or for arbitrary vectors w and z in \mathbb{R}^n

$$h; (F - f) w = h(F - f) z; w = h(F - f) z; (F - f) w.$$

Notice, incidentally, that just as P obeys, relative to $(; ;)$,

$$(; P;) = (P; ;) = (P; P;) = \mathbf{i}; P^2; \text{ } ,$$

so F obeys, relative to $h; ; i$

$$h; F; i = hF; ; i = hF; ; F; i = \text{ } ; F^2; \text{ } .$$

The same holds for f , $(I_n - f)$; $(F - f)$ and $(I_n - F)$, indeed any projection matrix which is orthogonal relative to $h; ; i$.

Turning back to S_T^2 , S_E^2 and S_U^2 , relative to the scalar product $h; ; i$; these become

$$S_T^2 = k(I_n - f) y k_{hi}^2,$$

$$S_E^2 = k(F - f) y k_{hi}^2,$$

$$S_U^2 = k(I_n - F) y k_{hi}^2.$$

It is natural to define R^n , relative $h; :i$, by analogy with R^2 relative to $(; :)$ in (3:18).

Thus from (3:18)

$$R^2 = \frac{(fP_i pg; fl_n i pg y)^2}{k(P_i p) y k^2 k(l_n i p) y k^2}$$

which is $\cos^2 \mu$ relative to the natural scalar product. When $h; :i$ replaces $(; :)$, (3:18)

is replaced by

$$R_{hi}^2 = \frac{h(F_i f) y; (l_n i f) y i^2}{k(F_i f) y k_{hi}^2 k(l_n i f) y k_{hi}^2} \quad (3.26)$$

which is $\cos^2 \mu$ relative to $h; :i$.

One way to look upon the scalar product $h; :i$ is as the natural scalar product in transformed vectors. If, for example, A is pd, then there exists a non-singular R of order n such that $R^>R = A$. Thus $hw; zi = (w; Az) = (w; R^>Rz) = (Rw; Rz)$. Thus when w and z are transformed into $Rw = r$ and $Rz = s$, then $hw; zi = (r; s)$ and any analysis of vectors relative to $h; :i$ may be recast, by a transformation matrix R , into a corresponding analysis of transformed vectors relative to the natural scalar product $(; :)$.

If the matrix A is an orthogonal projection matrix of rank K (and hence is $n \times n$), then there exists an $(n \times K)$ matrix B such that $A = BB^>$ and $B^>B = I_K$. In this case, $hw; zi = (w; Az) = \begin{matrix} i \\ B^>w; B^>z \end{matrix} = (p; q)$ with $B^>w = p; B^>z = q$. Note carefully that p and q will be K -tuples. The transformation $B^>$ is an example of a partial isometric matrix or partial isometry. Suppose that A is the orthogonal projection

matrix on a subspace M of dimension K (along M^\perp of dimension $(n - K)$), then for every $x \in M$, $Ax = x$ and $kB^>xk^2 = x^>BB^>x = x^>Ax = x^>x = kxk^2$. Thus $kB^>xk = kxk$; for every $z \in M^\perp$, $B^>z = 0$. For example consider $P = X(X^>X)^{-1}X^>$ onto L along L^\perp . Since X has full rank, $X^>X$ is pd and hence there exists a $(k \times k)$ non-singular Q such that $Q^>Q = (X^>X)^{-1}$. Then $P = XQ^>QX^> = BB^>$ with $B = XQ^>$ an $(n \times k)$ matrix. For any $x \in L$, $kB^>xk = \sqrt{x^>BB^>x} = \sqrt{x^>x} = kxk$ as required. For any $z \in L^\perp$, $B^>z = QX^>z = 0$ since $L^\perp \subset N(X^>)$. Partial isometries are important in econometrics since they are transformations which hold lengths of vectors in a subspace fixed while annihilating vectors lying orthogonal to that subspace.

3.4 Restricted Least Squares

By restricted least squares is meant the model (1:4) subject to $r < k$ linear restrictions on the estimation of β ; corresponding restrictions may be placed on the estimation of σ^2 . Let these restrictions be defined in terms of a known $(k \times r)$ matrix B of rank r and a corresponding $(n \times r)$ matrix A :

$$B^>\beta = 0 \quad A^>\sigma^2 = 0.$$

For given B and arbitrary $(n \times n)$ pd matrix Q ,

$$B^>\beta = 0 \quad B^>QX^>QX^>QX^>QX^> = 0$$

yielding for $\beta = X^{-1}A$ and $A = QX^i X^j QX^{i-1} B$, $A^{>1} = 0$. A is not unique for given B , but $PA = X^i X^j X^{i-1} B$ is. $B^{>-} = 0$ and $A^{>1} = 0$ are called restraint equations. Alternatively, if there are r linear restraints on k coefficients β , then there must be $(k - r)$ 'free' coefficients to estimate; let these be the $[(k - r) \times 1]$ vectors β° . Then β and β° must be related by

$$\beta = M\beta^{\circ}$$

where M is a $[k \times (k - r)]$ matrix of elements determined by $B^{>-} = 0$.

Example 5 Let $k = 4$ and $r = 2$, so that

$$\beta = \begin{matrix} 2 & 3 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{matrix}; B^{>-} = \begin{matrix} 2 & 3 \\ 6 & 2 & 1 & 0 & 0 \\ 4 & 0 & 0 & 1 & 1 \end{matrix} \begin{matrix} 2 & 3 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{matrix} = 0.$$

The second equations imply $\beta_1 = \beta_2$ and $\beta_3 = \beta_4$: Thus $(k - r) = 4 - 2 = 2$ β° is (2×1) . Let $\beta^{\circ}_1 = \beta_1$ and $\beta^{\circ}_2 = \beta_3$. Then

$$\beta = \begin{matrix} 2 & 3 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{matrix} = M\beta^{\circ} = \begin{matrix} 2 & 3 & 2 & 3 \\ \beta^{\circ}_1 & \beta^{\circ}_2 & \beta^{\circ}_1 & \beta^{\circ}_2 \\ 6 & 2 & 1 & 0 & 0 \\ 4 & 0 & 0 & 1 & 1 \end{matrix} \begin{matrix} 2 & 3 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{matrix} = \begin{matrix} 2 & 3 \\ \beta^{\circ}_1 & \beta^{\circ}_2 \\ 6 & 2 & 1 & 0 & 0 \\ 4 & 0 & 0 & 1 & 1 \end{matrix} \begin{matrix} 2 & 3 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{matrix}.$$

The equations $M\beta^{\circ} = \beta$ are called freedom equations. ■

The model (1:4) may now be expressed as

$$y = X\beta + \epsilon \quad (3.27)$$

subject to $B\beta = 0$ or as

$$y = X\beta + \epsilon \quad (3.28)$$

subject to $A\beta = 0$; $\beta \in L$, in restraint equation form; or as

$$y = XM\beta + \epsilon, \quad (3.29)$$

$M\beta = \beta$; in freedom equation form. From (3:28) and (3:29) the subspace $L = R[X]$ of dimension k is now restricted to $L_0 = R[XM]$ of dimension $(k - r)$. If the entire linear space R^n is considered, then this may be decomposed into orthogonal subspaces as follows

$$R^n = L_0 \oplus (L_0^\perp \setminus L) \oplus L^\perp. \quad (3.30)$$

The corresponding decomposition of y relative to $(; ; :)$ is

$$y = P_0 y + (P_1 - P_0) y + (I_n - P) y \quad (3.31)$$

or adjusting for the mean,

$$(I_n - P) y = (P_0 - P) y + (P_1 - P_0) y + (I_n - P) y \quad (3.32)$$

in which $P_0 = XM(XM)^\dagger XM^\dagger X^\dagger XM^\dagger X^\dagger$. For later reference, it will be helpful to find an expression for $(P_1 - P_0)$ in terms of the matrix A in equation (3:28) and P .

Notice that in (3:28), the restricted subspace L_0 must be

$$L_0 = L \setminus N^{\mathbf{f}} A^{\mathbf{a}}.$$

The subspace of interest, from (3:30) is $L_0^? \setminus L$:

$$\begin{aligned} L_0^? \setminus L &= L \setminus N^{\mathbf{f}} A^{\mathbf{a}} \setminus L \\ &= L^? + R[A]^{\mathbf{a}} \setminus L \\ &= L^? + PR[A] + (I_n - P)R[A]^{\mathbf{a}} \setminus L \\ &= L^? + R[PA]^{\mathbf{a}} \setminus L \\ &= R[PA]. \end{aligned}$$

The dimension of $R[PA]$ is the rank of PA and this will be the same as the rank of A if and only if (i[®]) $R[A] \setminus L^? = ;$. This last condition is reasonable since $A = QX^{\mathbf{i}} X^{\mathbf{a}} QX^{\mathbf{c}_i^{-1}} B$ for arbitrary pd Q and $PA = X^{\mathbf{i}} X^{\mathbf{a}} X^{\mathbf{c}_i^{-1}} B$ then has the rank of B , namely r . Although A is arbitrary up to the choice of Q , PA is unique; indeed $PA = A$ where $Q = I_n$ which is admissible.

Given $L_0^? \setminus L = R[PA]$ and $R[A] \setminus L^? = ;$, the unique orthogonal projection on $L_0^? \setminus L$ is from (3:30) and (3:31),

$$\begin{aligned} P_i P_0 &= PA^{\mathbf{i}} A^{\mathbf{a}} P A^{\mathbf{c}_i^{-1}} A^{\mathbf{a}} P \\ &= X^{\mathbf{i}} X^{\mathbf{a}} X^{\mathbf{c}_i^{-1}} B^{\mathbf{h}} B^{\mathbf{a}} X^{\mathbf{c}_i^{-1}} B^{\mathbf{i}_i^{-1}} B^{\mathbf{a}} X^{\mathbf{c}_i^{-1}} X^{\mathbf{a}}. \end{aligned} \quad (3.33)$$

Thus the decomposition (3:31) yields

$$(P_i - P_0)y = X_i' X' X_i^{-1} B' B^{-1} h + X_i' X' X_i^{-1} B' B^{-1} X_i^{-1} B' X_i' X' X_i^{-1} y. \quad (3.34)$$

Setting $X^\Delta = Py$ and $X_0^\Delta = P_0 y$, where Δ is the (unrestricted) estimator of β in (1:4) and Δ_0 is the corresponding restricted estimator in (3:27), (3:34) may be recast as

$$X^\Delta - X_0^\Delta = X_i' X' X_i^{-1} B' B^{-1} h + X_i' X' X_i^{-1} B' B^{-1} X_i^{-1} B' X_i' X' X_i^{-1} X^\Delta$$

leading to

$$\Delta_0 = \Delta - X_i' X' X_i^{-1} B' B^{-1} X_i^{-1} B' X_i' X' X_i^{-1} X^\Delta, \quad (3.35)$$

a familiar expression in least squares theory.

When $A^{>1} = 0$ in (3:28) is replaced by the (a± ne) restrictions $A^{>1} = \mu$, $\mu \in 0$, $1 \leq L$, then there must exist a solution β^{1^*} to $A^{>1} = \mu$, say $A^{>1^*} = \mu$, whereupon $A^{>1}(\beta - \beta^{1^*}) = 0$, $1 \leq L$. Then (3:28) is replaced by

$$y - X\beta^{1^*} = (y - X\beta^{1^*}) + u; \quad A^{>1}(\beta - \beta^{1^*}) = 0; \quad 1 \leq L,$$

β^{1^*} being regarded as a given constant vector. Repeating the same analysis, (3:35) is replaced by

$$\Delta_0 = \Delta - X_i' X' X_i^{-1} B' B^{-1} X_i^{-1} B' X_i' X' X_i^{-1} X^\Delta - X_i' X' X_i^{-1} B' B^{-1} X_i^{-1} B' X_i' X' X_i^{-1} X^\Delta. \quad (3.36)$$

Thus the formulation (3:28) is seen as being quite general in the sense that its least squares solution can accommodate the more general $A^{>1} = \mu$.

The relation between the freedom equation $y = M\beta$ and the restraint on equation $B'y = 0$ is straightforward to establish since $0 = B'y = B'M\beta$. Since β is chosen without restriction, $B'M = 0$, i.e. the vectors comprising M are selected from the null of B' . $N(B')$ has dimension $(k - r)$ and M has $(k - r)$ columns. Hence given B , M is unique.

Example 6 From example 5,

$$M = \begin{matrix} & \begin{matrix} 2 & 3 \end{matrix} \\ \begin{matrix} 6 \\ 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 4 \end{matrix} & \begin{matrix} 1 & 0 \\ 2 & 0 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{matrix} \end{matrix}; \quad B = \begin{matrix} & \begin{matrix} 2 & 3 \end{matrix} \\ \begin{matrix} 6 \\ 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 4 \end{matrix} & \begin{matrix} 2 & 0 \\ i & 1 \\ 0 & 1 \\ 0 & i \\ 0 & 1 \end{matrix} \end{matrix}.$$

Hence

$$B'M = \begin{matrix} & \begin{matrix} 2 & 3 \end{matrix} \\ \begin{matrix} 6 \\ 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 4 \end{matrix} & \begin{matrix} 1 & 0 \\ 2 & 0 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{matrix} \end{matrix} = \begin{matrix} & \begin{matrix} 2 & 3 \end{matrix} \\ \begin{matrix} 6 \\ 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 4 \end{matrix} & \begin{matrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{matrix} \end{matrix}$$

as is required by the relation $B'M = 0$. The nullity of B' is $k - r = 4 - 2$ in this case and M has two columns. Thus the solution is unique. ■

3.5 The Frisch-Waugh Theorem

It often happens in econometric work that the raw data available need to be adjusted to make them applicable to the problem under examination. If the application in-

volves a linear least squares regression and if the adjustments required are linear, the question arises: is it appropriate to adjust the data before applying them to OLS regression, or should the raw data be applied directly to the regression, augmented to make allowances for the needed adjustments? This question is answered in the Frisch-Waugh Theorem which was first applied in relation to linear adjustment for seasonality.

Consider the model (1:4) re-written as

$$y = X_1 \beta_1 + X_2 \beta_2 + \epsilon \quad (3.37)$$

in which X_1 is the first k_1 columns of X and X_2 the remaining $(k - k_1) = k_2$ columns. X_1 is to be thought of as the 'observations' on each of the variables entering into the linear adjustments to be made, while X_2 represents n observations on each of the k_2 variables which are of economic significance in 'explaining' y : Thus equation (3:37) is the regression augmented to include the required adjustment. For example, in terms of seasonality, y is the vector of raw data, that is, data that have not been de-seasonalized, X_1 represents the seasonal effects and X_2 the economic variables that theory predicts will explain the behavior of y . These latter observations are also in the form of raw data.

Now suppose y and X_2 are each adjusted for seasonality by fitting the following

artificial regression by OLS:

$$y = X_1\mu_1 + u_1 \quad (3.38)$$

$$X_2 = X_1\beta_2 + U_2. \quad (3.39)$$

In (3:38), μ_1 is $(k_1 \times 1)$ and u_1 is $(n \times 1)$; in (3:39) β_2 is a $(k_1 \times k_2)$ matrix of coefficients and U_2 is $(n \times k_2)$. The OLS regressions of (3:38) and (3:39) are, respectively,

$$y = P_1y + (I_n - P_1)y = X_1\hat{\mu}_1 + \hat{u}_1 \quad (3.40)$$

$$X_2 = P_1X_2 + (I_n - P_1)X_2 = X_1\hat{\beta}_2 + \hat{U}_2. \quad (3.41)$$

Here $P_i = X_i(X_i'X_i)^{-1}X_i'$, $i = 1, 2$; $(I_n - P_i)$ is denoted M_i . Having adjusted the variables for seasons, the true economic effects, devoid of seasonal influences, may be determined via OLS regression obtained by premultiplying (3:37) by M_1 :

$$M_1y = M_1X_2\beta_2 + v_1 \quad (3.42)$$

whereupon $X_2\hat{\beta}_2 = X_2(X_2'M_1X_2)^{-1}X_2'M_1y$. This is precisely the estimate of $X_2\beta_2$ obtained by applying OLS directly to (3:37). Of course $X_2\beta_2 = F_2y$, as indicated earlier. This simple result indicates that de-seasonalizing the variables linearly and then calculating the appropriate linear regression produces precisely the same estimate of β_2 as would have been produced by fitting (3:37) directly.

The Frisch-Waugh Theorem in fact provides the same answer via another route. The intuition here is to de-seasonalize y by fitting (3:38), yielding (3:40), and then

fitting

$$M_1 y = X_{1,1} + X_{2,2} + v_2. \quad (3.43)$$

The inclusion of X_2 in (3:43) is justified on the ground that the matrix X_2 has not yet been de-seasonalized. Consequently $\hat{X}_{2,2}$ should reflect only that part of X_2 which is orthogonal to $R[X_1]$ that is $M_1 X_2$. When estimated,

$$\hat{X}_{2,2} = X_2 (X_2^T M_1 X_2)^{-1} X_2^T M_1 y$$

and thus $\hat{X}_{2,2}$ from (3:42) = $\hat{X}_{2,2} = X_2^{\Delta} \hat{\mu}_2$ from (3:37) = $F_2 y$. Notice that in fitting (3:43),

$$\begin{aligned} \hat{X}_{1,1} &= X_1 (X_1^T M_2 X_1)^{-1} X_1^T M_2 M_1 y \\ &= F_1 y = X_1 (X_1^T M_2 X_1)^{-1} X_1^T M_2 X_1 (X_1^T X_1)^{-1} X_1^T y \\ &= F_1 y = P_1 y \\ &= X_1^{\Delta} \hat{\mu}_1 \end{aligned}$$

Thus the fitted regression equation (3:43) is

$$y = X_1^{\Delta} \hat{\mu}_1 + X_2^{\Delta} \hat{\mu}_2 + \hat{v}_2 \quad (3.44)$$

or using the projection matrices P_1 , F_2 and P ,

$$y = P_1 y = (F_1 + P_1) y + F_2 y + (I_n - P) y. \quad (3.45)$$

In (3:45), \hat{v}_2 of (3:44) has been written as $(I_n - P)y$. This may be seen from the following.

$$\begin{aligned}\hat{v}_2 &= (I_n - P)M_1y \\ &= (I_n - P_1 - P + P_1)y \\ &= (I_n - P)y \\ &= \hat{v}_2\end{aligned}$$

\hat{v}_2 coming from the fitting of equation (3:37)), as indicated in (3:45).

The Frisch-Waugh Theorem (which will not be stated as a proper theorem) will be summarized as follows. If it is desired to estimate β_2 by applying OLS to equation (3:37), then the following alternative methods are equivalent.

1. Direct estimation of (3:37) by OLS.
2. Linear adjustment by OLS of y and X_2 as in equations (3:38) and (3:39), followed by OLS estimation of the linear regression of the adjusted y on the adjusted X_2 as given in (3:40) and (3:41).
3. Linear adjustment of y as in (3:42), followed by the OLS estimation of β_2 as in (3:44) or the estimation of β_2 as in (3:43).

These calculations also reveal that:

the same? The answer is provided by a theorem which has come to be known in econometrics as Kruskal's Theorem. This will be proved when A is pd. The theorem is stated as:

Theorem 6 Let $y = \beta + \epsilon$ be a model in which y and ϵ lie in \mathbb{R}^n and $\beta \in L$, a subspace of dimension $k < n$. On \mathbb{R}^n is defined the natural scalar product $(\cdot; \cdot)$ and an alternative $h(\cdot; \cdot) = (\cdot; A\cdot)$, A being pd. The OLS estimate of β in L is the same as the GLS estimate of β in L if and only if the subspace L is invariant under A .

Proof. If L is invariant under A , then for every x in L , $Ax \in L$ and hence $AL \subseteq L$. Let L be invariant under A . Then for OLS, $(y - \hat{\beta}; Ax) = 0$ for every $x \in L$. But for GLS

$$0 = h(y - \hat{\beta}; x) = (y - \hat{\beta}; Ax) = (y - \hat{\beta}; w)$$

for every $w \in L$. Since orthogonality ensures uniqueness of $\hat{\beta}$ and β and each is such that

$$(y - \hat{\beta}; x) = (y - \beta; x) = 0$$

for all $x \in L$, $\hat{\beta} = \beta$.

