

# Appendix A

## Vector Spaces

**Abstract** This appendix reviews some of the basic definitions and properties of vector spaces. It presumes some basic knowledge of matrices.

**Definition A.1** A set  $\mathcal{M}$  is a *vector space* if, for any  $x, y, z \in \mathcal{M}$  and scalars  $\alpha, \beta$ , operations of vector addition and scalar multiplication are defined such that:

- (1)  $(x + y) + z = x + (y + z)$ .
- (2)  $x + y = y + x$ .
- (3) There exists a vector  $0 \in \mathcal{M}$  such that  $x + 0 = x = 0 + x$  for any  $x \in \mathcal{M}$ .
- (4) For any  $x \in \mathcal{M}$ , there exists  $y \equiv -x$  such that  $x + y = 0 = y + x$ .
- (5)  $\alpha(x + y) = \alpha x + \alpha y$ .
- (6)  $(\alpha + \beta)x = \alpha x + \beta x$ .
- (7)  $(\alpha\beta)x = \alpha(\beta x)$ .
- (8) There exists a scalar  $\xi$  such that  $\xi x = x$ . (Typically,  $\xi = 1$ .)

We rely on context to distinguish between the vector  $0 \in \mathcal{M}$  and the scalar  $0 \in \mathbf{R}$ . In most of our applications, we assume  $\mathcal{M} \subset \mathbf{R}^n$ . Vectors in  $\mathbf{R}^n$  will be considered as  $n \times 1$  matrices. The  $0$  vector referred to in Definition A.1 is just an  $n \times 1$  matrix of zeros.

**Definition A.2** Let  $\mathcal{M}$  be a vector space, and let  $\mathcal{N}$  be a set with  $\mathcal{N} \subset \mathcal{M}$ .  $\mathcal{N}$  is a *subspace* of  $\mathcal{M}$  if  $\mathcal{N}$  is a vector space using the same definitions of vector addition and scalar multiplication as for  $\mathcal{M}$ .

Thinking of vectors in three dimensions as  $(x, y, z)'$ , where  $w'$  denotes the *transpose* of a matrix  $w$ . The subspace consisting of the  $z$  axis is

$$\left\{ \begin{pmatrix} 0 \\ 0 \\ z \end{pmatrix} \mid z \in \mathbf{R} \right\}.$$

The subspace consisting of the  $x, y$  plane is

$$\left\{ \begin{pmatrix} x \\ y \\ 0 \end{pmatrix} \mid x, y \in \mathbf{R} \right\}.$$

The subspace consisting of the plane that is perpendicular to the line  $x = y$  in the  $x, y$  plane is

$$\left\{ \begin{pmatrix} x \\ -x \\ z \end{pmatrix} \mid x, z \in \mathbf{R} \right\}.$$

**Theorem A.3** Let  $\mathcal{M}$  be a vector space, and let  $\mathcal{N}$  be a nonempty subset of  $\mathcal{M}$ . If  $\mathcal{N}$  is closed under vector addition and scalar multiplication, then  $\mathcal{N}$  is a subspace of  $\mathcal{M}$ .

**Theorem A.4** Let  $\mathcal{M}$  be a vector space, and let  $x_1, \dots, x_r$  be in  $\mathcal{M}$ . The set of all linear combinations of  $x_1, \dots, x_r$ , i.e.,  $\{v \mid v = \alpha_1 x_1 + \dots + \alpha_r x_r, \alpha_i \in \mathbf{R}\}$ , is a subspace of  $\mathcal{M}$ .

**Definition A.5** The set of all linear combinations of  $x_1, \dots, x_r$  is called the *space spanned by*  $x_1, \dots, x_r$ . If  $\mathcal{N}$  is a subspace of  $\mathcal{M}$ , and  $\mathcal{N}$  equals the space spanned by  $x_1, \dots, x_r$ , then  $\{x_1, \dots, x_r\}$  is called a *spanning set* for  $\mathcal{N}$ .

The space spanned by the vectors

$$x_1 = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, \quad x_2 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$$

consists of all vectors of the form  $(a, b, b)'$ , where  $a$  and  $b$  are any real numbers.

Let  $A$  be an  $n \times p$  matrix. Each column of  $A$  is a vector in  $\mathbf{R}^n$ . The space spanned by the columns of  $A$  is called the *column space* of  $A$  and written  $C(A)$ . (Some people refer to  $C(A)$  as the *range space* of  $A$  and write it  $R(A)$ .) If  $B$  is an  $n \times r$  matrix, then  $C(A, B)$  is the space spanned by the  $p + r$  columns of  $A$  and  $B$ .