Now let $\hat{\beta} = \beta$. Then

$$(y - \hat{\beta}; x) = (y - \beta; Ax) = 0$$

for all $x \in L$. Since $\hat{\beta} = \beta$, by uniqueness, $Ax \in L$ and L is invariant under A . ■

Example 7 Let $X = [e; x_2; x_3; \dots; x_k]$ and $A = I_n(1 - \frac{1}{2}) + ee^T$. Show that $L = R[X]$ is invariant under A . Notice that if $AX\mu$, $\mu \in R^k$ may always be written as $XM\mu$ for some non-singular matrix M of order k , then L must be invariant under A . For any $x \in L$, $x = X\mu$, $\mu \in R^k$. Thus $Ax = AX\mu = XM\mu = X \cdot \mu \in L$. Now AX is given by

$$\begin{aligned}
 & (I_n(1 - \frac{1}{2}) + ee^T) [e; x_2; x_3; \dots; x_k] \\
 = & e(1 - \frac{1}{2}) + n \cdot \frac{1}{2} X(1 - \frac{1}{2}) + [enx_2; enx_3; \dots; enx_k]
 \end{aligned}$$

where $X = [x_2; x_3; \dots; x_k]$. But the RHS of the last equality may be re-written

$$\begin{aligned}
 AX &= \begin{bmatrix} 1 + (n-1)\frac{1}{2} & nx_2 & nx_3 & \dots & nx_k \\ 0 & I_{k-1}(1 - \frac{1}{2}) & & & \end{bmatrix} \\
 &= XM.
 \end{aligned}$$

Hence, $AX\mu = XM\mu = X \cdot \mu$ and $L = R[X]$ is invariant under A . ■

A useful theorem for proving invariance is

Theorem 7 If $L \subset R^n$ is invariant under the $(n \times n)$ matrix A , then $PAP = AP$ for every projection P on L . If $PAP = AP$ for some projection P on L , then L is invariant under A . ■

Example 8 From the previous example, note that $e \in L$. Hence $Pe = e$. In addition

$$\begin{aligned} PAP &= P \begin{bmatrix} 1 & \\ & \ddots \\ & & 1 \end{bmatrix} (1 \ j \ \frac{1}{2}) + ee^T \frac{1}{2} P \\ &= \begin{bmatrix} 1 & \\ & \ddots \\ & & 1 \end{bmatrix} P (1 \ j \ \frac{1}{2}) + ee^T \frac{1}{2} P \\ &= \begin{bmatrix} 1 & \\ & \ddots \\ & & 1 \end{bmatrix} (1 \ j \ \frac{1}{2}) + ee^T \frac{1}{2} P, \end{aligned}$$

which establishes invariance without the pain of manipulations. ■

When A is nnd, Kruskal's Theorem may be extended to include this condition. In Theorem 6, the condition of invariance is stated $AL \subseteq L$, meaning that AL implies L . If A is non-singular, then $AL = L$ and hence L is also invariant under A^{-1} i.e. $L = A^{-1}L$. When A is singular and L is invariant under A , then AL will lie in L but will not cover it, and so it is still correct to write $AL \subseteq L$.

4 The Multivariate Normal and Related Distributions

4.1 Probability and Random Variables

There is some mathematical difficulty in defining the concepts of probability and random variable. The conventional point of departure is a random experiment.

Definition 13 A process giving rise to observable outcomes is called a random experiment, denoted by E , if the following conditions are upheld.

1. All outcomes of E are known a priori.
2. In any particular performance of E , called a trial, the outcome cannot be known a priori.
3. E may be repeated under the same conditions. ■

Conditions 1{3 in Definition 13 may not suit every situation to which the concept of probability may apply. However, the definition of a random experiment is a starting point from which an expansion of ideas may take place.

The set of all possible outcomes of E forms a sample space denoted by S . The elements of S are called sample points or, more usually, elementary events. On the basis of the elementary events, more broadly based events comprising subsets of S may be formed. These subsets are called events in S and the way in which they are constructed is represented as an algebra called a σ -algebra.

Definition 14 Let \mathcal{F} denote the class of subsets of S which have the following two properties.

1. $A \in \mathcal{F} \Rightarrow A^c \in \mathcal{F}$, A^c being the complement of A in S ; that is, \mathcal{F} is closed under complementation.
2. $A_i \in \mathcal{F}, i = 1, 2, \dots; \Rightarrow \bigcup_{i=1}^{\infty} A_i \in \mathcal{F}$; that is, \mathcal{F} is closed under countable unions. The class \mathcal{F} is known as a σ -algebra or a σ -field, and $(S; \mathcal{F})$ is called a measurable space. $A \in \mathcal{F}$ is said to be an \mathcal{F} -measurable set. ■

One further element is required for a probability space. This is provided by:

Definition 15 Probability $P(\cdot) : \mathcal{F} \rightarrow \mathbb{R}$ is a set function defined on \mathcal{F} satisfying:

1. $P(A) \geq 0$ for every A in \mathcal{F} ;
2. $P(S) = 1$;
3. $P(\bigcup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} P(A_i)$ for every sequence of disjoint events in \mathcal{F} , that is for events having the property $A_i \cap A_j = \emptyset; i \neq j$.

The function $P(\cdot)$ is referred to as a probability measure. ■

Putting definitions 13{15 together permits:

Definition 16 Let a random experiment have sample space S and let S be endowed with a σ -field \mathcal{F} and a probability measure $P(\cdot)$. The triple $(S; \mathcal{F}; P(\cdot))$ constitutes a probability space. ■

The concept of a random variable may be developed from a probability space in the following way. A random variable $X(\cdot)$ is an \mathcal{F} -measurable function from S to the real line \mathbb{R} , that is to say, for every $s \in S$, $X(s) = x \in \mathbb{R}$. This raises a problem: while the events of the probability space are determined by the σ -field \mathcal{F} defined on S , the mapping $X : S \rightarrow \mathbb{R}$ takes points in S directly onto the real line, thereby 'avoiding' \mathcal{F} . Clearly, the mapping $X : S \rightarrow \mathbb{R}$ must be consistent with the definition of the σ -field \mathcal{F} . Briefly, there must be defined on \mathbb{R} a σ -algebra that is consistent with \mathcal{F} . The σ -algebra on the real line is referred to as a Borel σ -field and is denoted by \mathcal{B} . The elements of \mathcal{B} are Borel sets denoted by $B \in \mathcal{B}$. This arrangement leads to the following definition.

Definition 17 A random variable $X(\cdot)$ is an \mathcal{F} -measurable function from S to \mathbb{R} such that, for every Borel set $B \in \mathcal{B}$ on \mathbb{R} the set $X^{-1}(B) = \{s : X(s) \in B; s \in S\}$ satisfies $X^{-1}(B) \in \mathcal{F}$. ■

In Definition 17, suppose that $S_0 \subset S$ is a subset of S such that $X(S_0) = B$. Since $S_0 \in \mathcal{F}$ by definition, the statement $X^{-1}(B) \in \mathcal{F}$ is merely a requirement that the inverse image of B forms a proper subset S_0 of S , as defined by the σ -field \mathcal{F} .

Three points should be noted about the definition of a random variable. First, a random variable is defined for a specific σ -algebra, \mathcal{F} . Second, if $X(\cdot) : S \rightarrow \mathbb{R}$ is to be a random variable, then the order of argument is from the Borel σ -field \mathcal{B} to the σ -field \mathcal{F} , not the other way round. Finally, a random variable is, strictly

speaking, neither random nor a variable. $X(\cdot)$ is a real-valued function and the terms random and variable do not enter its definition. Randomness (or probability) enters the picture only after the definition of the random variable $X(\cdot)$ has been stated, as a means of completing the mathematical model induced by the definition of X .

In discussing random variables, it is often necessary to discuss pairs of random variables, or triplets or n -tuples. It is then standard to refer to a random vector or a vector-valued random variable. Before a random vector can be defined, it is necessary to introduce the concept of Borel sets of a finite-dimensional Euclidean space V on which the scalar product $(\cdot; \cdot)$ is defined. Let $\|x\| = (x; x)^{\frac{1}{2}}$. The open ball of radius r about the point z in V is defined as $B_r(z) = \{x : \|x - z\| < r\}$. The open ball is used to define a Borel field on V , to stand in place of the Borel field on \mathbb{R} used in defining a random variable.

Definition 18 The Borel field of V , denoted by B_V , is the smallest σ -algebra that contains all of the open balls. $(V; B_V)$ is a measurable space. ■

Note that B_V is independent of the scalar product of V . If two scalar products are defined on V , these will generate the same Borel field B_V .

In defining a random vector, it is necessary to have an appropriate probability space which, as before, may be regarded as the triple $(S; F; P(\cdot))$. With this in place, the definition of a random vector is stated as:

Definition 19 A random vector $X(\cdot)$ in V is a function mapping S into V such that, for each Borel set $B \in \mathcal{B}_V$, $X^{-1}(B) \in \mathcal{F}$, $X^{-1}(B) = \{s : X(s) \in B; s \in S\}$. ■

In Definition 19, let $V = \mathbb{R}^n$. Then the Borel set B may be regarded as comprising n subsets $B_1; B_2; \dots; B_n$ and $\mathcal{B}_V = \{B \in \mathcal{B} \mid B = B_1 \times B_2 \times \dots \times B_n\}$. Then $X^{-1}(B) = \{s : X_1(s) \in B_1; X_2(s) \in B_2; \dots; X_n(s) \in B_n; s \in S\}$, where $X_1; X_2; \dots; X_n$ are the co-ordinates of X corresponding to the sample point s .

The space on which a random vector is defined is seldom of much interest and so the argument of a random vector is ordinarily suppressed. What is of interest is the distribution of X on V . Let X be defined on S to V where $(S; \mathcal{F}; P(\cdot))$ is the probability space. For each set $B \in \mathcal{B}_V$, let $Q(B) = P(X^{-1}(B))$. Obviously $Q : \mathcal{B}_V \rightarrow \mathbb{R}$ is a probability measure on \mathcal{B}_V , and $(V; \mathcal{B}_V; Q)$ is also a probability space. Q is called the induced distribution of X , meaning that Q is induced by X and P . In summary, given the probability space $(S; \mathcal{F}; P(\cdot))$, a random vector defined by $X : S \rightarrow V$ such that $X^{-1}(B) = \{s : X(s) \in B; B \in \mathcal{B}_V; s \in S\}$, then $Q = P(X^{-1}(B))$ on \mathcal{B}_V is the induced distribution of X on \mathcal{B}_V . It is also possible to proceed the other way round. Suppose it is given that Q is a probability measure on \mathcal{B}_V (V being a finite dimensional vector space on which a scalar product is defined). Then there will exist a random variable X and a probability space $(S; \mathcal{F}; P(\cdot))$ such that Q is the induced distribution of X . The proof of this claim is achieved essentially by substitution. Setting $V = S$, $\mathcal{B}_V = \mathcal{F}$ and $P = Q$, let X be a random vector such that $X(v) = v$ for

$v \in V$. Then $P(X \in B) = Q$ for each $B \in \mathcal{B}_V$ and Q is the induced distribution on $\mathcal{B}_V = \mathcal{F}$.

Although Q is, properly speaking, the distribution induced by X and P on \mathcal{B}_V , it is unusual to use this phraseology. Rather, it is usual to say: " X is a random vector which ranges over V with distribution Q ." This form will be used throughout the sequel.

A function $f : V \rightarrow W$ is called Borel measurable if the inverse image of each set $B \in \mathcal{B}_W$ is in \mathcal{B}_V , i.e. if $f^{-1}(B) = \{x : f(x) \in B; x \in V\} \in \mathcal{B}_V$, in which W is a metric space, \mathcal{B}_W is a Borel σ -field on W and $(W; \mathcal{B}_W)$ is a Borel measurable space. If f is Borel measurable, then $f(X)$ is a random vector in W . When f is continuous, it is Borel measurable. If $W = \mathbb{R}$ and f is Borel measurable, then $f(X)$ is a real-valued random variable.

4.2 Expectations

Definition 20 Let X be a random vector which ranges over V with distribution Q and let f be a real-valued, Borel-measurable function defined on V . The expectation of $f(X)$, assumed to be finite, is $E f(X) = \int_V f(x) dQ$ the integral being a Lebesgue integral. ■

Let $V = \mathbb{R}^n$ on which is defined the natural scalar product. Let $Q(dx)$ denote the standard Lebesgue measure on \mathbb{R}^n . If q is a non-negative function on \mathbb{R}^n such that

the Riemann integral $\int_{\mathbb{R}^n} q(x) dx = 1$, then q is called a density function. $Q(B) = \int_B q(x) dx$ is a probability measure on \mathbb{R}^n implying that Q is the distribution of some random variable, say X . Let $e_1; e_2; \dots; e_n$ denote the orthonormal basis in \mathbb{R}^n (whose sum is the equiangular vector). Then $(e_i; X)$ represents X_i the i 'th co-ordinate of X . Assuming finite expectations, $EX_i = \int_{\mathbb{R}^n} (e_i; x) q(x) dx = \mu_i$, the mean of X_i . The vector μ comprises co-ordinates $\mu_1; \mu_2; \dots; \mu_n$ and is the mean vector of X . Notice that for any non-random vector $z \in \mathbb{R}^n$, $E(z; X) = E[\sum_i (z_i e_i; X)] = \sum_i [z_i E(e_i; X)] = \sum_i z_i \mu_i = (z; \mu)$ for all $z \in \mathbb{R}^n$, implying that μ is unique. This leads to the appealing equation $E(z; x) = (z; Ex)$ which is certainly valid in the co-ordinate case.

In referring to Euclidean subspaces, it is convenient to write the Euclidean vector space V with scalar product $(; ;)$ as $fV; (; ;)g$. This notation is used in:

Proposition 8 Let $X \in fV; (; ;)g$ with mean μ . Let $Y \in fW; h; ;ig$, $! \in W$ and A be the matrix of a linear transformation $V \rightarrow W$. The random variable $Y = AX + !$ has expectation $A\mu + !$.

Proof. Let $z \in \mathbb{R}^n$.

$$\begin{aligned}
 E\langle z, Y \rangle &= E\langle z, AX + \mu \rangle \\
 &= E\langle z, AX \rangle + \langle z, \mu \rangle \\
 &= E\langle A^T z, X \rangle + \langle z, \mu \rangle \\
 &= \langle A^T z, \mu \rangle + \langle z, \mu \rangle \\
 &= \langle z, A\mu \rangle + \langle z, \mu \rangle \\
 &= \langle z, A\mu + \mu \rangle. \blacksquare
 \end{aligned}$$

In addition to the mean, the most important moment is the variance-covariance matrix of a random vector. Consider the random vector X in \mathbb{R}^n . Using the basis matrix $I_n = [e_1, e_2, \dots, e_n]$, then $\langle e_i, X \rangle = X_i$. Let $E\langle X_i - \mu_i, X_j - \mu_j \rangle = \Sigma_{ij}$, $i, j = 1, 2, \dots, n$, and $\Sigma = [\Sigma_{ij}]$. Let z and w lie in \mathbb{R}^n . Then

$$\begin{aligned}
 \text{Cov}\langle z, X \rangle; \langle w, X \rangle &= \text{Cov}\left(\sum_i z_i X_i, \sum_j w_j X_j\right) \\
 &= \sum_i \sum_j z_i w_j \text{Cov}(X_i, X_j) \\
 &= \sum_i \sum_j z_i w_j \Sigma_{ij} \\
 &= \langle z, \Sigma w \rangle.
 \end{aligned}$$

Alternatively, let $E \begin{pmatrix} \mathbf{h} \\ (\mathbf{X}_i - \bar{\mathbf{X}}) (\mathbf{X}_i - \bar{\mathbf{X}})^{\mathbf{i}} \end{pmatrix} = \mathcal{S}$. Then

$$\begin{aligned} \text{Cov } f(\mathbf{z}; \mathbf{X}); (\mathbf{w}; \mathbf{X})\mathbf{g} &= \text{Cov } f(\mathbf{z}; \mathbf{X}_i - \bar{\mathbf{X}}); (\mathbf{w}; \mathbf{X}_i - \bar{\mathbf{X}})\mathbf{g} \\ &= E f(\mathbf{z}; \mathbf{X}_i - \bar{\mathbf{X}}) : (\mathbf{X}_i - \bar{\mathbf{X}}; \mathbf{w})\mathbf{g} \\ &= E \begin{pmatrix} \mathbf{h} \\ \mathbf{z}^{\mathbf{h}} (\mathbf{X}_i - \bar{\mathbf{X}}) (\mathbf{X}_i - \bar{\mathbf{X}})^{\mathbf{i}} \end{pmatrix} \mathbf{w}^{\mathbf{i}} \\ &= \mathbf{z}^{\mathbf{h}} \mathcal{S} \mathbf{w}^{\mathbf{i}} \\ &= (\mathbf{z}; \mathcal{S} \mathbf{w}). \end{aligned}$$

The matrix \mathcal{S} (which is nnd or pd) is called the variance-covariance matrix of \mathbf{X} , or the dispersion of \mathbf{X} relative to $(; ;)$. In place of $\text{Cov } f(\mathbf{z}; \mathbf{X}); (\mathbf{w}; \mathbf{X})\mathbf{g}$ could be written $D f(\mathbf{z}; \mathbf{X}); (\mathbf{w}; \mathbf{X})\mathbf{g}$. The notation $D(\mathbf{X}) = \mathcal{S}$ is also used.

Definition 21 The unique non-negative definite matrix \mathcal{S} on $\mathbf{fV}; (; ;)\mathbf{g}$ to $\mathbf{fV}; (; ;)\mathbf{g}$ that satisfies

$$D f(\mathbf{z}; \mathbf{X}); \mathbf{w}; \mathbf{X}\mathbf{g} = (\mathbf{z}; \mathcal{S} \mathbf{w})$$

is called the dispersion of \mathbf{X} relative to $(; ;)$. ■

It should be emphasized that the dispersion depends on the scalar product specified. The next result shows how the dispersion changes as a function of the scalar product.

Proposition 9 Let X be a random vector which ranges over $fV; (; ;)g$ with dispersion ξ . Let $h; ; i$ be another scalar product given by $h; ; i = (; ; A;)$, A being a positive-definite matrix. The dispersion of X in $(V; h; ; i)$ is ξA .

Proof. It is necessary to establish that

$$D f h z ; X i ; h w ; X i g = h z ; \xi A w i$$

for all $z; w \in V$. Now

$$\begin{aligned} D f h z ; X i ; h w ; X i g &= D f (z ; A X) ; (w ; A X) g \\ &= D f (A z ; X) ; (A w ; X) g \\ &= (A z ; \xi A w) \\ &= (z ; A \xi A w) \\ &= h z ; \xi A w i . \blacksquare \end{aligned}$$

There are two consequences of Proposition 9. First, if $D(X)$ exists in one scalar product, it exists in all scalar products. Second, if $D(X) = \xi$ in $fV; (; ;)g$, then the dispersion relative to $h z ; w i = (z ; \xi^{-1} w)$ is the identity.

Another result which is a generalization of Proposition 9 is:

Proposition 10 Let X be a random vector which ranges over $fV; (; ;)g$. If A is the matrix of a linear transformation $V \rightarrow W$, where $fW; h; ; i g$ is a Euclidean subspace, then $D(A X + !) = A \xi A^>$ for all $! \in W$.

Proof.

$$\begin{aligned}
 & D f(hz; AX + ! i; hw; AX + ! i)g \\
 = & D f(hz; AXi + hz; ! i; hw; AXi + hw; ! i)g \\
 = & D f(hz; AXi; hw; AXi)g \\
 = & D \begin{matrix} \textcircled{i} & \textcircled{z} \\ A & X \end{matrix} ; \begin{matrix} \textcircled{i} & \textcircled{w} \\ A & X \end{matrix} \\
 = & \begin{matrix} \textcircled{i} & \textcircled{z} \\ A & X \end{matrix} ; \begin{matrix} \textcircled{i} & \textcircled{w} \\ A & X \end{matrix} \\
 = & \begin{matrix} - & \textcircled{z} \\ A & X \end{matrix} ; A \begin{matrix} \textcircled{i} & \textcircled{w} \\ A & X \end{matrix}
 \end{aligned}$$

Thus

$$D f(AX + !)g = A \begin{matrix} \textcircled{i} & \textcircled{w} \\ A & X \end{matrix}. \blacksquare$$

4.3 Special Covariance Structures

Relative to the Euclidean space $fV; (; ;)g$, the group of orthogonal transformations $V \rightarrow V$ will be denoted by $O(V)$.