**Definition A.6** Let  $x_1, \dots, x_r$  be vectors in  $\mathcal{M}$ . If there exist scalars  $\alpha_1, \dots, \alpha_r$  not all zero so that  $\sum \alpha_i x_i = 0$ , then  $x_1, \dots, x_r$  are *linearly dependent*. If such  $\alpha_i$ s do not exist, i.e., if  $\sum \alpha_i x_i = 0$  implies that  $\alpha_i = 0$  for all  $i$ , then  $x_1, \dots, x_r$  are *linearly independent*.

**Definition A.7** If  $\mathcal{N}$  is a subspace of  $\mathcal{M}$  and if  $\{x_1, \dots, x_r\}$  is a linearly independent spanning set for  $\mathcal{N}$ , then  $\{x_1, \dots, x_r\}$  is called a *basis* for  $\mathcal{N}$ .

**Theorem A.8** If  $\mathcal{N}$  is a subspace of  $\mathcal{M}$ , all bases for  $\mathcal{N}$  have the same number of vectors.

**Definition A.9** The *rank* of a subspace  $\mathcal{N}$  is the number of elements in a basis for  $\mathcal{N}$ . The rank is written  $r(\mathcal{N})$ . If  $A$  is a matrix, the rank of  $C(A)$  is called the rank of  $A$  and is written  $r(A)$ , i.e.,  $r(A) \equiv r[C(A)]$ .

Exercise B.24 establishes that  $r(A) = r(A')$ .

**Theorem A.10** If  $v_1, \dots, v_r$  is a basis for  $\mathcal{N}$ , and  $x \in \mathcal{N}$ , then the characterization  $x = \sum_{i=1}^r \alpha_i v_i$  is unique.

*Proof* Suppose  $x = \sum_{i=1}^r \alpha_i v_i$  and  $x = \sum_{i=1}^r \beta_i v_i$ . Then  $0 = \sum_{i=1}^r (\alpha_i - \beta_i) v_i$ . Since the vectors  $v_i$  are linearly independent,  $\alpha_i - \beta_i = 0$  for all  $i$ .  $\square$

The vectors

$$x_1 = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, \quad x_2 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \quad x_3 = \begin{pmatrix} 2 \\ 3 \\ 3 \end{pmatrix}$$

are linearly dependent because  $0 = 3x_1 - x_2 - x_3$ . Any two of  $x_1, x_2, x_3$  form a basis for the space of vectors with the form  $(a, b, b)'$ . This space has rank 2.

**Definition A.11** Let  $\mathcal{N}_1$  and  $\mathcal{N}_2$  be vector subspaces. The sum of  $\mathcal{N}_1$  and  $\mathcal{N}_2$  is  $\mathcal{N}_1 + \mathcal{N}_2 \equiv \{x | x = x_1 + x_2 \text{ for some } x_1 \in \mathcal{N}_1, x_2 \in \mathcal{N}_2\}$ . This is sometimes called the *direct sum* and written  $\mathcal{N}_1 \oplus \mathcal{N}_2$ .

**Theorem A.12**  $\mathcal{N}_1 + \mathcal{N}_2$  is a vector space and  $C(A, B) = C(A) + C(B)$ .

**Theorem A.13** Let  $\mathcal{N}_1$  and  $\mathcal{N}_2$  be vector subspaces. If  $\mathcal{N}_1 \cap \mathcal{N}_2 = \{0\}$ , then any vector  $x \in \mathcal{N}_1 + \mathcal{N}_2$  has a unique decomposition  $x = x_1 + x_2$  with  $x_1 \in \mathcal{N}_1$  and  $x_2 \in \mathcal{N}_2$ .

*Proof* Let  $x = x_1 + x_2$ ,  $x_1 \in \mathcal{N}_1$ ,  $x_2 \in \mathcal{N}_2$  and  $x = y_1 + y_2$ ,  $y_1 \in \mathcal{N}_1$ ,  $y_2 \in \mathcal{N}_2$ . Then  $x_1 + x_2 = x = y_1 + y_2$ , so  $x_1 - y_1 = y_2 - x_2 \equiv v$ . But  $v \equiv x_1 - y_1 \in \mathcal{N}_1$  and  $v = y_2 - x_2 \in \mathcal{N}_2$ , so  $v \in \mathcal{N}_1 \cap \mathcal{N}_2 = \{0\}$  and  $x_1 = y_1$  and  $y_2 = x_2$ .  $\square$

**Theorem A.14** Let  $\mathcal{N}_1$  and  $\mathcal{N}_2$  be vector subspaces. If  $\mathcal{N}_1 \cap \mathcal{N}_2 = \{0\}$ , then  $r(\mathcal{N}_1 + \mathcal{N}_2) = r(\mathcal{N}_1) + r(\mathcal{N}_2)$ .