Definition 22 A random vector X which ranges over $fV; (; ;)g$ with distribution Q has an orthogonally invariant, or spherical, distribution if its distribution and the distribution of MX are identical for all $B \in B_V$ and $M \in O(V)$. \blacksquare

An implication of a spherical distribution is that the corresponding dispersion is $I\lambda^2$ for some $\lambda^2 > 0$. Thus if X is a random vector ranging over $fV; (; ;)g$ with

spherical distribution, then $D(X) = I\lambda^2$ for some $\lambda^2 > 0$ and $D(MX) = D(X)$ for all $M \in O(V)$.

Definition 23 If $X \in \mathcal{FV}$; $(\cdot, \cdot)_g$ and

$$D(X) = I\lambda^2$$

for some $\lambda^2 > 0$, then X is said to have a weakly spherical distribution. This is equivalent to the condition $D(X) = D(MX)$ for all $M \in O(V)$. ■

The difference between a spherical and a weakly spherical distribution is that, in the former case, the whole distribution is invariant under orthogonal transformation whereas in the latter case, the dispersion is a positive multiple of the identity and this dispersion is invariant under orthogonal transformation as a consequence of Proposition 10.

4.4 The Multivariate Normal Distribution

In classical terminology, a probability distribution induced by a random variable or random vector is normally written as a distribution function (or cumulative density function). This was referred to as Q in section 4.1. If X is the random variable of interest with induced continuous distribution on \mathbb{R} then $P(X \leq x) = Q(x)$. The corresponding probability density function (pdf) in this case is $q(x)$ where $dQ(x) =$

$q(x) dx$. The interpretation of this last expression is $P\{x_i - \frac{1}{2} \leq X_i \leq x_i + \frac{1}{2}\} = \int_{x_i - \frac{1}{2}}^{x_i + \frac{1}{2}} q(x) dx$. Clearly $\int_{-\infty}^x q(t) dt = Q(x)$.

The multivariate normal distribution occupies a central place in the distribution theory now to be introduced. In elementary statistics, the univariate normal distribution is introduced as $X \sim N(\mu; \sigma^2)$, meaning that X is distributed normally with mean μ and variance σ^2 . In this case, the pdf corresponding to $X \sim N(\mu; \sigma^2)$ is

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}. \quad (4.1)$$

Here the integral

$$\int_{-\infty}^{\infty} g(x) dx = \int_{-\infty}^{\infty} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\} dx = \sigma\sqrt{2\pi}$$

and hence, setting $f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$, the integral $\int_{-\infty}^{\infty} f(x) dx = 1$, as required. If

a random sample of n is drawn from $N(\mu; \sigma^2)$, then the sample may be written as an n -tuple of mutually independent random variables $X_1; X_2; \dots; X_n$. If A_i represents the event $x_i - \frac{1}{2} \leq X_i \leq x_i + \frac{1}{2}$, $i = 1; 2; \dots; n$; then $P(A_1 \cap A_2 \cap \dots \cap A_n)$

is given by

$$\prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} \exp\left\{-\frac{(x_i - \mu)^2}{2\sigma^2}\right\} dx_1 dx_2 \dots dx_n. \quad (4.2)$$

If $I_n = \{x; x > [x_1; x_2; \dots; x_n]\}$, then the pdf (4.2) may be expressed more compactly as

$$f(x) dx = \frac{1}{(\sigma\sqrt{2\pi})^n} \exp\left\{-\frac{1}{2\sigma^2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right\} dx \quad (4.3)$$

where $f(x)$ now stands for $f(x_1) f(x_2) \dots f(x_n)$ of (4:1) and dx for $dx_1; dx_2 \dots dx_n$. If $X_1; X_2; \dots; X_n$ form a dependent set of normal random variables, such that for $i; j = 1; 2; \dots; n$

$$\int_{\mathbb{R}^n} x_i f(x) dx = EX_i = \mu_i$$

$$\int_{\mathbb{R}^n} (x_i - \mu_i)(x_j - \mu_j) f(x) dx = E(X_i - \mu_i)(X_j - \mu_j) = \sigma_{ij},$$

then Σ in (4:3) becomes $\Sigma = [\sigma_{ij}]$ and e^1, \dots, e^n being scalar in (4:3), is replaced by the vector μ , defined by $\mu^T = [\mu_1; \mu_2; \dots; \mu_n]$. Then the pdf for n dependent X_i 's is:

$$f(x) = (2\pi)^{-n/2} \det^{-1/2} \Sigma \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}. \quad (4.4)$$

In equation (4:4) Σ is presumed to be pd whereupon $\det \Sigma > 0$ and Σ^{-1} exists. In this case the random vector X has pdf according to (4:4) and this is indicated by writing $X \sim N(\mu; \Sigma)$.

Notice carefully that, in view of the positive definiteness of Σ ; there will exist a non-singular $(n \times n)$ matrix A such that $A \Sigma A^T = I_n$, $A^T A = \Sigma^{-1}$. Let

$$Z = A(X - \mu) \quad (4.5)$$

or

$$X = A^{-1}Z + \mu. \quad (4.6)$$

The Jacobian of the transformation (4:6) is

$$\begin{aligned} |J| &= \text{mod det} \frac{\partial (X)}{\partial (Z)} \\ &= \text{mod det}^{-1} A. \end{aligned}$$

Thus, corresponding to (4:4), the pdf of the vector Z is

$$g(z) = (2\pi)^{-\frac{n}{2}} \text{det}^{\frac{1}{2}} \Sigma \exp \left\{ -\frac{1}{2} z^T \Sigma^{-1} z \right\} |J|.$$

But $\text{det}^{-1} \Sigma = \text{det}^{-1} (A^T A) = \text{det}^{-2} A \Rightarrow \text{det}^{\frac{1}{2}} \Sigma = \text{mod det}^{-1} A$. Thus

$$\text{det}^{\frac{1}{2}} \Sigma |J| = \text{mod det}^{-1} A \text{mod det}^{-1} A = 1.$$

Hence

$$g(z) = (2\pi)^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2} z^T z \right\}. \quad (4.7)$$

Moreover from (4:5) $EZ = 0$ and $DZ = A \Sigma A^T = I_n$. Thus Z in (4:7) obeys $Z \sim N(0; I_n)$ and (4:5) represents the standardization of X; that is, the transformation of $X \sim N(\mu; \Sigma)$ to $Z \sim N(0; I_n)$. Notice that, for the normal distribution, independence and zero covariance are synonymous, Usually, independence implies zero covariance but not vice-versa. Also notice that the affine transformation (4:5) or (4:6) (meaning that the hyperplane in which X lies has a different origin from the hyperplane in which Z lies because $\mu \neq 0$) allows inheritance of normality. This is true of any affine transformation of a normal vector: an affine transformation of a normal variate is normal.

The standardized multivariate normal distribution of $Z \gg N(0; I_n)$ implies a set of n independent standardized normal variates $Z_i \gg N(0; 1)$ with $\text{cov}(Z_i; Z_j) = 0$ for all $i \neq j$. This is commonly written $Z_i \gg \text{NID}(0; 1) \quad i = 1; 2; \dots; n$, NID meaning normally and independently distributed. The distribution of Z plays an important role when the dispersion of X is singular, that is Σ is nnd. In this case, $\det \Sigma = 0$ and Σ^{-1} does not exist. When this is the case, the distribution of X is said to be degenerate. Moreover, (4:4) is ill defined. What can be made of a degenerate distribution? The secret to understanding degeneracy lies with the standardized normal distribution and a generalization of (4:6).

Consider m non-homogeneous linear functions of Z where $Z^> = [Z_1; Z_2; \dots; Z_n]$:

$$X = A^>Z + \mu^1. \quad (4.8)$$

In equation (4:8) $A^>$ is an $(m \times n)$ constant matrix, μ^1 is an $(m \times 1)$ vector and X is an $(m \times 1)$ vector comprising elements $X_1; X_2; \dots; X_m$. Since $Z \gg N(0; I_n)$, X has mean μ^1 and dispersion $A^>A = \Sigma$. If A has rank m , then $m < n$, and $A^>A = \Sigma$ is $(m \times m)$ and positive definite. In this case it is clearly in order to write $X \gg N(\mu^1; \Sigma)$. Now consider the case when the rank of A is $q < m$. In this case it is possible to show that there will exist a constant matrix $B^>$ of order $(q \times m)$ such that $B^>A^>AB = I_q$. Setting $W = B^>(x - \mu^1)$, W will have the $N(0; I_q)$ distribution. This implies that, if X has a degenerate distribution, in the sense that it has mean μ^1 and dispersion Σ of rank $q < m$ (m being the number of random variables comprising X), then X is

distributed on a q -dimensional subspace of the m -dimensional subspace on which X is defined.

Another degenerate case arises when $m > n$ and A has rank n . It is clear from (4:8) that $S = A^T A$ with rank $n < m$ and order $(m \times m)$. Thus X has a degenerate distribution. Nevertheless AA^T will be $(n \times n)$ and pd whereupon

$$i_{AA^T}^{-1} A(X_{i-1}) = Z \gg N(0; I_n)$$

and so the underlying distribution upon which X is based lies on an n -dimensional hyperplane, whereas X itself, being an m -tuple with $m > n$, appears to lie in R^m .

Commonly with a linear hypothesis on R^n , $y \gg N(1; I_n \frac{1}{2})$. Here the dispersion is non-singular and $1 \in L$, a k -dimensional subspace of R^n . The estimate of 1 is $\hat{1} = Py$ where P is the orthogonal projection on L of rank $k < n$. How is Py distributed? Evidently like $N(1; P \frac{1}{2})$, except that P is nnd and hence is singular. Thus Py has a degenerate distribution | not surprisingly, because $Py \in L$ and L has dimension k less than the number n of rows and columns of P . It is surely reasonable to expect that Py is distributed on L , even though it is expressed as if it is distributed in R^n (and is thereby degenerately distributed). In this case it is known that $P = BB^T$ where B^T is the $(k \times n)$ partial isometry such that $B^T B = I_k$. Using B^T ,

$$B^T (Py_{i-1}) = B^T (y_{i-1})$$

which has distribution $N(0; I_k \frac{1}{2})$. Thus Py has a degenerate distribution on R^n , but

a non-degenerate distribution on \mathbb{R}^k as anticipated.

The notation $X \in \mathbb{R}^n$ and $X \sim N(\mu; \Sigma)$ is taken to mean that there exists a $Z \in \mathbb{R}^q$ where q is the rank of Σ , and an $(n \times q)$ matrix A of rank q , such that

$$X = AZ + \mu$$

with $Z \sim N(0; I_q)$. This implies that $\Sigma = AA^T$ which is nnd. Under any linear transformation B , B being an $(m \times n)$ matrix of constants, if $X \sim N(\mu; \Sigma)$, $BX \sim N(B\mu; B\Sigma B^T)$.

4.5 The Non-central Chi-square Distribution

Definition 24 If $X \in \mathbb{R}^m$ has the $N(\mu; I_m)$ distribution, the random variable $X^T X$ has a non-central chi-square distribution with m degrees of freedom and non-centrality parameter $\lambda = \mu^T \mu$, written $\chi^2(m; \lambda)$. ■

The ordinary, or central, chi-square distribution is the special case of the non-central chi-square distribution when the non-centrality parameter is zero. If the distribution is said simply to be chi-square with m degrees of freedom, without specifying the non-centrality parameter, then it is understood that $\lambda = 0$. Notice that λ is a simple quadratic form in the mean of the present normal distribution i.e. $\lambda = \mu^T \mu$. Some authors define the non-centrality parameter as $\frac{\lambda}{\mu^T \mu}$ and others as $\frac{\lambda}{2}$. The important point is that if $\mu \neq 0$; $\lambda \neq 0$.

The mean and variance of $X \gg \hat{A}^2(m; \pm)$ are given by

$$EX = m + \pm$$

$$VX = 2m + 4\pm$$

Thus if $\pm = 0$, the mean is m and the variance is $2m$.

From definition 24, it follows that if $Q_1 \gg \hat{A}^2(m_1; \pm_1)$ independently of $Q_2 \gg \hat{A}^2(m_2; \pm_2)$, then $Q_1 + Q_2 \gg \hat{A}^2(m_1 + m_2; \pm_1 + \pm_2)$.

Theorem 11 If X ranges over \mathbb{R}^m according to the $N(1; \xi)$ distribution, ξ being pd, then

$$(X_i^{-1})^{\>} \xi^{-1} (X_i^{-1}) \gg \hat{A}^2(m).$$

Proof. There always exists an $(m \times m)$ non-singular matrix Q such that $Q\xi Q^{\>} = I_m$, $Q^{\>}Q = \xi^{-1}$. Let $X \gg N(1; \xi)$ as indicated. Then $Z = Q(X_i^{-1}) \gg N(0; I_m)$ implying that $Z^{\>}Z = (X_i^{-1})^{\>} Q^{\>}Q (X_i^{-1}) \gg \hat{A}^2(m)$. ■

Theorem 12 Let X range over \mathbb{R}^m according to $N(0; I_m)$ and let A be a fixed, symmetric matrix of order $(m \times m)$ and rank q . A necessary and sufficient condition for $X^{\>}AX \gg \hat{A}^2(q)$ is that $A = A^{\>} = A^2$.

Proof. Let $x \in \mathbb{R}^m$ be a realization of X : Clearly $x^{\>}x \gg \hat{A}^2(m)$ and

$$x^{\>}x = x^{\>}Ax + x^{\>}(I_m - A)x.$$

Let M be an orthogonal matrix such that $MM^T = D$, D being diagonal with q non-zero and $(m - q)$ zero elements. Since $M^{-1} = M^T$, $MM^T = I_n$ and

$$x^T M^T A M x + x^T M^T (I_m - A) M x = x^T x$$

or

$$x^T D x + x^T (I_m - D) x = x^T x.$$

D has $(m - q)$ zero diagonal elements, so the corresponding elements of $(I - D)$ must each be unity. Hence the non-zero elements of D must also be unity. Since a necessary and sufficient condition for A to be symmetric and idempotent is that its characteristic values are unity or zero, A must be symmetric and idempotent, and hence so must $(I_m - A)$.

On the other hand, assume $A = A^T = A^2$. Then there exists an orthogonal matrix M such that

$$x^T M^T A M x = x^T D x = x_1^T x_1$$

in which x_1 is a $(q \times 1)$ subvector of x . But $x_1 \in N(0; I_q)$. Hence $x_1^T x_1 \in \hat{A}^2(q)$. ■

Theorem 13 Let $x \in N(0; I_m)$. A necessary and sufficient condition for $x^T A_1 x$ and $x^T A_2 x$; of rank r_1 and r_2 respectively, to be independent $\hat{A}^2(r_1)$ and $\hat{A}^2(r_2)$ variates is that $A_1 A_2 = 0$. ■

4.6 Non-central F-distribution

Definition 25 If Q_1 and Q_2 are independent random variables and $Q_1 \sim \hat{A}^2(m_1; \pm)$, $Q_2 \sim \hat{A}^2(m_2)$, the distribution of the quotient $\frac{Q_1/m_1}{Q_2/m_2}$ is called the non-central F-distribution with m_1 and m_2 degrees of freedom and non-centrality parameter \pm , written $F(m_1; m_2; \pm)$. $F(m_1; m_2; 0)$ is the central F-distribution with m_1 and m_2 degrees of freedom, written $F(m_1; m_2)$. Occasionally, one encounters the doubly non-central F-distribution in which $Q_1 \sim \hat{A}^2(m_1; \pm_1)$ and $Q_2 \sim \hat{A}^2(m_2; \pm_2)$ independently. This is written $F(m_1; m_2; \pm_1; \pm_2)$. ■

4.7 Non-central t-distribution

Definition 26 If $Q_1 \sim \hat{A}^2(1; \pm_1)$ independently of $Q_2 \sim \hat{A}^2(m_2)$ then $\frac{Q_1}{Q_2/m_2} \sim F(1; m_2; \pm_1)$ and its square root $\frac{Q_1}{Q_2/m_2}^{1/2}$ has the non-central t-distribution with m_2 degrees of freedom and non-centrality parameter $\frac{\pm_1}{m_2}$. Clearly $t^2(m_2; \frac{\pm_1}{m_2}) = F(1; m_2; \pm_1)$. ■

4.8 Cochran's Theorem

Theorem 14 Let $X \sim N(0; I_m)$ and let $X^T X = \sum_{i=1}^k X^T A_i X$, the A_i being fixed matrices of rank r_i . The following conditions are equivalent (i.e. any one implies the

other two):

$$\mathbf{P}_{i=1}^{i=k} r_i = m;$$

Each $X^>A_iX \gg \hat{A}^2(r_i)$;

The $X^>A_iX$ are mutually independent. ■

Cochran's theorem is stated without proof. It is the central distribution theorem for the linear hypothesis.

5 Gauss-Markov Estimation and the Linear Hypothesis

5.1 The Gauss-Markov Theorem

The model to be considered is of a random vector y which ranges over \mathbb{R}^n on which the natural scalar product $(\cdot; \cdot)$ has been defined. For fixed w and z in \mathbb{R}^n ,

$$E(w; y) = (w; \beta) \quad \beta \in L \subseteq \mathbb{R}^n; \dim L = k < n; \quad (5.1)$$

$$D[(w; y); (z; y)] = (w; z) \sigma^2. \quad (5.2)$$

Thus $Ey = \beta \in L$, a k -dimensional subspace of \mathbb{R}^n , and $Dy = I_n \sigma^2$. A basis for L is available in the form of an $(n \times k)$ matrix X of rank k , and so β may be expressed as $\beta = X\gamma$, where γ is a $(k \times 1)$ vector which lies in some subset of \mathbb{R}^k . This subset of \mathbb{R}^k is easily shown to be \mathbb{R}^k . For suppose the rank of X is $r < k$, then in \mathbb{R}^k there are two orthogonal subspaces: \mathbb{R}^r of dimension r and $N[X]$ of dimension $(k - r)$. Thus γ may be uniquely decomposed as

$$\gamma = \gamma^* + (\gamma - \gamma^*)$$

with $\gamma^* \in \mathbb{R}^r$ and $(\gamma - \gamma^*) \in N[X]$. Now

$$X\gamma = X\gamma^* + 0.$$

because $N[X] = \{\mu : X\mu = 0; \mu \in \mathbb{R}^k\}$. Thus $X\gamma = X\gamma^*$. Moreover, if $r = k$, $N[X] = \{0\}$; whereupon $\gamma = \gamma^* \in \mathbb{R}^k$. In the model described above, $r = k$ and hence

$\beta \in \mathbb{R}^k$. Note also that

$$\mathbb{E} \beta = \mathbb{E} X'X^{-1} \sum_{i=1}^n h_i \epsilon_i$$

The model $y \in \mathbb{R}^n$, $\beta \in \mathbb{R}^k$, may also be written as

$$y = X\beta + \epsilon; \quad \mathbb{E}\epsilon = 0; \quad D\epsilon = I_n \quad (5.3)$$

with the X matrix considered as a fixed basis of L . In particular, for any $x \in \mathbb{R}^n$, $\mathbb{E}(x; \beta) = \mathbb{E}(x; y | \beta) = (x; 0) = 0$.

It is convenient and useful to consider the estimation of β in terms of the parameter vector β . Estimation of β is considered in terms of linear combinations of its elements.

Definition 27 A parametric function is a linear function of unknown parameters. If the unknown parameters are represented by the vector β in (5.3), then for a fixed vector $c \in \mathbb{R}^k$; $a = (c; \beta)$ is a parametric function. ■

The estimability of a is defined in terms of unbiased linear estimators.

Definition 28 The parametric function $a = (c; \beta)$ is said to be estimable if there exists a fixed vector $a \in \mathbb{R}^n$ such that the linear estimator $(a; y)$ is unbiased for a ; that is if $\mathbb{E}(a; y) = (c; \beta)$. ■

Theorem 15 The parameter function a is estimable if and only if there exists a fixed vector $a \in \mathbb{R}^n$ such that $c' = a'X$.

Proof. If α is estimable, then

$$E(\alpha; y) = (c; \bar{1}). \quad (5.4)$$

But $E(\alpha; y) = (\alpha; \bar{1})$ and $\bar{1} = X^-$. Hence

$$\alpha^>X^- = c^>\bar{1} \quad (5.5)$$

must hold identically in $\bar{1}$ implying $\alpha^>X = c^>$. If $c^> = \alpha^>X$, then (5:5) holds, which may be written as (5:4). Thus α is estimable if and only if $c^> = \alpha^>X$. ■

Lemma 16 If α is estimable, then there exists a unique, unbiased linear estimator $(\alpha^u; y)$ with $\alpha^u \in L$.

Proof. If α is estimable, there exists a vector a such that $E(\alpha; y) = (c; \bar{1})$ and $c^> = \alpha^>X$. Let $a = \alpha^u + (a_j - \alpha^u)$ with $\alpha^u \in L$; $(a_j - \alpha^u) \in L^\perp$. Then

$$(\alpha; y) = (\alpha^u; y) + (a_j - \alpha^u; y).$$

Taking expectations,

$$(\alpha; \bar{1}) = (\alpha^u; \bar{1}) + 0$$

because $(a_j - \alpha^u) \in L^\perp$ and $\bar{1} \in L$. Thus $E(\alpha; \bar{1}) = E(\alpha^u; \bar{1}) = (c; \bar{1})$.

To prove uniqueness, let the vector $\alpha^u \in L$ satisfy $E(\alpha^u; y) = (c; \bar{1})$. Then $0 = E(\alpha^u; y) - E(\alpha; y) = (\alpha^u - \alpha; \bar{1})$ implying, since $\bar{1} \in L$, that $(\alpha^u - \alpha) \in L^\perp$. But α^u and α both lie in L by construction. Thus $(\alpha^u - \alpha)$ lies in both L and L^\perp , implying $(\alpha^u - \alpha) = 0$, that is $\alpha^u = \alpha$ and α^u is shown to be unique. ■

Theorem 17 (The Gauss-Markov Theorem) Under the assumptions of model (5:1), every estimable function a has a unique linear estimator which has minimum variance in the class of unbiased linear estimators. The unique estimator is $a; \hat{1}^{\Delta} = c; \hat{\Delta}$ where $\hat{\Delta}$ is the least squares estimator of β in (5:3) and $\hat{1}^{\Delta} = X^{\Delta}$.

Proof. Let $(a; y)$ be an unbiased linear estimator of a , then $a = a^{\beta} + (a_j - a^{\beta})$; $a^{\beta} \in L; (a_j - a^{\beta}) \perp L^{\perp}$. Now $V(a; y) = k a k^2 + k (a_j - a^{\beta}) k^2$ and

$$k a k^2 = k a^{\beta} k^2 + k (a_j - a^{\beta}) k^2.$$

Since $k (a_j - a^{\beta}) k^2 \geq 0$ and $V(a^{\beta}; y) = k a^{\beta} k^2$,

$$V(a; y) \geq V(a^{\beta}; y).$$

But $(a; y)$ is any unbiased linear estimator a while $(a^{\beta}; y)$ is unique. Moreover, if P is the orthogonal projection matrix on L , $(a^{\beta}; y) = (P a; y) = (a; P y) = a; X^{\Delta} = X^{\Delta} a; \hat{\Delta} = c; \hat{\Delta}$, in view of theorem 15. This proves the Theorem. ■

A minimum variance unbiased linear estimator is often referred to as a best linear unbiased estimator or a BLUE.

5.2 The Linear Hypothesis

5.2.1 The Setting

To test a linear hypothesis, a distributional assumption is required. This is done by augmenting (5:3) with the maintained hypothesis

$$H : y \sim N^{i_1}(\bar{y}); I_n^{i_2} \quad (5.6)$$

$\bar{y} \in L$ as is required under (5:3). The null hypothesis to be tested may be written as r linearly independent equations

$$H_0 : (a_i; \bar{y}) = 0 \quad i = 1; 2; \dots; r; r < k; \bar{y} \in L. \quad (5.7)$$

Let A be the matrix $[a_1; a_2; \dots; a_r]$ which, by the linear independence of the a_i , is of rank r . Thus (5:7) may be re-written

$$H_0 : A\bar{y} = 0; \quad \bar{y} \in L, \quad (5.8)$$

whereupon it is clear that, if $L_0 = L \cap N(A)$, then

$$H_0 : \bar{y} \in L_0 \subseteq L, \quad (5.9)$$

L_0 being of dimension $(k - r)$.

The statement that $\bar{y} \in L$ is a linear hypothesis because L is a linear subspace. Similarly H_0 is a linear hypothesis because $L_0 \subseteq L$ is a linear subspace.

In (5:7){(5:9), H_0 is formulated as an hypothesis concerning \bar{y} as a function of the parameter vector β . The same hypothesis may also be expressed simply in terms

of β . As has been demonstrated in chapter 3, the a_i and hence the matrix A are not unique. What is unique is the matrix $PA = X'X^{-1}X'X^{-1}B$, where B is a $(k \times r)$ matrix of rank r such that

$$H_0 : B\beta = 0 \quad \beta \in \mathbb{R}^k \quad (5.10)$$

is consistent with (5:7)-(5:9).

The formulations of H_0 in (5:7)-(5:9) are in restraint form. In (5:10), the implied subspace in which β must lie on H_0 in $\mathbb{R}^k \setminus N(B)$, having dimension $(k - r)$. Moreover, $\dim N(B) = (k - r)$. Let a basis of $N(B)$ in \mathbb{R}^k be M a $(k \times (k - r))$ matrix satisfying

$$BM = 0.$$

It follows that if α is a $((k - r) \times 1)$ vector satisfying

$$\beta = M\alpha \quad (5.11)$$

and $\beta \in \mathbb{R}^k \setminus N(B)$; then (5:11) represents the freedom equation form of H_0 which appears in restraint form in (5:10).

The setting defined by (5:6) and one of (5:7)-(5:11) covers all of the known linear tests in econometrics, including the a_{\pm} ne hypothesis H_0^a :

$$H_0^a : B\beta = \mu \quad \mu \neq 0; \quad \beta \in \mathbb{R}^k. \quad (5.12)$$

Since $B^>$ has rank $r < k$, with r rows and k columns, it is clear that (5:12) has more than one solution. If any one of these is β , then $B^> \beta = \mu$ and

$$B^> (\beta - \mu) = 0. \quad (5.13)$$

The maintained hypothesis in (5:6) may now be re-formulated, using $y^* = y - X\beta$ and $\beta^* = \beta - \mu$, as

$$H^* : y^* \sim N(0, \sigma^2 I_n); \beta^* \in L. \quad (5.14)$$

If (5:14) replaces (5:6) and H_0 in (5:10) is replaced by

$$H_0^* : B^> \beta^* = 0 \quad (5.15)$$

where $\beta^* = (\beta - \mu)$, then (5:6) and the alternative hypothesis (5:12) has been accommodated into the linear hypothesis framework, (5:14) corresponding to (5:6), (5:12) corresponding to (5:10).

5.2.2 Decomposition

From chapter 3, the natural decomposition of R^n is

$$R^n = L_0 \oplus L_0^\perp \quad (5.16)$$

with $L = L_0 \oplus L_0^\perp$. The orthogonal projection from R^n on L_0 is P_0 and $PP_0 = P_0P = P_0$; also $(I_n - P)(I_n - P_0) = (I_n - P_0)(I_n - P) = (I_n - P); P(I_n - P) =$

since $(I_n - P)^{-1} = 0$ on H and on H_0 . It follows from this discussion that

$$F = \frac{k(P - P_0)y'k^2}{k(I_n - P)y'k^2} \cdot \frac{n - k}{r} \quad (5.19)$$

has the $F(r; n - k)$ -distribution, central on H_0 , non-central on H_a . Notice that both H_0 and H_a are consistent with H . In (5.19), the evidence that y has a 'large' component in $L_0^\perp \setminus L^\perp$ relative to its component in L^\perp is taken as counter to H_0 , whereupon F will be 'large'; and the evidence that y in $L_0^\perp \setminus L^\perp$ is 'small' relative to its component in L^\perp is taken as counter to H_a and hence not counter to H_0 , whereupon F will be 'small' and H_0 'acceptable' as a working hypothesis.

5.2.3 Alternative Forms of the F-test

Various alternative expressions for F in (5.19) may be developed.

$$\begin{aligned} F &= \frac{y'(P - P_0)y \cdot (n - k)}{y'(I_n - P)y \cdot r} \\ &= \frac{y'(I_n - P_0)y - y'(I_n - P)y \cdot (n - k)}{y'(I_n - P)y \cdot r} \end{aligned} \quad (5.20)$$

whereupon, if $\hat{\beta} = (I_n - P)y$ and $\hat{\beta}_0 = (I_n - P_0)y$, the unrestricted and restricted estimators of β respectively, then

$$F = \frac{\hat{\beta}'\hat{\beta}_0 - \hat{\beta}'\hat{\beta} \cdot (n - k)}{\hat{\beta}'\hat{\beta} \cdot r} \quad (5.21)$$

Equation (5.21) is easily calculated directly from a computer output by noting the restricted and unrestricted sums of squared errors.

ization of the Wald statistic. In fact, the test is based on the Wald Principle which is a principle for developing statistics to test restrictions on a parameter vector.

The Wald Principle is quite general. Let μ be a vector-valued parameter and $\hat{\mu}_n$ a consistent unrestricted estimator of it, based on a random sample of n observations. As is commonly the case for sufficiently large n , $\sqrt{n}(\hat{\mu}_n - \mu) \rightarrow N(0; \Sigma)$, Σ being a positive definite-dispersion. The null hypothesis, $h(\mu) = 0$, is an $(r \times 1)$ vector of functions $h_i(\mu)$; $i = 1, 2, \dots, r$. When n is large and there is sufficient regularity, $\sqrt{n}h(\hat{\mu}_n) \rightarrow N(0; H_\mu \Sigma H_\mu')$ where $H_\mu = \frac{\partial h(\mu)}{\partial \mu}$. Then

$$\sqrt{n}h(\hat{\mu}_n) \xrightarrow{d} N(0; H_\mu \Sigma H_\mu') \quad \sqrt{n}h(\hat{\mu}_n) \xrightarrow{d} \hat{A}^2(r; \pm)$$

with $\pm = 0$ on the null hypothesis: $h(\mu) = 0$. The null hypothesis may be tested using the statistic

$$W = nh(\hat{\mu}_n)' H_{\hat{\mu}_n}^{-1} \hat{\Sigma}_{\hat{\mu}_n}^{-1} h(\hat{\mu}_n) \xrightarrow{d} \hat{A}^2(r; \pm)$$

because $H_{\hat{\mu}_n}$ and $\hat{\Sigma}_{\hat{\mu}_n}$ will be 'close' to H_μ and Σ , for large n . In the case of the linear hypothesis examined above, (5:21) gives the form of W with $\sqrt{n}h(\hat{\mu}_n) = \sqrt{n}B\hat{\mu}_n$ and $B = \frac{\partial h(\mu)}{\partial \mu}$. In all cases, only unrestricted estimates are applied to determine W . Thus the Wald Principle is to be seen as a principle for developing tests of general hypotheses, of the kind $h(\mu) = 0$, on the basis of unrestricted, asymptotically normal estimators $\hat{\mu}_n$ of μ , using a standardized quadratic form in $h(\hat{\mu}_n)$.

There are now three forms of the F-test of (5:18) for a linear hypothesis: (5:20), (5:21) and (5:23). The form (5:20) is a difference-in-regression formulation: $\frac{y' (P_i - P_0) y}{r \mathbb{K}^2}$; (5:21) is a difference-in-residual-sum-of-squares formulation: $\frac{u_i' - u_i^{*k}}{r \mathbb{K}^2}$; and (5:23) is the Wald formulation based on unrestricted estimation of β and \mathbb{K}^2 . All three forms are identical. A fourth form is an equivalent test involving a different distribution. This comes through the following result.

Theorem 18 Let y be a random vector ranging over \mathbb{R}^n according to (5:6). Then a test for H_0 in (5:8) is $F = \frac{y' (P_i - P_0) y}{r \mathbb{K}^2}$ where $\mathbb{K}^2 = \frac{y' (I_{n-i} - P) y}{(n-i-k)}$, P being the orthogonal projection on L and P_0 being the orthogonal projection $L_0 \perp L \perp \mathbb{R}^n$, L and L_0 having dimensions k and $(k-i-r)$. A test equivalent to F is $M = \frac{y' (P_i - P_0) y}{y' (P_i - P_0) y + y' (I_{n-i} - P_0) y}$ or

$$M = \frac{y' (P_i - P_0) y}{y' (I_{n-i} - P_0) y}. \quad (5.24)$$

This has the $F_{i-r, n-i-k}$ distribution, central on H_0 , non-central otherwise. ■

Theorem 18 is stated without proof. Both the F- and M-tests are exact, small-sample tests; one is just another form of the other. Thus F and M will never yield conflicting results. As will now be demonstrated, M is a test based on the Lagrange Multiplier Principle of testing restrictions.

5.2.4 The Lagrange Multiplier Principle

Consider the Lagrange expression

$$L(\beta; \lambda) = (y - X\beta)'(y - X\beta) + 2\lambda'AX\beta; \quad \beta \in L$$

in which λ is the Lagrange Multiplier enforcing the restriction $AX\beta = 0$. The first-order conditions for minimizing $(y - X\beta)'(y - X\beta)$ subject to $AX\beta = 0; \beta \in L$, are the stationary conditions for $L(\beta; \lambda)$:

$$-2X'y + 2X'X\beta + 2A'\lambda = 0$$

$$AX\beta = 0$$

where $\hat{\beta}$ and $\hat{\lambda}$ are the restricted estimates of β and λ . Since $\beta \in L$ and $\lambda \in L_0 \perp L$

$$P_X y - P_X X\hat{\beta} = P_X A'\hat{\lambda} \tag{5.25}$$

$$A'X\hat{\beta} = (A'X)\hat{\beta} = 0. \tag{5.26}$$

Premultiplying equation (5:25) by A' and noting that $A'X\hat{\beta} = A'P_X\hat{\beta} = 0$, then $A'Py = A'PA'\hat{\lambda}$ and

$$\hat{\lambda} = (A'PA')^{-1}A'Py. \tag{5.27}$$

Also note that, since we know that $\beta = P_0y$, (5:25) may also be written

$$(P - P_0)y = P_X A'\hat{\lambda}$$

implying, from (5:27), that

$$(P_i - P_0)y = PA^i A^>PA^{\zeta_i - 1} A^>Py$$

as is to be expected. From (5:27)

$$\hat{\zeta}_i \approx N^{-3} \mathbf{E} A^>PA^{\alpha_i - 1} A^>1; \mathbf{E} A^>PA^{\alpha_i - 1} \zeta_i^2$$

which on H_0 is $N^{-3} \mathbf{E} A^>PA^{\zeta_i - 1} \zeta_i^2$ since then $A^>1 = 0$. Hence

$$\frac{y^>PA^i A^>PA^{\zeta_i - 1} A^>PA^i A^>PA^{\zeta_i - 1} A^>Py}{\zeta_i^2} \approx \hat{A}^2(r; \pm).$$

Substituting in the estimator of ζ_i^2 on H_0

$$\zeta_i^2 = \frac{y^>(I_{n_i} - P_0)y}{(n_i - k + r)}$$

there emerges

$$\begin{aligned} LM &= \hat{\zeta}_i^{-3} \mathbf{h}^3 \hat{\zeta}_i^{-1} \zeta_i^2 \\ &= \frac{y^>(P_i - P_0)y}{\frac{y^>(I_{n_i} - P_0)y}{(n_i - k + r)}} \end{aligned}$$

which is the large-sample Lagrange Multiplier statistic for testing H_0 . Moreover

$$\frac{LM}{(n_i - k + r)} = M$$

and M is seen to be a small-sample specialization of the Lagrange Multiplier statistic.

5.3 Applications of the F-test

In considering various examples of the F-test, it is worth noting that, of the various expressions (5:20), (5:21) and (5:23), it is usually (5:23) that is easiest to apply to find specific theoretical results, while (5:21) is the most practical form. Sometimes it is beneficial to adopt the freedom-equation approach.

5.3.1 Testing the Difference Between Two Means

The situation to be examined is that of two independent samples of n observations, in the $(n \times 1)$ vector y_1 , and m observations in the $(m \times 1)$ vector y_2 . The maintained hypothesis is written in terms of $y = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}$ an $((n+m) \times 1)$ vector. The common population mean of the first sample is μ_1 and its variance is σ^2 ; the common population mean and variance for the second sample are μ_2 and σ^2 . Thus the variances are the same. Under H_0

$$y \sim N \left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}; I_{n+m} \sigma^2 \right)$$

where $\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$ and L has basis

$$E = \begin{pmatrix} I_n & 0 \\ 0 & I_m \end{pmatrix}$$

e_n being the equiangular vector on R^n and e_m being the equiangular vector on R^m .

Thus $L = R[E]$ and

$$\mu = E^{-1} \mu$$

where $\bar{y} = [\bar{y}_1 \ \bar{y}_2]$; \bar{y} being (2×1) . In this formulation, then, the X-matrix is replaced by E. It is desired to test

$$H_0: \bar{y}_1 = \bar{y}_2 \quad [1 \ -1] \bar{y} = 0.$$

Thus $B = [1 \ -1]$ and $r = 1$. On H_0 , $\bar{y}_1 = \bar{y}_2 = \bar{y}$ and

$$1 = e_{n+m}.$$

Let $P_n = e_n e_n' + e_{n+1} e_{n+1}' + \dots + e_{n+i-1} e_{n+i-1}'$ with P_m the same using e_m ; using e_{n+m} , the orthogonal projection onto its range is P_0 . Thus

$$P = \begin{bmatrix} P_n & 0 \\ 0 & P_m \end{bmatrix}$$

and

$$F = \frac{y' (P - P_0) y}{y' (I_n - P) y} \cdot \frac{n + m - 1}{1}.$$

$$y' (P - P_0) y = n(\bar{y}_1 - \bar{y})^2 + m(\bar{y}_2 - \bar{y})^2$$

$$= n\bar{y}_1^2 + m\bar{y}_2^2 - (n + m)\bar{y}^2$$

where \bar{y}_1 ; \bar{y}_2 are the sample means and $(n + m)\bar{y} = n\bar{y}_1 + m\bar{y}_2$. Hence

$$(n + m)\bar{y}^2 = \frac{n^2\bar{y}_1^2 + m^2\bar{y}_2^2 + 2nm\bar{y}_1\bar{y}_2}{n + m}$$

and

$$y' (P - P_0) y = \frac{nm(\bar{y}_1 - \bar{y}_2)^2}{n + m}.$$

Now

$$\begin{aligned} \frac{y' (I_{n+m} - P) y}{n + m - 2} &= \frac{(n - 1) \sigma_1^2 + (m - 1) \sigma_2^2}{(n + m - 2)} \\ &= \sigma^2. \end{aligned}$$

It follows that

$$\begin{aligned} F &= \frac{nm (\bar{y}_1 - \bar{y}_2)^2}{\frac{n+m}{\sigma^2}} \\ &= \frac{(\bar{y}_1 - \bar{y}_2)^2}{\frac{\sigma^2}{n} + \frac{\sigma^2}{m}}. \end{aligned}$$

$F \gg F(1; n + m - 2)$ centrally on H_0 , non-centrally otherwise. Hence

$$t = \frac{\bar{y}_1 - \bar{y}_2}{\sqrt{\frac{\sigma^2}{n} + \frac{\sigma^2}{m}}}$$

which has the $t(n + m - 2)$ -distribution, centrally on H_0 non-centrally otherwise.

5.3.2 Testing a Regression Coefficient

In the linear regression (5:3) and (5:6), the null hypothesis is $H_0: \beta_i = 0$, β_i being the i 'th element of β . In this case $B^{\beta} = \beta_i = 0$ and

$$B^{\beta} = \begin{pmatrix} 0 & 0 & \dots & 0 & 1 & 0 & \dots & 0 & 0 \end{pmatrix}$$

with unity in the i 'th position. Clearly $B^{\beta} X' X^{-1} X' y = \hat{\beta}_i$ and $B^{\beta} X' X^{-1} B = m^{-1}$ the $(i; i)$ 'th element of $X' X^{-1}$. Since $r = 1$, applying (5:23)

$$F = \frac{\hat{\beta}_i^2}{m^{-1} \sigma^2} \gg F(1; n - k).$$

Hence

$$t = \frac{\hat{\beta}_i}{SE_{\hat{\beta}_i}}$$

where $SE_{\hat{\beta}_i} = \sqrt{\frac{\sigma^2}{n} m^{ii}}$ and $t \gg t(n_i - k)$, centrally on H_0 , non-centrally otherwise.