*Proof* Let  $v_1, \dots, v_r$  be a basis for  $\mathcal{N}_1$  and  $w_1, \dots, w_s$  be a basis for  $\mathcal{N}_2$ . It suffices to show that  $v_1, \dots, v_r, w_1, \dots, w_s$  is a basis for  $\mathcal{N}_1 + \mathcal{N}_2$ . Clearly,  $v_1, \dots, v_r, w_1, \dots, w_s$  is a spanning set for  $\mathcal{N}_1 + \mathcal{N}_2$ . If the vectors are linearly independent, the result is proven.

Suppose  $0 = \sum_{i=1}^r \alpha_i v_i + \sum_{j=1}^s \beta_j w_j$  which implies that

$$\sum_{i=1}^r \alpha_i v_i = \sum_{j=1}^s -\beta_j w_j.$$

The left-hand side of the equation is a vector in  $\mathcal{N}_1$  and the right-hand side is a vector in  $\mathcal{N}_2$ , so it is a vector in both  $\mathcal{N}_1$  and  $\mathcal{N}_2$ , hence must be in the intersection which, by assumption, is the 0 vector. Since both  $0 = \sum_{i=1}^r \alpha_i v_i$  and  $0 = \sum_{j=1}^s \beta_j w_j$  and the  $v_i$ s and  $w_j$ s are each a basis, we must have  $0 = \alpha_1 = \dots = \alpha_r = \beta_1 = \dots = \beta_s$ , hence  $v_1, \dots, v_r, w_1, \dots, w_s$  are linearly independent and a basis. It follows immediately that  $r(\mathcal{N}_1 + \mathcal{N}_2) = r + s = r(\mathcal{N}_1) + r(\mathcal{N}_2)$ .  $\square$

**Definition A.15** The (Euclidean) *inner product* between two vectors  $x$  and  $y$  in  $\mathbf{R}^n$  is  $x'y$ . Two vectors  $x$  and  $y$  are *orthogonal* (written  $x \perp y$ ) if  $x'y = 0$ . Two subspaces  $\mathcal{N}_1$  and  $\mathcal{N}_2$  are orthogonal if  $x \in \mathcal{N}_1$  and  $y \in \mathcal{N}_2$  implies that  $x'y = 0$ .  $\{x_1, \dots, x_r\}$  is an *orthogonal basis* for a space  $\mathcal{N}$  if  $\{x_1, \dots, x_r\}$  is a basis for  $\mathcal{N}$  and for  $i \neq j$ ,  $x_i'x_j = 0$ .  $\{x_1, \dots, x_r\}$  is an *orthonormal basis* for  $\mathcal{N}$  if  $\{x_1, \dots, x_r\}$  is an orthogonal basis and  $x_i'x_i = 1$  for  $i = 1, \dots, r$ . The terms *orthogonal* and *perpendicular* are used interchangeably. The *length* of a vector  $x$  is  $\|x\| \equiv \sqrt{x'x}$ . The *distance* between two vectors  $x$  and  $y$  is the length of their difference, i.e.,  $\|x - y\|$ .

The lengths of the vectors given earlier are

$$\|x_1\| = \sqrt{1^2 + 1^2 + 1^2} = \sqrt{3}, \quad \|x_2\| = 1, \quad \|x_3\| = \sqrt{22} \doteq 4.7.$$

If  $x = (2, 1)'$ , its length is  $\|x\| = \sqrt{2^2 + 1^2} = \sqrt{5}$ . If  $y = (3, 2)'$ , the distance between  $x$  and  $y$  is the length of  $x - y = (2, 1)' - (3, 2)' = (-1, -1)'$ , which is  $\|x - y\| = \sqrt{(-1)^2 + (-1)^2} = \sqrt{2}$ .

Just prior to Section B.4 and in Sections 2.7 and 6.3 we discuss more general versions of the concepts of inner product and length. In particular, a more general version of Definition A.15 is given in Subsection 6.3.5. The remaining results and definitions in this appendix are easily extended to general inner products.