When the null is $\beta_i = \beta_i^0$, then the t-test becomes

$$t = \frac{\hat{\beta}_i - \beta_i^0}{SE_{\hat{\beta}_i}}$$

5.3.3 Testing a Linear Combination of Regression Coefficients

As an example, consider, in the context of (5:3) with (5:6), $H_0 : \beta_2 + \beta_3 = 1$ as a null hypothesis to test. In this case, $B^{\beta} = 1$ with $B^{\beta} = [0 \ 2 \ 3 \ 0 \ 0 \ \dots \ 0]$.

Let β_0 be a solution to $B^{\beta} \beta_0 = 1$, e.g. $\beta_0 = [0 \ 1 \ 1 \ 0 \ 0 \ \dots \ 0]$. Then

$B^{\beta}(\beta - \beta_0) = 0$; $r = 1$ and

$$\frac{1}{\sigma^2} F = B^{\beta} \hat{\beta} (\beta - \beta_0)^{\beta} \frac{1}{\sigma^2} = \frac{1}{\sigma^2} \left(2\hat{\beta}_2 + 3\hat{\beta}_3 - 1 \right)^2 \frac{1}{4m^{22} + 12m^{23} + 9m^{33}}$$

where m^{ij} is the $(i; j)$ 'th element of $(X^{\beta})^{-1}$. F will follow the $F(1; n_i - k)$ distribution and hence

$$t = \frac{2\hat{\beta}_2 + 3\hat{\beta}_3 - 1}{\sqrt{\frac{\sigma^2}{n} (4m^{22} + 12m^{23} + 9m^{33})}} \gg t(n_i - 1)$$

since $m^{23} = m^{32}$ centrally on H_0 , non-centrally otherwise. The denominator of t is $SE_{2\hat{\beta}_2 + 3\hat{\beta}_3}$.

5.3.4 Testing a Block of Regression Coefficients

Again using (5:3) and (5:6) as the maintained hypothesis, consider the special case

$$y = X_1 \beta_1 + X_2 \beta_2 + \epsilon$$

in which X_1 has k_1 columns and X_2 has k_2 columns. The null hypothesis is $H_0 :$

$\beta_2 = 0$. In this case, $L = R[X]$ where $X = [X_1 : X_2]$, $L_0 = R[X_1]$ and hence

$$P_0 = X_1 (X_1' X_1)^{-1} X_1' = P_1. \text{ Now}$$

$$\begin{aligned} P &= F_1 + F_2 \\ &= X_1 (X_1' X_1)^{-1} X_1' M_2 X_2 (X_2' X_2)^{-1} X_2' M_1 + X_2 (X_2' X_2)^{-1} X_2' M_1. \end{aligned}$$

Then

$$\begin{aligned} P - P_0 &= (I - P_0) P \\ &= (I - P_1) P \\ &= M_1 P \\ &= M_1 F_2 \\ &= M_1 X_2 (X_2' X_2)^{-1} X_2' M_1. \end{aligned}$$

Thus

$$\begin{aligned}
 F &= \frac{y' (P_i - P_0) y}{y' (I_n - P) y} \cdot \frac{n_i - k}{r} \\
 &= \frac{y' (P_i - P_1) y}{y' (I_n - P) y} \cdot \frac{n_i - k}{k_2} \\
 &= \frac{y' M_1 X_2' X_2 M_1 X_2' X_2 M_1 y}{k_2 \mathbb{A}^2}
 \end{aligned}$$

Moreover

$$k_2 F = \frac{y' M_1 X_2' X_2 M_1 X_2' X_2 M_1 y}{\mathbb{A}^2}.$$

F here is distributed as $F(k_2; n_i - k)$ and $k_2 F$ is approximately $\hat{A}^2(k_2)$ for 'large' n .

5.3.5 Durbin-Hausman Testing

Referring to section 5.3.5, where $\beta_2 = 0$, the maintained model is

$$y = X_1 \beta_1 + X_2 \beta_2 + \epsilon \tag{5.28}$$

and, on H_0 , the model is

$$y = X_1 \beta_1 + \epsilon. \tag{5.29}$$

It will be convenient to write $\beta_1 = \beta_1^*$ in (5:29) and to test

$$y = X_1 \beta_1^* + \epsilon \tag{5.30}$$

as a model related to (5:28) when $\beta_1 = \beta_1^*$. This permits a re-interpretation of H_0 from $H_0 : \beta_2 = 0$ to $H_0^* : \beta_1 - \beta_1^* = 0$ where β_1 is taken to be from (5:28). A natural

way to test H_0^a is to consider a standardized quadratic form based upon the difference

$$\begin{aligned}
 \Delta_{1i}^a &= \mathbf{y}'_1 \mathbf{M}_2 \mathbf{X}_1 \mathbf{C}_{i-1}^{-1} \mathbf{X}'_1 \mathbf{M}_2 \mathbf{y}_i - \mathbf{y}'_1 \mathbf{X}_1 \mathbf{C}_{i-1}^{-1} \mathbf{X}'_1 \mathbf{y}_i \\
 &= \mathbf{y}'_1 \mathbf{M}_2 \mathbf{X}_1 \mathbf{C}_{i-1}^{-1} \mathbf{X}'_1 \mathbf{M}_2 \mathbf{y}_i - \mathbf{y}'_1 \mathbf{X}_1 \mathbf{C}_{i-1}^{-1} \mathbf{X}'_1 \mathbf{X}_1 \mathbf{C}_{i-1}^{-1} \mathbf{X}'_1 \mathbf{y}_i \\
 &= \mathbf{y}'_1 \mathbf{M}_2 \mathbf{X}_1 \mathbf{C}_{i-1}^{-1} \mathbf{X}'_1 \mathbf{M}_2 (\mathbf{I}_i - \mathbf{P}_1) \mathbf{y}_i \\
 &= \mathbf{y}'_i \mathbf{X}'_1 \mathbf{M}_2 \mathbf{X}_1 \mathbf{C}_{i-1}^{-1} \mathbf{X}'_1 \mathbf{P}_2 \mathbf{M}_1 \mathbf{y}_i.
 \end{aligned} \tag{5.31}$$

Clearly, (5.31) is zero whenever L_1 and L_2 are orthogonal; to proceed it is presumed L_1 is not orthogonal to L_2 . Let $A = \mathbf{X}'_1 \mathbf{M}_2 \mathbf{X}_1$ and $\mathbf{P}_2 \mathbf{X}_1 = \mathbf{X}_2 \mathbf{Q}$, $\mathbf{Q} = \mathbf{X}'_2 \mathbf{X}_2 \mathbf{C}_{i-1}^{-1} \mathbf{X}_2' \mathbf{X}_1$. A is always non-singular. When $k_1 = k_2$, \mathbf{Q} is square and non-singular and

$$\Delta_{1i}^a = \mathbf{y}'_i \mathbf{A}^{-1} \mathbf{Q}' \mathbf{X}'_2 \mathbf{M}_1 \mathbf{y}_i \gg N_{i-1} \mathbf{y}'_i \mathbf{A}^{-1} \mathbf{Q}' \mathbf{X}'_2 \mathbf{M}_1 \mathbf{y}_i; \mathbf{A}^{-1} \mathbf{Q}' \mathbf{X}'_2 \mathbf{M}_1 \mathbf{X}_2 \mathbf{Q} \mathbf{A}^{-1} \mathbf{I}_{i-1} \mathbf{C}_{i-1}^{-1} \tag{5.32}$$

with $\mathbf{y}'_i \mathbf{A}^{-1} \mathbf{Q}' \mathbf{X}'_2 \mathbf{M}_1 \mathbf{y}_i = 0$ when $\mathbf{y}_i \in L_1$, i.e. when H_0 and H_0^a coincide. Thus, given L_1 and L_2 are not orthogonal,

$$\begin{aligned}
 \Delta_{1i}^a &\gg \mathbf{y}'_i \mathbf{A}^{-1} \mathbf{Q}' \mathbf{X}'_2 \mathbf{M}_1 \mathbf{X}_2 \mathbf{Q} \mathbf{A}^{-1} \mathbf{I}_{i-1} \mathbf{C}_{i-1}^{-1} \mathbf{y}_i \\
 &= \frac{\mathbf{y}'_i \mathbf{M}_1 \mathbf{X}_2 \mathbf{Q} \mathbf{A}^{-1} \mathbf{A} \mathbf{Q}' \mathbf{X}'_2 \mathbf{M}_1 \mathbf{X}_2 \mathbf{C}_{i-1}^{-1} \mathbf{X}'_2 \mathbf{M}_1 \mathbf{y}_i}{\mathbf{I}_{i-1} \mathbf{C}_{i-1}^{-1} \mathbf{A}^{-1} \mathbf{A} \mathbf{Q}' \mathbf{X}'_2 \mathbf{M}_1 \mathbf{y}_i} \\
 &= \frac{\mathbf{y}'_i \mathbf{M}_1 \mathbf{X}_2 \mathbf{X}'_2 \mathbf{M}_1 \mathbf{X}_2 \mathbf{C}_{i-1}^{-1} \mathbf{X}'_2 \mathbf{M}_1 \mathbf{y}_i}{\mathbf{I}_{i-1} \mathbf{C}_{i-1}^{-1}} \gg \hat{A}^2(k_2; \pm)
 \end{aligned} \tag{5.33}$$

with $\pm = 0$ when $\mathbf{y}_i \in L_1$. Moreover,

$$\frac{\mathbf{y}'_i \mathbf{M}_1 \mathbf{X}_2 \mathbf{X}'_2 \mathbf{M}_1 \mathbf{X}_2 \mathbf{C}_{i-1}^{-1} \mathbf{X}'_2 \mathbf{M}_1 \mathbf{y}_i}{\mathbf{I}_{i-1} \mathbf{C}_{i-1}^{-1}} \gg F(k; n_i - k) \tag{5.34}$$

where \mathbb{A}^2 is calculated from (5:28). Thus the testing of $H_0^a : \beta_1 = 0$, when $k_1 = k_2$ and L_1 is not orthogonal to L_2 , is seen to lead to precisely the same test as $H_0 : \beta_2 = 0$. Although $k_1 = k_2$ seems to be a very special case, in fact it is quite usual in specification testing using instrumental variables (IV's).

When $k_2 < k_1$, (5:32) still holds except that now the dispersion is singular, being of order $(k_1 - k_2)$ and rank k_2 . The easiest way then to proceed is to standardize the distribution of β_1 on L_2 . Note that Q is $(k_2 \times k_1)$ of rank k_2 , whereupon QQ' has a unique inverse. Also let C be a non-singular matrix of order $(k_2 \times k_2)$ such that $CX_2'M_1X_2C' = I_{k_2}$. Then

$$C'QQ'Q^{-1}QA^{-1}\beta_1 = CX_2'M_1y \sim N(0; I_{k_2}\mathbb{A}^2)$$

when $\beta_1 \in L_1$ or when H_0^a holds and L_1 is not orthogonal to L_2 . It follows that

$$\frac{\beta_1'QA^{-1}Q'Q^{-1}C'CX_2'M_1X_2C'Q^{-1}QA^{-1}\beta_1}{\mathbb{A}^2} \sim \hat{A}^2(k_2; \pm)$$

or

$$\begin{aligned} & \frac{y'M_1X_2QA^{-1}Q'Q^{-1}C'CX_2'M_1X_2C'Q^{-1}QA^{-1}Q'X_2'M_1y}{\mathbb{A}^2} \\ &= \frac{y'M_1X_2X_2'M_1X_2C'X_2M_1y}{\mathbb{A}^2} \sim \hat{A}^2(k_2; \pm) \end{aligned}$$

i.e. (5:33), and hence (5:34). Thus when $k_1 \geq k_2$ and L_1 lies strictly oblique to L_2 , then an F-test for H_0 is equivalent to an F-test for H_0^a . In Durbin-Hausman testing, the difference in estimators is not between unrestricted and restricted estimators, but

rather between an IV estimator of the form $\hat{\beta} = (X'AX)^{-1}X'AY$, where A is an orthogonal projection matrix, and the least squares estimator of the same coefficient $\hat{\beta} = (X'X)^{-1}X'Y$. This

$$\hat{\beta} = (X'AX)^{-1}X'AY$$

$M = I_n - P$, with dispersion

$$D_{\hat{\beta}} = (X'AX)^{-1}X'AMAX(X'AX)^{-1} \frac{1}{k}$$

The standardized quadratic form is

$$\frac{Y'AMAX(X'AX)^{-1}X'AY}{\frac{1}{k}} = \frac{Y'AMAX(X'AX)^{-1}X'AY}{\frac{1}{k}}$$

with corresponding F-statistic given by

$$F = \frac{Y'AMAX(X'AX)^{-1}X'AY}{k}$$

where $\frac{1}{k}$ has been calculated from the equation

$$y = X\beta + \epsilon$$

Under appropriate conditions, $F \approx F(k; n - k)$. Thus the testing of $H_0: \beta = 0$ is the same as testing $\beta = 0$ or $AX\beta = 1_2 = 0$; 1_2 otherwise belonging to a k -dimensional subspace of $R[A]$. Notice that AX is of the same order and rank as X .

5.3.7 Testing for Structural Change

Consider two time periods comprising n_1 and n_2 observations which follow the linear regression

$$y_1 = X_1\beta_1 + \epsilon_1 \quad (5.35)$$

$$y_2 = X_2\beta_2 + \epsilon_2$$

or

$$y = X\beta + \epsilon \quad (5.36)$$

where $y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$, $X = \begin{bmatrix} X_1 & 0 \\ 0 & X_2 \end{bmatrix}$, $\beta = \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}$ and

$$X = \begin{bmatrix} 2 & 3 \\ 6 & 7 \\ 4 & 5 \\ 0 & X_2 \end{bmatrix}$$

In this situation β_1 and β_2 are presumed to be different. The null hypothesis is $H_0 : \beta_1 = \beta_2 = \beta_0$ and it is assumed that both n_1 and n_2 are greater than k . The model under H_0 is

$$y = X_0\beta_0 + \epsilon \quad (5.37)$$

where

$$X_0 = \begin{bmatrix} 2 & 3 \\ 6 & 7 \\ 4 & 5 \\ X_2 \end{bmatrix}$$

β_0 having half the number of elements as β . Notice that H_0 may be written

$$I_k - \frac{1}{k} \mathbf{1} \mathbf{1}' = 0$$

or

$$M^{-1} = \frac{1}{k} \begin{bmatrix} 2 & 3 \\ 6 & 7 \\ 4 & 5 \end{bmatrix} = M^{-1} \quad (5.38)$$

Thus using (5:38) in (5:36)

$$\begin{aligned} y &= XM^{-1} + \epsilon \\ &= X_0 \beta_0 + \epsilon \end{aligned}$$

where $X_0 = XM$ in (5:37).

The orthogonal projection matrix onto $R[X]$ is

$$P = \frac{1}{k} \begin{bmatrix} P_1 & 0 \\ 0 & P_2 \end{bmatrix}$$

where $P_i = X_i' X_i^{-1} X_i$; $i = 1, 2$; $P_0 = X_0' X_0^{-1} X_0$. Then the F-statistic is

$$F = \frac{y' (P_1 - P_0) y}{k \sigma^2}$$

where $\sigma^2 = \frac{y' (I_n - P) y}{(n_1 + n_2 - 2k)} = \frac{n_1 \sigma_1^2 + n_2 \sigma_2^2}{(n_1 + n_2 - 2k)}$. It is an assumption of equation (5:36) that

$\epsilon_1 \gg N(0; I_{n_1} \sigma^2)$ and $\epsilon_2 \gg N(0; I_{n_2} \sigma^2)$, so that ϵ_1 and ϵ_2 have a common mean 0 and a common variance σ^2 .

If the two regressions contain constant terms, it may be desirable to consider an unrestricted model (5:36) in which every coefficient may differ from one period to the next, and a null in which corresponding slope coefficients are the same but the constants differ. More generally, it may be desirable to consider restrictions on only a subset of coefficients. For example, let the two regressions from (5:36) be

$$\begin{aligned}
 y &= \begin{matrix} 2 \\ 6 \\ 4 \end{matrix} \begin{matrix} X_{11} & X_{12} & 0 & 0 \\ 0 & 0 & X_{21} & X_{22} \end{matrix} \begin{matrix} 3 \\ 6 \\ 7 \\ 5 \\ 6 \\ 6 \\ 6 \\ 4 \end{matrix} \begin{matrix} - \\ - \\ - \\ - \\ - \\ - \\ - \\ - \end{matrix} \begin{matrix} 11 \\ 12 \\ 21 \\ 22 \end{matrix} \begin{matrix} 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 5 \end{matrix} + \epsilon \\
 &= X_0 \beta_0 + \epsilon,
 \end{aligned} \tag{5.39}$$

the columns of X_{11} and X_{21} referring to the same variables, the columns of X_{12} and X_{22} referring to the same or possibly different variables. The null hypothesis is $H_0 : \beta_{11} = \beta_{21} = \beta_1$ and the null model is then

$$\begin{aligned}
 y &= \begin{matrix} 2 \\ 6 \\ 4 \end{matrix} \begin{matrix} X_{11} & X_{12} & 0 & 0 \\ X_{21} & 0 & X_{22} & 0 \end{matrix} \begin{matrix} 3 \\ 6 \\ 7 \\ 5 \\ 6 \\ 6 \\ 6 \\ 4 \end{matrix} \begin{matrix} - \\ - \\ - \\ - \\ - \\ - \\ - \\ - \end{matrix} \begin{matrix} 1 \\ 12 \\ 22 \end{matrix} \begin{matrix} 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 7 \\ 5 \end{matrix} + \epsilon \\
 &= X_0 \beta_0 + \epsilon
 \end{aligned} \tag{5.40}$$

in an obvious notation. P is the orthogonal projection for (5:39) and

$$P = \begin{matrix} 2 & 3 \\ 6 & 0 \\ 4 & 7 \\ & 5 \\ & 0 \\ & P_2 \end{matrix}$$

having ranks of, on the left-hand side, $(n_1 - k_1 - k_2)$ and, on the right-hand side, $(n_1 - k_1)$ and $(n_2 - k_2)$; $n = n_1 + n_2$. Thus \mathcal{Y}_1^2 and \mathcal{Y}_2^2 , each divided by \mathcal{Y}^2 , are independent \hat{A}^2 -variates, and

$$F = \frac{y_1^2 (I_{n_1 - k_1 - k_2}) y_1}{y_2^2 (I_{n_2 - k_2}) y_2} \cdot \frac{n_2 - k_2}{n_1 - k_1},$$

and $F \sim F(n_1 - k_1 - k_2; n_2 - k_2)$. F is suitable for testing the null hypothesis $H_0 : \mathcal{Y}_1^2 = \mathcal{Y}_2^2 = \mathcal{Y}^2$. This should be done prior to any test for structural change since $\mathcal{Y}_1^2 = \mathcal{Y}_2^2$ is a condition for application of the F -distribution.

5.3.8 Recursive Residuals

The aim with so-called recursive residuals is to calculate residuals sequentially to facilitate the testing of the constancy of a regression. Thus a regression may be fitted to the first k observations, then to $k + 1$; $k + 2$, and so on, until all n observations are used. At each stage, the additional observation is tested to determine whether it is significantly different from the previous regression. To formulate the various tests, the following notation is used.

X represents the $(n \times k)$ matrix of n observations on k explanatory variables, the first k observations being represented by X_k . The orthogonal projection onto $R[X_k] = X_k (X_k^T X_k)^{-1} X_k^T = P_k = I_k$ since X_k is $(k \times k)$ of rank k . If r observations are added to X_k there results X_{k+r} and hence $P_{k+r} = X_{k+r} (X_{k+r}^T X_{k+r})^{-1} X_{k+r}^T$. Using

1. For any two distinct values of r , say m and q ; $m < q$, $\hat{P}_q \hat{P}_m = \hat{P}_q$ and hence

$$\begin{aligned} \sum_{i=1}^3 \hat{P}_{q,i} \hat{P}_{q+1,i} - \sum_{i=1}^3 \hat{P}_{m,i} \hat{P}_{m+1,i} &= \hat{P}_{q,i} \hat{P}_{q,i} \hat{P}_{q+1,i} + \hat{P}_{q+1,i} \\ &= 0. \end{aligned}$$

Thus all the quadratic forms are naturally independent and the ranks of the right-hand side quadratic forms sum to $k + (k + 1) + k + (k + 2) + \dots + (k + 1) + \dots + (n - 1) + (n - 2) + \dots + n + (n - 1) = k + 1 + 1 + \dots + 1 + 1 = k + (n - k) = n =$ the rank of the left-hand side quadratic form. Division by $\frac{1}{2}$ ensures that the left-hand side is $\hat{A}^2(n; \pm)$ while each of the terms on the right-hand side is an independent $\hat{A}^2(1; \pm_r)$ except the first which is $\hat{A}^2(k; \pm_r)$, $r = 1; 2; \dots; (n - k)$ and $\xi_{\pm_r} = \pm$.

2. Let $y^> \hat{P}_{q,i} \hat{P}_{q+1,i} y = y^> I_{n,i} \hat{P}_{q+1,i} y_i$ $y^> I_{n,i} \hat{P}_q y = S_{q+1,i}^2 S_q^2 = U_q^2$; U_q is called the q 'th recursive residual.

3. Moving forward from the k 'th observation,

$$\begin{aligned} \frac{y^> \hat{P}_{1,i} \hat{P}_{2,i} y}{y^> I_{n,i} \hat{P}_{2,i} y} &\gg -\frac{\mu_1}{2}; 1, \\ \frac{y^> \hat{P}_{2,i} \hat{P}_{3,i} y}{y^> I_{n,i} \hat{P}_{3,i} y} &\gg -\frac{\mu_1}{2}; \frac{3}{2}, \\ \frac{y^> \hat{P}_{3,i} \hat{P}_{4,i} y}{y^> I_{n,i} \hat{P}_{4,i} y} &\gg -\frac{\mu_1}{2}; 2, \end{aligned}$$

and so on; these statistics represent ways of testing whether the $(k + r)$ 'th observation is significantly different from the preceding $(k + r - 1)$ observations, in respect of the linear regression taking the general form of $y = X\beta + \epsilon$. These tests are based on the Lagrange Multiplier Principle. Alternatively an F-test may be used at each stage, beginning with

$$\frac{y' \hat{P}_1 y}{y' I_n y} \gg F(1; 1).$$

If this test does not reject $H_0 : \beta_{k+1} = \beta_{k+2} = 0$ (where β_{k+1} refers to the first $k + 1$ observations and β_{k+2} to the first $k + 2$ observations) then the numerator and denominator are added to yield $y' \hat{P}_2 y$ and the next test is

$$\frac{y' \hat{P}_2 y}{y' I_n y} \gg F(1; 2).$$

Again if this test does not reject $H_0 : \beta_{k+2} = \beta_{k+3} = 0$ then we proceed to add $y' \hat{P}_2 y + y' \hat{P}_3 y = y' \hat{P}_3 y$ and use

$$\frac{y' \hat{P}_3 y}{y' I_n y} \gg F(1; 3)$$

and so on. The method is referred to as "adding insignificant sums of squares" (as is done at each stage).

A full discussion of the tests and associated tests is found in R. L. Brown, J. Durbin and J. M. Evans, "Techniques for testing the constancy of relationships over

time" (with discussion), *Journal of the Royal Statistical Society, Series B*, 37(1975), 149-92.

6 Limits, Continuity and Convergence

6.1 Introduction

The discussion in this chapter is decidedly terse. The aim is to provide a minimum of understandable material, sufficient to tackle the econometric problems that arise later on. Many of these problems are provided with a solution that rests on reasonable assumptions about the data being analyzed. To understand the assumptions it is necessary to understand certain definitions. The definitions may appear forbidding, while in fact being natural to the task at hand. To ease the way forward, each definition is followed by an example. A good reference in the material covered (and many extensions) is White (1984).

6.2 Real Numbers

Definition 29 Let $\{b_n\}$ be a sequence of real numbers, $n = 1; 2; 3; \dots$. If for every real $\epsilon > 0$ there exists an integer $N(\epsilon)$ and a real number b such that

$$|b_n - b| < \epsilon \quad \text{for all } n \geq N(\epsilon)$$

then b is said to be the limit of $\{b_n\}$ and $\{b_n\}$ is said to converge to b as $n \rightarrow \infty$; this is written $b_n \rightarrow b$ as $n \rightarrow \infty$ or $\lim_{n \rightarrow \infty} b_n = b$. ■

Example 9 Let $b_n = b + \frac{1}{n}$; $n = 1; 2; 3; \dots$. Then $b_1 = b + 1$; $b_2 = b + \frac{1}{2}$; $b_3 = b + \frac{1}{3}$; \dots . Clearly, $\lim_{n \rightarrow \infty} b_n = b$.

Example 10 Consider the expansion of $(1 + \frac{a}{n})^n$:

$$(1 + \frac{a}{n})^n = 1 + n \frac{a}{n} + \frac{n(n-1)}{2!} \frac{a^2}{n^2} + \dots + \frac{n(n-1)\dots(n-r+1)}{r!} \frac{a^r}{n^r} + \dots$$

As $n \rightarrow \infty$

$$(1 + \frac{a}{n})^n = 1 + a + \frac{a^2}{2!} + \frac{a^3}{3!} + \dots = e^a.$$

Hence $\lim_{n \rightarrow \infty} (1 + \frac{a}{n})^n = e^a$.

Turning to vectors of real elements, let b_n be a $(k \times 1)$ vector with typical element b_n^i ($i = 1; 2; \dots; k$). Let $b_n^i \rightarrow b^i$ as $n \rightarrow \infty$ for each i . Then $\lim_{n \rightarrow \infty} b_n = b$ where b has elements b^i . Matrices are treated similarly, except it is usually necessary to make an assumption about the properties of the limiting matrix. For example, if X is an $(n \times k)$ real matrix, then $X^T X = \sum_{q=1}^k x_{iq} x_{jq}$ $i; j = 1; 2; \dots; k$; $q = 1; 2; \dots; n$. If X has rank k , then $X^T X$ has rank k and is positive definite. Let

$$\frac{1}{n} X^T X \rightarrow M_{XX} \quad \text{as } n \rightarrow \infty$$

where $\frac{1}{n} \sum_{q=1}^k x_{iq} x_{jq} \rightarrow m_{ij}$. Commonly it will be assumed that the m_{ij} are finite and M_{XX} is a positive-definite matrix.

6.3 Continuity

Definition 30 Let $g(\cdot)$ be a vector-valued function $\mathbb{R}^p \rightarrow \mathbb{R}^k$ and let $b \in \mathbb{R}^p$. The function $g(\cdot)$ is continuous at b if for any sequence $\{b_n\}$ whose limit as $n \rightarrow \infty$ is b ,

$$\lim_{n \rightarrow \infty} g(b_n) = g(b).$$

Equivalently, $g(\cdot)$ is continuous at b if, for every $\epsilon > 0$, there exists a $\delta(\epsilon) > 0$ such that for $a \in \mathbb{R}^p$ and $|a_i - b_i| < \delta(\epsilon); i = 1, 2, \dots, p; |g_j(a) - g_j(b)| < \epsilon; j = 1, 2, \dots, k$. If $B \subset \mathbb{R}^p$, then g is continuous on B if it is continuous for every point of B . ■

Example 11 If $a_n \rightarrow a$ and $b_n \rightarrow b; a_n + b_n \rightarrow a + b$ and $a_n b_n \rightarrow ab$.

Example 12 The matrix inverse function is continuous at every point that represents a non-singular matrix. Thus if $\frac{1}{n}X \rightarrow X \in M_{XX}$, then $\frac{1}{n}X^{-1} \rightarrow X^{-1} \in M_{XX}^{-1}$, so long as M_{XX} is finite and non-singular.

6.4 The Order of a Sequence

Definition 31 The sequence $\{b_n\}$ is at most of order $n^{-\phi}$, written $O(n^{-\phi})$ if, for some finite real number ϕ there exists an integer N such that $|n^{-\phi} b_n| < \phi$ for all $n \geq N$. ■

Definition 32 The sequence $\{b_n\}$ is of order smaller than $n^{-\phi}$, written $o(n^{-\phi})$, if, for every $\pm > 0$, there exists an $N(\pm)$ such that $|n^{-\phi} b_n| < \pm$ for all $n \geq N(\pm)$. ■

Clearly $fb_{ng} = O(n^a)$ if $n^i \cdot b_n$ is eventually bounded; and $b_n = o(n^a)$ if $n^i \cdot b_n \rightarrow 0$. If fb_{ng} is $o(n^a)$, then it must be $O(n^a)$. Also, if $fb_{ng} = O(n^a)$, then for every $\epsilon > 0$, $fb_{ng} = o(n^{a+\epsilon})$. If $fb_{ng} = O(n^0) = O(1)$, then it is eventually bounded and it may or may not have a limit. If $fb_{ng} = o(1)$, then $b_n \rightarrow 0$.

Example 13 Let $\sum x_i = O(n)$. Then $\frac{1}{n} \sum x_i = O(1)$ and this implies that, for all $\epsilon > 0$; $\sum x_i = o(n^{1+\epsilon})$.

Example 14 Let $b_n = 4 + 3n + 6n^2$. Then $fb_{ng} = O(n^2)$ and $fb_{ng} = o(n^{2+\epsilon})$ for all $\epsilon > 0$. ■

In regard to matrices and vectors, these are of $O(n^a)$ or $o(n^a)$ if each element is $O(n^a)$ or $o(n^a)$.

6.5 Almost Sure Convergence

It is often of interest in statistics generally, and hence in econometrics, to consider averages of sequences of random variables $\{Y_i\}$. Let Ω represent the entire random sequence $\{Y_i\}$. Let $b_n(\Omega) = \frac{1}{n} \sum_{i=1}^n Y_i$.

Definition 33 Let $\{b_n(\Omega)\}$ represent a sequence of random variables. $b_n(\Omega)$ converges almost surely to b , written $b_n(\Omega) \xrightarrow{a.s.} b$ if there exists a real number b such that $P_\Omega [b_n(\Omega) \rightarrow b] = 1$ where the probability measure P_Ω refers to the distribution of Ω .

which determines the joint distribution function for the entire sequence of random variables. ■

What $P\{b_n \rightarrow 1\} = 1$ means is that some sequences in $\{b_n\}$ may not converge, but these have probability measure zero.

If $\{b_n\} \rightarrow b$ and $g(\cdot)$ is continuous at b , then $g(b_n) \rightarrow g(b)$.

Example 15 If $y = X\beta + \epsilon$ is the familiar least squares model in n observations and k variables, and these k variables are random variables, then let $\frac{1}{n}X'X \rightarrow M_{XX}$ a finite pd matrix and $\frac{1}{n}X'\epsilon \rightarrow 0$. It follows that $\frac{1}{n}X'X^{-1}X'y = \hat{\beta}$ may be written as $\frac{1}{n}X'X^{-1}X'y = \frac{1}{n}X'X^{-1}X'(X\beta + \epsilon)$. Thus $\hat{\beta} = \beta + \frac{1}{n}X'X^{-1}X'\epsilon \rightarrow \beta + M_{XX}^{-1}0 = \beta$.

Example 16 Consider now $\hat{\beta} = (X'AX)^{-1}X'AY$. If it is assumed that $\frac{1}{n}X'AX \rightarrow M_A$ which is finite and positive definite and $\frac{1}{n}X'AY \rightarrow 0$, then $\hat{\beta} \rightarrow \beta$ by the same reasoning as above. ■

6.6 Convergence in Probability

Definition 34 If there exists a real number b such that, for every $\epsilon > 0$,

$$P\{|b_n - b| < \epsilon\} \rightarrow 1 \quad \text{as } n \rightarrow \infty,$$

then b_n is said to converge in probability to b , which is written $b_n \xrightarrow{P} b$ or $P\text{-}\lim_{n \rightarrow \infty} b_n = b$. If $b_n \xrightarrow{P} b$ and $g(\cdot)$ is continuous at b , then $g(b_n) \xrightarrow{P} g(b)$.

g(b). Let $g(\cdot)$ be a vector-valued function $\mathbb{R}^p \rightarrow \mathbb{R}^k$ and $b_n(\cdot)$ and b be $(p \times 1)$ vectors. If $g(\cdot)$ is continuous at b , then $b_n(\cdot) \rightarrow^p b$ implies that $g_j(b_n(\cdot)) \rightarrow^p g_j(b)$ $i = 1; 2; \dots; p; j = 1; 2; \dots; k$. ■

6.7 Convergence in Distribution

From elementary statistics it is known that if $X \sim (1; \frac{1}{4^2}); 0 < \frac{1}{4^2} < 1$, and $x_1; x_2; \dots; x_n$ is a sequence of random selections from this distribution then, by the central limit theorem, for sufficiently large n ,

$$\frac{\sum_{j=1}^n x_j - n}{\sqrt{n}} \rightarrow N(0; 1)$$

or $P_n(x_j - 1) \rightarrow N(0; \frac{1}{4^2})$.

Definition 35 Let $\{F_n\}$ be a sequence of random variables with joint distribution function $fF_n(\cdot)$. If $F_n(y) \rightarrow F(y)$ as $n \rightarrow \infty$ for every point y , where $F(y)$ is the distribution function of the random variable Y , then $\{F_n\}$ converges in distribution to the distribution of the random variable Y which is written $b_n \rightarrow^d Y$. In this definition vector may be interchanged with variable. ■

Let y range over \mathbb{R}^p according to a distribution with mean \bar{y} , a $(p \times 1)$ vector, and dispersion Σ , a $(p \times p)$ pd matrix. This is written $y \sim (1; \Sigma)$. Let $\{y_n\}$ be a random sequence of n from $(1; \Sigma)$. Then $y_1; y_2; \dots; y_n$ are n $(p \times 1)$ vectors each

with mean μ and as $n \rightarrow \infty$

$$\sqrt{n} \left(\bar{y} - \mu \right) \xrightarrow{d} N(0, \Sigma)$$

or

$$\sqrt{n} \bar{y} \xrightarrow{d} N(\sqrt{n}\mu, \Sigma)$$

where $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ is a $(p \times 1)$ vector of sample means. By an abuse of notation, it is common to write

$$\bar{y} \xrightarrow{d} N\left(\mu, \frac{1}{n}\Sigma\right).$$

Let μ be a parameter vector in \mathbb{R}^k and $\hat{\mu}_n$ an estimate of μ based upon a random sample of n such that $\hat{\mu}_n \xrightarrow{d} N\left(\mu, \frac{1}{n}\Sigma\right)$. Let $h(\mu)$ be a set of $r < k$ functions of μ which are real and continuous at μ . Then

$$h(\hat{\mu}_n) \xrightarrow{d} N\left(h(\mu), n^{-1}H_\mu^> H_\mu^<\right),$$

where $H_\mu^> = \frac{\partial h(\mu)}{\partial \mu}$ evaluated at μ . Now $\sqrt{n}(\hat{\mu}_n - \mu) \xrightarrow{d} N(0; \Sigma)$ and hence $p\text{-lim } \hat{\mu}_n = \mu$. Moreover, by Taylor's theorem,

$$h(\hat{\mu}_n) = h(\mu) + \frac{\partial h(\mu)}{\partial \mu} \sqrt{n}(\hat{\mu}_n - \mu) + R$$

where $R \rightarrow 0$ as $\hat{\mu}_n \rightarrow \mu$. Hence

$$\sqrt{n} \left(h(\hat{\mu}_n) - h(\mu) \right) = \sqrt{n} \frac{\partial h(\mu)}{\partial \mu} (\hat{\mu}_n - \mu) + o_p(1) \xrightarrow{d} N\left(0; H_\mu^> \Sigma H_\mu^<\right).$$

On the null hypothesis $h(\mu) = 0$,

$$h(\hat{\mu}_n) \stackrel{d}{\approx} N(\mathbf{0}; n^{-1} H_{\mu}^{-1} \Sigma H_{\mu}).$$

Assuming that $\hat{\mu} \xrightarrow{p} \mu$; $H_{\hat{\mu}} \xrightarrow{p} H_{\mu}$; and $\hat{\Sigma} \xrightarrow{p} \Sigma$ then

$$H_{\hat{\mu}} \hat{\Sigma} H_{\hat{\mu}} \xrightarrow{p} H_{\mu} \Sigma H_{\mu}$$

whereupon, if $h(\mu) = 0$, then

$$W = n h(\hat{\mu})' \hat{\Sigma}^{-1} h(\hat{\mu}) \xrightarrow{d} \chi^2(r; 0).$$

7 Maximum-likelihood Estimation Procedures and Associated Tests of Significance

7.1 The General Problem

There are n independent observations x_1, x_2, \dots, x_n from a known probability density function $f(x; \mu)$ and it is also known that the hypothesis $H: \mu \in \Omega$ applies, Ω being some subset (or all) of some k -dimensional Euclidean space \mathbb{R}^k . It is desired to test whether the true value of μ , denoted by μ_0 , is an element of Ω , a $(k - r)$ -dimensional subspace of \mathbb{R}^k . There are two ways of expressing Ω when \mathbb{R}^k is known: in the form of restraint equations, $H_0: h(\mu) = 0$, where h represents r independent functions, or in the form of freedom equations, written $H_0: \mu = \mu^{(a)}$, in which $\mu^{(a)}$ is a set of $(k - r)$ independent parameters: $\mu^{(a)} \in \mathbb{R}^{k-r}$. It is sometimes the case that a combination of restraint and freedom equations is specified. There is always a freedom specification corresponding to a restraint equation specification and vice-versa; but the corresponding relationship is often difficult to derive in practice.

Using the restraint form, $\Omega = \{\mu: h(\mu) = 0; \mu \in \mathbb{R}^k\}$. Using the freedom equation specification, $\Omega = \{\mu: \mu = \mu^{(a)}; \mu^{(a)} \in \mathbb{R}^{k-r}\}$.

7.2 General Justification

There are essentially two general justifications for using maximum-likelihood methods.

The first comes from Bayes' Theorem which states that

$$P(\mu \in \Omega \mid x_1; x_2; \dots; x_n) / P(x_1; x_2; \dots; x_n \mid \mu \in \Omega) P(\mu \in \Omega) \quad (7.1)$$

The probability that $\mu \in \Omega$, given that the sample $x_1; x_2; \dots; x_n$ has been observed, is proportional to the probability that $x_1; x_2; \dots; x_n$ will be observed, given that μ is indeed an element of Ω , times the probability that μ lies in Ω . $P(x_1; x_2; \dots; x_n \mid \mu \in \Omega)$ is in principle a probability but, since the x_i 's are observed and so are to be treated as parameters, and $\mu \in \Omega$ is to be treated as variable, $P(x_1; x_2; \dots; x_n \mid \mu \in \Omega)$ is called the likelihood corresponding to n observations, written $L_n(x; \mu)$. $P(\mu \in \Omega)$ is called the prior, or a priori, probability since it represents the probability that μ is an element of Ω before the sample $x_1; x_2; \dots; x_n$ has been observed. $P(\mu \in \Omega \mid x_1; x_2; \dots; x_n)$ is the probability that μ is an element of Ω , given that the sample $x_1; x_2; \dots; x_n$ has been observed; hence it is called the posterior or a posteriori probability.

$P(\mu \in \Omega)$ can only be assigned subjectively and this is not universally acceptable. However, it can be assumed that $P(\mu \in \Omega)$ is uniformly distributed over Ω , that is, that the investigator cannot distinguish, in terms of prior beliefs, any values of μ that are more likely than others. Such uniform distribution of prior probabilities over Ω is known as the Doctrine of Equal Ignorance. Applying this doctrine to finding the

most likely value of μ in Ω , there results

$$\sup_{\mu \in \Omega} P(\mu | x_1, x_2, \dots, x_n) = \sup_{\mu \in \Omega} \{ L_n(x; \mu) \}; \quad (7.2)$$

$\{$ being the product of the constant of proportionality from (7.1) and the constant probability covering every $\mu \in \Omega$. Thus maximizing the likelihood for variations in $\mu \in \Omega$ yields the same value of μ as maximizing $P(\mu | x_1, x_2, \dots, x_n)$ and the principle of maximum-likelihood is given a justification as a consequence of (i) the calculus of probability and (ii) the Doctrine of Equal Ignorance.

A second justification for maximum-likelihood methods is found in the properties that maximum-likelihood estimators possess under general conditions. Given that the probability density function, from which that n sample observations have been chosen, is known then the maximum-likelihood estimator of $\mu \in \Omega$, say $\hat{\mu} \in \Omega$, will generally have the following properties.

1. $\hat{\mu}$ is consistent.
2. $\hat{\mu}$ is a BAN (Best Asymptotic Normal) estimator, that is, $\sqrt{n}(\hat{\mu} - \mu_0) \xrightarrow{d} N(0; B_{\mu_0}^{-1})$ where $B_{\mu_0}^{-1}$ is the $(k \times k)$ positive-definite matrix which attains the Cramér-Rao lower bound, to be discussed below.
3. In view of property 2, $\hat{\mu} \approx N(\mu_0; \frac{1}{n} B_{\mu_0}^{-1})$ for large n , and hence standardized quadratic forms in $\hat{\mu}$, and in functions of $\hat{\mu}$ like $h(\hat{\mu})$ for example, are \hat{A}^2 -

distributed. Thus if $H_\mu = \frac{\partial h(\mu)}{\partial \mu}$ then, given sufficient regularity,

$$nh(\hat{\mu} - H_{\hat{\mu}}^{-1} H_{\hat{\mu}} \hat{\mu}) \xrightarrow{d} \chi^2(r; \pm)$$

and $\pm = 0$ on $H_0 : \mu \in \Theta_0$.

4. $\hat{\mu}$ is a function of a sufficient statistic for μ , if a sufficient statistic for μ exists.
5. Maximum-likelihood estimators are invariant to reparametrization in the following sense. Let μ and ζ both be $(k \times 1)$ vectors such that for $g : \Theta \rightarrow \mathbb{R}^k$ is a smooth mapping that transforms $\mu \in \Theta$ uniquely into $g(\mu) = \zeta \in T$. Let the new likelihood based on the same n observations be $L_n^g(x; \zeta)$. Then $L_n(x; \mu) = L_n^g(x; \zeta)$. Now $L_n(x; \hat{\mu}) > L_n(x; \mu)$ for all $\mu \in \Theta$. It follows that $L_n^g(x; \hat{\zeta}) = L_n^g(x; g(\hat{\mu})) = L_n(x; \hat{\mu}) > L_n(x; \mu) = L_n^g(x; \zeta)$ for all $\mu \in \Theta$ and all $\zeta \in T$. Thus the maximum-likelihood estimator $\hat{\zeta} = g(\hat{\mu})$.

7.3 Notation

$L_n(x; \mu) = \prod_{i=1}^n f(x_i; \mu)$ represents the likelihood corresponding to n observations. $\log L_n(x; \mu)$ is the loglikelihood corresponding to n observations and

$$\frac{1}{n} \log L_n(x; \mu) = \frac{1}{n} \sum_{i=1}^n \log f(x_i; \mu) \quad (7.3)$$

denotes the loglikelihood corresponding to a single observation. A justification for this terminology is based on the following argument. $f(x; \mu)$ is the probability density

function for the single observation x and

$$B_{\mu} = E_{\mu} \left[\frac{\partial^2 \log f(x; \mu)}{\partial \mu \partial \mu'} \right] \quad (7.4)$$

is a $(k \times k)$ symmetric matrix called the information matrix. B_{μ} clearly applies to a single observation. Now

$$\begin{aligned} \frac{1}{n} E_{\mu} \left[\frac{\partial^2 \log L_n(x; \mu)}{\partial \mu \partial \mu'} \right] &= \frac{1}{n} \sum_{i=1}^n E_{\mu} \left[\frac{\partial^2 \log f(x_i; \mu)}{\partial \mu \partial \mu'} \right] \\ &= \frac{1}{n} n B_{\mu} \\ &= B_{\mu}. \end{aligned}$$

Thus $\frac{1}{n} \log L_n(x; \mu)$ is justified as the loglikelihood corresponding to a single observation because it contains precisely the amount of information contained in a single observation, as given in (7.4). That (7.4) may be interpreted as measuring information contained in a single observation needs some justification; this follows in section 7.5 below (Cramér-Rao Inequality).

$\hat{\mu}$ will represent the maximum-likelihood estimator of μ in \mathcal{X} ; μ_0 the maximum-likelihood estimator of μ in \mathcal{I} .

The symbol $z(\mu)$ is used for $E_{\mu_0} [\log f(x; \mu)]$, that is

$$z(\mu) = \int \log f(x; \mu) f(x; \mu_0) dx = E_{\mu_0} [\log f(x; \mu)] \quad (7.5)$$

E_{μ_0} denoting expectation corresponding to the distribution of the random variable x when the true value of μ is μ_0 .

H_μ (as defined in section 7.2 above), $z(\mu)$ and B_μ may be evaluated at various values of μ . At $\hat{\mu}$ for example, H_μ , $z(\mu)$ and B_μ become $H_{\hat{\mu}}$, $z(\hat{\mu})$ and $B_{\hat{\mu}}$. $\hat{\mu}$ is the point in Ω at which $z(\mu)$ is greatest.

7.4 Regularity

Certain regularity conditions will be required. Those given below are not the weakest that could be used, but they are fairly simple and sufficient for the purposes of this exposition.

A.1 Potential derivatives with respect to μ of integrals over the sample space of functions of x and μ are the same as the integrals over the sample space of the partially differentiated functions with respect to μ .

The main implications of assumption A.1 are that the range of x should be independent of μ and that the tails of the distribution of x should allow convergence of the differentiated integral.

A.2 $z(\mu)$ exists for all $\mu \in \Omega$.

The whole problem of maximum-likelihood estimation, restricted or unrestricted, is closely bound up with the behavior of the function z . This is because the Law of Large Numbers ensures that, for each μ , the sequence $\frac{1}{n} \log L_n(x; \mu)$ converges, for almost all x , to $z(\mu)$. If this convergence is uniform with respect to μ , then for large n and most x , $\frac{1}{n} \log L_n(x; \mu)$ will be uniformly near z and, under appropriate

conditions, will reach its supremum in \mathcal{I} near the point where z attains its supremum in \mathcal{I} . The assumptions given below are designed to achieve this situation.

A.3 \mathcal{I} is a convex compact subset of \mathbb{R}^k .

A.4 $\log f(x; \cdot)$ is continuous on \mathcal{I} for almost all $x \in \mathbb{R}^n$.

A.5 For almost all $x \in \mathbb{R}^n$ and for every $\mu \in \mathcal{I}$, $\frac{\partial \log f(x; \mu)}{\partial \mu_i}$ $i = 1, 2, \dots, k$ exists and its absolute value is less than some function $g(x)$ whose expectation, relative to the distribution at μ_0 , is finite.

A.6 The function h is continuous on \mathcal{I} .

A.7 There exists a point $\mu^* \in \mathcal{I}$ such that $z(\mu^*) > z(\mu)$ when $\mu \in \mathcal{I}$ and $\mu \neq \mu^*$.

Assumptions A.3{A.5 ensure that, for almost all x , the sequence $\frac{1}{n} \log L_n(x; \mu)$ converges to $z(\mu)$ uniformly with respect to $\mu \in \mathcal{I}$. Assumptions A.3 and A.6 ensure that \mathcal{I} is a compact subset of \mathbb{R}^k and consequently that any continuous function on \mathcal{I} attains its supremum at some point in \mathcal{I} . In particular, the function $\frac{1}{n} \log L_n(x; \mu)$ attains its supremum in \mathcal{I} at some point $\hat{\mu}(x)$ in \mathcal{I} , for almost all x , the sequence of values of $\hat{\mu}(x)$ as n increases converges to μ^* . If $\mu_0 \in \mathcal{I}$, then usually μ_0 will satisfy the conditions required of μ^* , i.e. $z(\mu_0) > z(\mu)$ if $\mu \neq \mu_0$.

A.8 μ^* is an interior point of \mathcal{I} .

Assumption A.8 ensures that for large n and most x , $\hat{\mu}(x)$ will be an interior point of \mathcal{I} and consequently will emerge as a solution to the restricted maximum-likelihood equations when the function h is differentiable.

The method by which the asymptotic distribution of maximum-likelihood estimates is usually derived, involves expanding the likelihood function by Taylor's Theorem. To adopt this method requires three more assumptions.

A.9 The functions h_i ($i = 1; 2; \dots; r$) possess both first- and second-order partial derivatives which are continuous on \mathcal{R}^n .

A.10 For almost all $x \in \mathcal{R}^n$ the function $\log f(x; \cdot)$ possesses continuous second-order partial derivatives in a neighborhood of μ^0 . If μ belongs to this neighborhood, then the absolute value of $\frac{\partial^2 \log f(x; \mu)}{\partial \mu_i \partial \mu_j}$ ($i = 1; 2; \dots; k$) is less than a function of x whose expectation relative to the distribution of x at μ_0 is finite.

A.10 ensures that $\frac{\partial^2 z(\mu)}{\partial \mu_i \partial \mu_j}$ at μ^0 exists and that the sequence of second-order partial derivatives of $\log f(x; \mu)$ converges for almost all x to $\frac{\partial^2 z(\mu)}{\partial \mu_i \partial \mu_j}$ at $\mu = \mu^0$.

A.11 For almost all $x \in \mathcal{R}^n$ the function $f(x; \cdot)$ possesses third-order partial derivatives in the neighborhood of μ^0 and if μ is in this neighborhood then the absolute value of $\frac{\partial^3 \log f(x; \mu)}{\partial \mu_i \partial \mu_j \partial \mu_k}$ ($i; j; k = 1; 2; \dots; k$) is less than some function of x whose expectation relative to the distribution of x at μ_0 is finite.

A.12 The matrix B_{μ_0} is positive definite and the matrix H_{μ_0} has rank r .

7.5 The Cramér-Rao Inequality

Let $\frac{\partial \log f(x; \mu)}{\partial \mu} = s(\mu)$; $s(\mu)$ is a $(k \times 1)$ vector of functions of x , and hence is a random variable. Now

$$\int f(x; \mu) dx = 1$$

and hence by assumption A.1

$$\int \frac{\partial f(x; \mu)}{\partial \mu} \frac{1}{f(x; \mu)} f(x; \mu) dx = 0;$$

or

$$\int \frac{\partial \log f(x; \mu)}{\partial \mu_i} f(x; \mu) dx = 0 \quad i = 1; 2; \dots; k; \quad (7.6)$$

Equation (7.6) may be represented as a $(k \times 1)$ vector of integrals. These are written

$$\int \frac{\partial \log f(x; \mu)}{\partial \mu} f(x; \mu) dx = 0;$$

Thus $E_{\mu}[s(\mu)] = 0$, implying that

$$E_{\mu} [s(\mu) s(\mu)^{\prime}] = D_{\mu} [s(\mu)] = A_{\mu} \quad (7.7)$$

in which D_{μ} is the operation denoting dispersion relative to the distribution of x when the parameter is μ . Differentiating (7.6) a second time with respect to μ yields

$$\int \frac{\partial \log f(x; \mu)}{\partial \mu} \frac{\partial f(x; \mu)}{\partial \mu^{\prime}} \frac{1}{f(x; \mu)} f(x; \mu) dx + \int \frac{\partial^2 \log f(x; \mu)}{\partial \mu \partial \mu^{\prime}} f(x; \mu) dx = 0 \quad (7.8)$$

where the integrals now represent $(k \times k)$ matrices of integrals. Since $\frac{\partial f(x; \mu)}{\partial \mu} = \frac{1}{f(x; \mu)} \frac{\partial \log f(x; \mu)}{\partial \mu}$ it follows from (7:7) and (7:8) that the first term on the left-hand side of (7:8) is precisely (7:7) while the second term is $-B_\mu$. Hence

$$A_\mu = B_\mu = E_\mu \left[\frac{\partial^2 \log f(x; \mu)}{\partial \mu \partial \mu} \right] \quad (7.9)$$

Let $t(x)$ be an unbiased estimator of μ . Then

$$\int t(x) f(x; \mu) dx = \mu$$

in which $t(x)$ and μ are each $(k \times 1)$ and hence the integral represents the integral of k separate functions. $E_\mu [t_i(x)] = \mu_i$. It follows, again from assumption A.1, that

$$\int t(x) \frac{\partial \log f(x; \mu)}{\partial \mu} f(x; \mu) dx = \frac{\partial \mu}{\partial \mu} = I_k$$

Thus

$$E_\mu [t(x) s^a(\mu)] = I_k$$

But since $E_\mu [s(\mu)] = 0$, $E_\mu [t(x) | \mu] s^a(\mu) = E_\mu [t(x) s^a(\mu)] + \mu E_\mu [s^a(\mu)] = I_k + 0$. Thus the covariance matrix between $t(x)$ and s is I_k . Now consider

$$\begin{aligned} D_\mu \begin{pmatrix} t(x) \\ s(\mu) \end{pmatrix} &= \begin{pmatrix} D_\mu [t(x)] & E_\mu [t(x) | \mu] s^a(\mu) \\ E_\mu [s(\mu) | t(x) | \mu] & D_\mu [s(\mu)] \end{pmatrix} \\ &= \begin{pmatrix} D[t(x)] & I_k \\ I_k & B_\mu \end{pmatrix} \end{aligned}$$

which, being a dispersion, is required to be at least non-negative definite. But I_k and B_μ are each positive definite whereupon, for every $q \in \mathbb{R}^k$,

$$q^T \left(I_k - \frac{D_\mu[t(x)]^T D_\mu[t(x)]}{\text{tr}(D_\mu[t(x)])} \right) q \geq 0 \quad (7.10)$$

Thus $\frac{D_\mu[t(x)]^T D_\mu[t(x)]}{\text{tr}(D_\mu[t(x)])}$ in (7.10) is non-negative definite, or, for any fixed $c \in \mathbb{R}^k$, the variance of $c^T t(x)$ is

$$V_\mu \{c^T t(x)\} = c^T D_\mu[t(x)] c \geq c^T B_\mu^{-1} c \quad (7.11)$$

It follows that $c^T B_\mu^{-1} c$ represents a lower bound below which the variance of any unbiased estimator of $c^T \mu$ cannot go. Another way of expressing (7.10) or (7.11) is

$$\det D_\mu[t(x)] \geq \det B_\mu^{-1}$$

Moreover, the smaller (greater) is $\det B_\mu^{-1}$, the greater (smaller) is $\det B_\mu$ implying that the greater (smaller) is the information contained in a single observation.

7.6 Maximum-likelihood Estimation in -

Many of the assumptions A.1{A.10 are concerned with establishing the existence of a supremum of $z(\mu)$ at the point $\mu^* \in \Omega$. Ultimately interest will concentrate on the case of estimation in Ω (called restricted maximum-likelihood), but to begin with, unrestricted (or free) estimation in Ω is considered.

It is, of course, presumed on the maintained hypothesis H_m that $\mu_0 \in \Theta$. Assumption A.3 ensures that the supremum of any continuous function on Θ is attained at some interior point of Θ while assumptions A.3 and A.4 together ensure that $\frac{1}{n} \log L_n(x; \mu)$, for almost all x , attains its supremum at a point $\hat{\mu} \in \Theta$. Finally, as has already been noted, assumptions A.3{A.5 ensure that, for almost all x , the sequence $\frac{1}{n} \log L_n(x; \mu)$ converges to $z(\mu)$ uniformly with respect to μ . Now $z(\mu)$ assumes a maximum at the point μ_0 . If the distributions on Θ corresponding to different values of μ are essentially different, then for no other μ is $z(\mu)$ equal to $z(\mu_0)$. The fact that $z(\mu)$ attains a maximum at μ_0 in Θ together with the assumptions that ensure that $\frac{1}{n} \log L_n(x; \mu)$ is uniformly near $z(\mu)$ guarantee that $\frac{1}{n} \log L_n(x; \mu)$ assumes its maximum at $\hat{\mu}$ which is near μ_0 . Thus as $n \rightarrow \infty$ the Strong Law of Large Numbers (SLLN) will ensure that $\hat{\mu} \xrightarrow{p} \mu_0$.

Since $\hat{\mu}$ is interior to Θ , it can be obtained as a solution to

$$\frac{\partial}{\partial \mu} \log L_n(x; \hat{\mu}) = 0 \quad (7.12)$$

Using assumption A.10, a Taylor's expansion around μ_0 yields, from (7.12),

$$0 = \frac{\partial}{\partial \mu} \log L_n(x; \hat{\mu}) = \frac{\partial}{\partial \mu} \log L(x; \mu_0) + \frac{\partial^2}{\partial \mu^2} \log L(x; \mu_0) (\hat{\mu} - \mu_0) + o_p(1) \quad (7.13)$$

By the Weak Law of Large Numbers (WLLN)

$$p \lim \frac{\partial^2}{\partial \mu^2} \log L(x; \mu_0) = -I_{\mu_0} \quad (7.14)$$

which by assumption A.12 is positive definite. Hence

$$\hat{\mu}_n - \mu_0 = B_{\mu_0}^{-1} \frac{1}{n} \frac{\partial \log L_n(x; \mu_0)}{\partial \mu} + o_p(1) \quad (7.15)$$

Now

$$\frac{1}{n} \frac{\partial \log L_n(x; \mu_0)}{\partial \mu} = s_n(\mu) = \frac{1}{n} \sum_{i=1}^n \frac{\partial \log f(x_i; \mu_0)}{\partial \mu}$$

and by (7:6) and (7:9) each $\frac{\partial \log f(x_i; \mu_0)}{\partial \mu}$ is an independent random vector having mean zero and dispersion B_{μ_0} . Hence $s_n(\mu_0)$ will have mean zero and dispersion $n^{-1} B_{\mu_0}$ and, by the Multivariate Central Limit Theorem,

$$P_n \left(\frac{1}{n} \frac{\partial \log L_n(x; \mu_0)}{\partial \mu} \right) \rightarrow P_{ns_n(\mu_0)} \approx N(0; B_{\mu_0})$$

Since B_{μ_0} does not depend on n , $P_n \left(\hat{\mu}_n - \mu_0 \right)$ in (7:15) is $O_p(1)$ and hence

$$P_n \left(\hat{\mu}_n - \mu_0 \right) \rightarrow^d N(0; B_{\mu_0}^{-1}) \quad (7.16)$$

Notice that the procedure that yields (7:16) as its final conclusion involves four distinct steps.

1. The imposition of sufficient regularity.
2. A Taylor's expansion applied to the maximum-likelihood equation (7:12).
3. The WLLN applied to $\frac{1}{n} \frac{\partial^2 \log L_n(x; \mu)}{\partial \mu \partial \mu^T}$.
4. The Multivariate Central Limit Theorem (MCLT) applied to $P_n s_n(\mu_0)$.

This procedure is characteristic of all maximum-likelihood theory. Finally note that $\sqrt{n}(\hat{\mu} - \mu_0) \xrightarrow{d} N(0; B_{\mu_0}^{-1})$ implies that approximately, for large n ,

$$\hat{\mu} \approx N\left(\mu_0; \frac{1}{n} B_{\mu_0}^{-1}\right) \quad (7.17)$$

In practice $B_{\mu_0}^{-1}$ is not known, but it can be estimated consistently since, as $n \rightarrow \infty$, the distribution of $\hat{\mu}$ collapses on μ_0 . Thus $B_{\mu_0}^{-1}$ is consistently estimated as $B_{\hat{\mu}}^{-1}$, as is suggested by (7.14).

7.7 Maximum-likelihood Estimation in \mathcal{R}^k

To find the maximum-likelihood estimator of μ subject to the r restrictions $h(\mu) = 0$, the method of Lagrange multipliers is employed. The Lagrange function is L and

$$L = \frac{1}{n} \log L_n(x; \mu) + h^{\lambda}(\mu),$$

where λ is the $(r \times 1)$ vector of Lagrange multipliers. If $\hat{\lambda}$ and $\hat{\mu}$ are the restricted maximum-likelihood estimators of λ and μ , then these vectors are obtained at a stationary point on L satisfying

$$\begin{aligned} \frac{1}{n} \frac{\partial \log L_n(x; \mu)}{\partial \mu} + H_{\mu}^{\lambda} \hat{\lambda} &= 0 \\ h(\hat{\mu}) &= 0; \end{aligned} \quad (7.18)$$

the first of these implying $s_n(\hat{\mu}) = -H_{\mu}^{\lambda} \hat{\lambda}$.

Since $\hat{\mu} \xrightarrow{p} \mu^0$ for almost all x , applying Taylor's theorem to the first equation in

(7:18) yields

$$\frac{1}{n} \frac{\partial \log L_n(x; \mu)}{\partial \mu} = \frac{1}{n} \frac{\partial \log L_n(x; \mu^a)}{\partial \mu} + \frac{1}{n} \frac{\partial^2 \log L_n(x; \mu^a)}{\partial \mu \partial \mu^a} (\mu - \mu^a) + o_p(1) \quad (7.19)$$

for almost all x with probability near 1. Moreover

$$H_{\mu^a}^{\mu} = H_{\mu^a}^{\mu^a} + o_p(1) : \quad (7.20)$$

Equation (7:20) holds because when μ is near μ^a , it is also near μ_0 (and hence $\hat{\mu}$) so that $\mu - \mu^a$ is relatively small. Thus expanding $H_{\mu^a}^{\mu}$ around μ^a , the first-order terms in the expansion involve $\mu - \mu^a$ and $\mu - \mu^a$ and so are of smaller order than those that have been included in (7:20). Equations (7:19) and (7:20) together enable the re-writing of the first equation of (7:18), to $o_p(1)$, as

$$i \frac{1}{n} \frac{\partial^2 \log L_n(x; \mu)}{\partial \mu \partial \mu^a} \mu - \mu^a = i H_{\mu^a}^{\mu} \mu - \mu^a = P_{ns_n}(\mu^a) : \quad (7.21)$$

The argument regarding the expansion (7:20) may also be applied to the second equation of (7:18):

$$h_{\mu^a}^{\mu} = h(\mu^a) + H_{\mu^a}^{\mu} \mu - \mu^a + o_p(1)$$

in which, since $\mu^a \neq 0$, $h(\mu^a) = 0$. This last expression along with (7:21) leads to the matrix equation, to $o_p(1)$,

$$\begin{bmatrix} 2 \\ 6 \\ 4 \end{bmatrix} i \frac{1}{n} \frac{\partial^2 \log L_n(x; \mu)}{\partial \mu \partial \mu^a} \mu - \mu^a = i H_{\mu^a}^{\mu} \begin{bmatrix} 3 \\ 7 \\ 5 \end{bmatrix} \mu - \mu^a = \begin{bmatrix} 3 \\ 7 \\ 5 \end{bmatrix} P_{ns_n}(\mu^a) \begin{bmatrix} 3 \\ 7 \\ 5 \end{bmatrix} \quad (7.22)$$

Before proceeding to apply the WLLN it is useful to consider the role of the function $z(\mu)$. As with unrestricted estimation, $z(\mu)$ plays an important role in restricted maximum-likelihood estimation.

The regularity conditions ensure that $z(\mu)$ reaches a maximum at the point μ_0 in Ω , and reaches a supremum at $z(\mu^*)$ when μ is constrained to be chosen from Ω_0 . If $\mu_0 \in \Omega_0$, then μ_0 replaces μ^* and $z(\mu_0)$ is the maximum of $z(\mu)$ in the set Ω_0 . In this case $\frac{\partial z(\mu^*)}{\partial \mu} = \frac{\partial z(\mu_0)}{\partial \mu} = 0$ and the Lagrange multiplier that is used to constrain the choice of μ to μ in Ω_0 in any sample, must converge to zero as $n \rightarrow \infty$. This is because in the limit, as $n \rightarrow \infty$, the estimator $\hat{\mu}$ of μ in Ω_0 must converge to its true value μ_0 in Ω_0 and restraining it to lie there eventually becomes unnecessary. In this case, then, the solutions, $\hat{\mu}$ and $\hat{\lambda}$, to the restrained maximum-likelihood equations (7:18) will tend to μ_0 and zero, since the null hypothesis $H_0 : h(\mu) = 0$ is undoubtedly satisfied.

What happens in a practical case when n is large and μ_0 , while not belonging to Ω_0 , is nevertheless near to this set? In this case, $z(\mu_0)$ will be $\sup_{\mu \in \Omega_0} z(\mu)$ and $\mu^* \in \Omega_0$ will be near to μ_0 ; that is, $H_0 : \mu_0 \in \Omega_0$, while not strictly true, is in fact very close to being true. In particular, $\frac{\partial z(\mu^*)}{\partial \mu}$ will be near to the same expression for $\mu = \mu_0$ and $\frac{\partial z(\mu_0)}{\partial \mu}$ is indeed zero. Then $\hat{\lambda}$ will be near to zero and moreover $\frac{\partial^2 z(\mu^*)}{\partial \mu \partial \mu'}$ will have elements near to the corresponding element of $\frac{\partial^2 z(\mu_0)}{\partial \mu \partial \mu'} = -I_{\mu_0}$. In short, convergence of expressions determined at μ^* will be toward the corresponding expressions determined at μ_0 .

These arguments will now be applied to equation (7:22) in which it may be noted

that, by applying the MCLT,

$$p_{n_{S_n}}(\mu_0) = \frac{1}{n} \frac{\partial \log L_n(x; \mu_0)}{\partial \mu} \Big|_{\mu_0} \quad (0; B_{\mu_0}) \quad (7.23)$$

as before. Thus it is in order to write (7:22) as

$$\frac{1}{n} \frac{\partial \log L_n(x; \mu_0)}{\partial \mu} \Big|_{\mu_0} = \frac{1}{n} \frac{\partial \log L_n(x; \mu_0)}{\partial \mu} \Big|_{\mu_0} \quad (7.24)$$

Now

$$\begin{aligned} & \frac{1}{n} \frac{\partial \log L_n(x; \mu_0)}{\partial \mu} \Big|_{\mu_0} \\ &= \frac{1}{n} \frac{\partial \log L_n(x; \mu_0)}{\partial \mu} \Big|_{\mu_0} \end{aligned}$$

whereupon, from (7:22), (7:23) and (7:24)

$$p_{n_{S_n}}(\mu_0) = \frac{1}{n} \frac{\partial \log L_n(x; \mu_0)}{\partial \mu} \Big|_{\mu_0}$$

implying that

$$p_{n_{S_n}}(\mu_0) = \frac{1}{n} \frac{\partial \log L_n(x; \mu_0)}{\partial \mu} \Big|_{\mu_0}$$

or

$$p_{n_{S_n}}(\mu_0) = \frac{1}{n} \frac{\partial \log L_n(x; \mu_0)}{\partial \mu} \Big|_{\mu_0} \quad (7.25)$$

In addition, $\rho_{\beta_i, \mu_0}^{-3} = B_{\mu_0}^{i-1} H_{\mu_0}^{i-1} H_{\mu_0}^{i-1} H_{\mu_0}^{i-1} \zeta_{i-1} H_{\mu_0} B_{\mu_0}^{-1} \rho_{\beta_i, \mu_0}^{-3}$ and hence

$$\rho_{\beta_i, \mu_0}^{-3} \stackrel{d}{\sim} N\left(0; \begin{matrix} B_{\mu_0}^{i-1} H_{\mu_0}^{i-1} H_{\mu_0}^{i-1} H_{\mu_0}^{i-1} \zeta_{i-1} H_{\mu_0} B_{\mu_0}^{-1} & 0 \\ 0 & B_{\mu_0} \end{matrix} \right)$$

or

$$\rho_{\beta_i, \mu_0}^{-3} \stackrel{d}{\sim} N\left(0; \begin{matrix} B_{\mu_0}^{i-1} H_{\mu_0}^{i-1} H_{\mu_0}^{i-1} H_{\mu_0}^{i-1} \zeta_{i-1} H_{\mu_0} B_{\mu_0}^{-1} & 0 \\ 0 & B_{\mu_0} \end{matrix} \right) \quad (7.26)$$

Finally, the covariance term between $\rho_{\beta_i, \mu_0}^{-3}$ and $\rho_{\beta_i, \mu_0}^{-1}$ is

$$\begin{aligned} & \begin{matrix} h \\ B_{\mu_0}^{i-1} H_{\mu_0}^{i-1} H_{\mu_0}^{i-1} H_{\mu_0}^{i-1} \zeta_{i-1} H_{\mu_0} B_{\mu_0}^{-1} \end{matrix} \begin{matrix} h \\ B_{\mu_0}^{i-1} H_{\mu_0}^{i-1} H_{\mu_0}^{i-1} H_{\mu_0}^{i-1} \zeta_{i-1} H_{\mu_0} B_{\mu_0}^{-1} \end{matrix} \\ &= \begin{matrix} B_{\mu_0}^{i-1} H_{\mu_0}^{i-1} H_{\mu_0}^{i-1} H_{\mu_0}^{i-1} \zeta_{i-1} H_{\mu_0} B_{\mu_0}^{-1} \\ B_{\mu_0}^{i-1} H_{\mu_0}^{i-1} H_{\mu_0}^{i-1} H_{\mu_0}^{i-1} \zeta_{i-1} H_{\mu_0} B_{\mu_0}^{-1} \end{matrix} \\ &= 0; \end{aligned} \quad (7.27)$$

Putting (7:25), (7:26) and (7:27) into one equation

$$\begin{pmatrix} \rho_{\beta_i, \mu_0}^{-3} \\ \rho_{\beta_i, \mu_0}^{-1} \end{pmatrix} \stackrel{d}{\sim} N\left(0; \begin{matrix} \begin{matrix} B_{\mu_0}^{i-1} H_{\mu_0}^{i-1} H_{\mu_0}^{i-1} H_{\mu_0}^{i-1} \zeta_{i-1} H_{\mu_0} B_{\mu_0}^{-1} & 0 \\ 0 & B_{\mu_0} \end{matrix} & 0 \\ 0 & \begin{matrix} B_{\mu_0}^{i-1} H_{\mu_0}^{i-1} H_{\mu_0}^{i-1} H_{\mu_0}^{i-1} \zeta_{i-1} H_{\mu_0} B_{\mu_0}^{-1} \\ B_{\mu_0} \end{matrix} \end{matrix} \right)$$

where

$$\begin{matrix} \begin{matrix} B_{\mu_0}^{i-1} H_{\mu_0}^{i-1} H_{\mu_0}^{i-1} H_{\mu_0}^{i-1} \zeta_{i-1} H_{\mu_0} B_{\mu_0}^{-1} & 0 \\ 0 & B_{\mu_0} \end{matrix} & 0 \\ 0 & \begin{matrix} B_{\mu_0}^{i-1} H_{\mu_0}^{i-1} H_{\mu_0}^{i-1} H_{\mu_0}^{i-1} \zeta_{i-1} H_{\mu_0} B_{\mu_0}^{-1} \\ B_{\mu_0} \end{matrix} \end{matrix}$$

By assumption A.12 B_{μ_0} is positive definite. Hence there will exist a non-singular

matrix Q_{μ_0} such that $Q_{\mu_0}^> Q_{\mu_0} = B_{\mu_0}^{i-1}$. Writing $G_{\mu_0} = Q_{\mu_0}^> H_{\mu_0}^{i-1}$, then

$$\begin{matrix} \begin{matrix} Q_{\mu_0}^> I_k & G_{\mu_0} \\ G_{\mu_0}^> G_{\mu_0} & \zeta_{i-1} G_{\mu_0}^> Q_{\mu_0} \end{matrix} & 0 \\ 0 & \begin{matrix} G_{\mu_0}^> G_{\mu_0} \\ G_{\mu_0}^> G_{\mu_0} \end{matrix} \end{matrix} \zeta_{i-1}$$

and, if $P_{\mu_0} = G_{\mu_0}^{-1} G_{\mu_0}^{\prime} G_{\mu_0}^{-1}$ and $M_{\mu_0} = (I_{k-r} \quad P_{\mu_0})$, then ξ reduces to

$$\xi = \begin{pmatrix} Q_{\mu_0}^{\prime} M_{\mu_0} Q_{\mu_0} & 0 \\ 0 & G_{\mu_0}^{\prime} G_{\mu_0} \end{pmatrix}$$

Since $Q_{\mu_0}^{\prime} M_{\mu_0} Q_{\mu_0}$ must be non-negative definite (since Q_{μ_0} is non-singular and M_{μ_0} is non-negative definite), ξ must be non-negative definite. This is reasonable since μ is a k -dimensional subset of \mathbb{R}^k which is restricted by r restraint equations; thus μ has $k-r$ elements, $(k-r)$ of which are independently estimated. Thus the dispersion of ρ_{μ}^3 would be expected to have rank $(k-r)$. Since Q_{μ_0} is non-singular, and by assumption A.12 H_{μ_0} has rank r , $Q_{\mu_0}^{\prime} M_{\mu_0} Q_{\mu_0}$ has rank $(k-r)$.

7.8 Associated Tests of Significance

Notice from (7:18) that

$$H_{\mu}^3 = \frac{-2 \log L_n(x; \mu)}{n} = \rho_{\mu}^3 \quad (7.28)$$

Since for large n the approximate distribution of ρ_{μ}^3 is $N(0; H_{\mu_0} B_{\mu_0}^{-1} H_{\mu_0}^{\prime})$, then

$$\rho_{\mu}^3 \approx D \left(\frac{-2 \log L_n(x; \mu)}{n} \right) = n^{-1/2} H_{\mu} B_{\mu}^{-1} H_{\mu}^{\prime} \quad (7.29)$$

is approximately $\hat{A}^2(r; \pm)$ with $\pm = 0$ on H_0 . This is called the Lagrange Multiplier (LM) test. Equation 7:29 is also equivalent, by (7:28), to

$$n s_n^2 \mu B_{\mu}^{-1} s_n^2 \mu \stackrel{d}{=} \hat{A}^2(r; 0) \quad (7.30)$$

where $\rho_{ns_n}^{-3} \hat{\mu} = \frac{1}{n} \frac{\partial \log L_n(x; \mu)}{\partial \mu}$ is the score statistic. This implies that the LM and the score methods lead to the same test.

Also, since $h(\hat{\mu}) = h(\mu_0) + H_{\mu_0}(\hat{\mu} - \mu_0)$,

$$\rho_{ns_n}^{-3} h(\hat{\mu}) = \frac{1}{n} \frac{\partial \log L_n(x; \hat{\mu})}{\partial \mu} = H_{\mu_0} \rho_{ns_n}^{-3}(\hat{\mu} - \mu_0) \stackrel{d}{\rightarrow} N(0; H_{\mu_0} B_{\mu_0}^{-1} H_{\mu_0}^{-1}) \quad (7.31)$$

and hence

$$W = nh(\hat{\mu})' H_{\mu_0}^{-1} B_{\mu_0}^{-1} H_{\mu_0}^{-1} h(\hat{\mu}) \stackrel{d}{\rightarrow} \chi^2(r; \pm) \quad (7.32)$$

and $\pm = 0$ on H_0 . W is known as the Wald statistic.

The Wald test is associated with estimation of μ in $\hat{\mu}$ and the LM test with estimation of μ in μ_0 . The likelihood ratio (LR) test, in contrast, is essentially a test based upon a ratio of the value of the likelihood for estimation in $\hat{\mu}$ relative to the value of the likelihood for estimation in μ_0 . If μ_0 is the LR, then

$$LR = \frac{L_n(x; \mu_0)}{L_n(x; \hat{\mu})} \quad (7.33)$$

Since $L_n(x; \mu_0) \leq L_n(x; \hat{\mu})$ it is clear that $0 < LR \leq 1$. The test is based on the statistic $-2 \log LR$, thus ensuring that the statistic is never negative. The reason for the multiple 2 will become clear as the test statistic is developed below.

The first step is to expand in a Taylor's series $\log L_n(x; \mu)$ around the point $\hat{\mu}$:

$$\log L_n(x; \mu) = \log L_n(x; \hat{\mu}) + \frac{\partial \log L_n(x; \hat{\mu})}{\partial \mu} (\mu - \hat{\mu}) + \frac{1}{2} (\mu - \hat{\mu})' \frac{\partial^2 \log L_n(x; \hat{\mu})}{\partial \mu \partial \mu'} (\mu - \hat{\mu}) + \dots$$

Noting that the second term on the right-hand side is zero,

$$\log L_n(x; \hat{\mu}) - \log L_n(x; \mu_0) = \frac{1}{2} \sum_{i=1}^n \frac{\partial^2 \log L_n(x; \hat{\mu})}{\partial \mu_i^2} \mu_0^i$$

and this may be expressed as

$$\sum_{i=1}^n \log \pi_i(\mu_0) = \sum_{i=1}^n \log \pi_i(\hat{\mu}) - B_{\hat{\mu}} \sum_{i=1}^n \mu_0^i \quad (7.34)$$

Several previously developed results will now become important: from (7:15) and

$$Q_{\mu_0}^> Q_{\mu_0} = B_{\mu_0}^{-1}$$

$$\sum_{i=1}^n \log \pi_i(\mu_0) = B_{\mu_0}^{-1} \sum_{i=1}^n \mu_0^i = Q_{\mu_0}^> Q_{\mu_0} \sum_{i=1}^n \mu_0^i; \quad (7.35)$$

from the solution to (7:24)

$$\sum_{i=1}^n \log \pi_i(\mu_0) = \sum_{i=1}^n \log \pi_i(\mu_0) + \sum_{i=1}^n \log \frac{\pi_i(\mu_0)}{\pi_i(\mu_0)}$$

which reduces to

$$\sum_{i=1}^n \log \pi_i(\mu_0) = Q_{\mu_0}^> M_{\mu_0} Q_{\mu_0} \sum_{i=1}^n \mu_0^i; \quad (7.36)$$

in which $M_{\mu_0} = \sum_{i=1}^n \frac{\partial^2 \log \pi_i(\mu_0)}{\partial \mu_i^2} = \sum_{i=1}^n \frac{\partial^2 \log \pi_i(\mu_0)}{\partial \mu_i^2} = \sum_{i=1}^n \frac{\partial^2 \log \pi_i(\mu_0)}{\partial \mu_i^2} = \sum_{i=1}^n \frac{\partial^2 \log \pi_i(\mu_0)}{\partial \mu_i^2} = \sum_{i=1}^n \frac{\partial^2 \log \pi_i(\mu_0)}{\partial \mu_i^2}$. Since $\sum_{i=1}^n \log \pi_i(\mu_0) \sim N(0; B_{\mu_0})$, $Q_{\mu_0} \sum_{i=1}^n \mu_0^i \sim N(0; I_k)$ and hence taking (7:35) from (7:36)

$$\sum_{i=1}^n \log \pi_i(\mu_0) = Q_{\mu_0}^> (M_{\mu_0} + I_k) Q_{\mu_0} \sum_{i=1}^n \mu_0^i$$

and the last expression reduces to

$$Q_{\mu_0}^> (I_k + M_{\mu_0} + I_k) Q_{\mu_0} \sum_{i=1}^n \mu_0^i = \sum_{i=1}^n \log \pi_i(\mu_0) = \sum_{i=1}^n \log \pi_i(\mu_0)$$

Hence for (7:34)

$$i \cdot 2 \log \alpha \cdot P_{\bar{n}S_n(\mu_0)} \mathbf{f} Q_{\mu_0}^> P_{\mu_0} Q_{\mu_0} B_{\hat{\mu}} Q_{\mu_0}^> P_{\mu_0} Q_{\mu_0} \alpha S_n(\mu_0) P_{\bar{n}}$$

which reduces to

$$i \cdot 2 \log \alpha \cdot P_{\bar{n}S_n^>(\mu_0)} Q_{\mu_0}^> \mathbf{h} Q_{\mu_0} H_{\mu_0}^> \mathbf{i} H_{\mu_0} B_{\mu_0}^i H_{\mu_0}^> \mathbf{c}_i^{-1} H_{\mu_0} Q_{\mu_0}^> \mathbf{i} Q_{\mu_0} P_{\bar{n}S_n(\mu_0)}$$

which is clearly approximately $\hat{A}^2(r; \pm)$ with $\pm = 0$ on $H_0 : h(\mu) = 0$, since the expression in square brackets is a projection matrix of rank r and $Q_{\mu_0} P_{\bar{n}S_n(\mu_0)} \gg N(0; I_k)$.