Our emphasis on orthogonality and our need to find orthogonal projection matrices make both the following theorem and its proof fundamental tools in linear model theory:

**Theorem A.16 The Gram–Schmidt Theorem.**

Let  $\mathcal{N}$  be a space with basis  $\{x_1, \dots, x_r\}$ . There exists an orthonormal basis for  $\mathcal{N}$ , say  $\{y_1, \dots, y_r\}$ , with  $y_s$  in the space spanned by  $x_1, \dots, x_s$ ,  $s = 1, \dots, r$ .

*Proof* The Gram–Schmidt algorithm defines the  $y_i$ s inductively:

$$\begin{aligned} y_1 &= x_1 / \sqrt{x_1' x_1}, \\ w_s &= x_s - \sum_{i=1}^{s-1} (x_s' y_i) y_i, \\ y_s &= w_s / \sqrt{w_s' w_s}. \end{aligned}$$

See Exercise A.1. □

This result is written using the Euclidean inner product, but terms like  $x_s' y_i$  can be replaced with any general inner product, say,  $\langle x_s, y_i \rangle$ .

You can also apply the Gram–Schmidt algorithm to an arbitrary spanning set, rather than to a basis. It still gives an orthonormal basis for the space originally spanned but it involves bookkeeping issues. If the spanning set has linear dependencies, some of the vectors  $w_j$  in the algorithm will be 0 vectors, so they need to be dropped from the orthonormal basis.

The vectors

$$x_1 = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, \quad x_2 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$$

are a basis for the space of vectors with the form  $(a, b, b)'$ . To orthonormalize this basis, take  $y_1 = x_1 / \sqrt{3}$ . Then take

$$w_2 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} - \frac{1}{\sqrt{3}} \begin{pmatrix} 1/\sqrt{3} \\ 1/\sqrt{3} \\ 1/\sqrt{3} \end{pmatrix} = \begin{pmatrix} 2/3 \\ -1/3 \\ -1/3 \end{pmatrix}.$$

Finally, normalize  $w_2$  to give

$$y_2 = w_2 / \sqrt{6/9} = (2/\sqrt{6}, -1/\sqrt{6}, -1/\sqrt{6})'.$$

Note that another orthonormal basis for this space consists of the vectors

$$z_1 = \begin{pmatrix} 0 \\ 1/\sqrt{2} \\ 1/\sqrt{2} \end{pmatrix}, \quad z_2 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}.$$

The result of Gram–Schmidt depends on the order in which you list the vectors in the basis. If you change the order, typically you get a different orthonormal basis. In fact, the  $z_i$ s are the vectors you get if you change the order of the two  $x_i$ s.

**Definition A.17** For  $\mathcal{N}$  a subspace of  $\mathcal{M}$ , let  $\mathcal{N}_{\mathcal{M}}^{\perp} \equiv \{y \in \mathcal{M} \mid y \perp \mathcal{N}\}$ .  $\mathcal{N}_{\mathcal{M}}^{\perp}$  is called the *orthogonal complement* of  $\mathcal{N}$  with respect to  $\mathcal{M}$ . If  $\mathcal{M}$  is taken as  $\mathbf{R}^n$ , then  $\mathcal{N}^{\perp} \equiv \mathcal{N}_{\mathbf{R}^n}^{\perp}$  is simply referred to as the orthogonal complement of  $\mathcal{N}$ .

**Corollary A.18** For any subspace  $\mathcal{N}$ ,  $\mathcal{N} \cap \mathcal{N}^{\perp} = \{0\}$ .

*Proof* For any  $x \in \mathcal{N} \cap \mathcal{N}^{\perp}$ , the vector must be orthogonal to itself, so  $x'x = 0$  and  $x = 0$ . □

**Theorem A.19** Let  $\mathcal{M}$  be a vector space, and let  $\mathcal{N}$  be a subspace of  $\mathcal{M}$ .  $\mathcal{N}_{\mathcal{M}}^{\perp}$  is a subspace of  $\mathcal{M}$  and  $\mathcal{M} = \mathcal{N} + \mathcal{N}_{\mathcal{M}}^{\perp}$ .

Before proving Theorem A.19 we state and prove the following corollary and give examples.

**Corollary A.20** Any vector  $x \in \mathcal{M}$  can be written uniquely as  $x = x_1 + x_2$  with  $x_1 \in \mathcal{N}$  and  $x_2 \in \mathcal{N}_{\mathcal{M}}^{\perp}$ . Moreover,  $r(\mathcal{M}) = r(\mathcal{N}) + r(\mathcal{N}_{\mathcal{M}}^{\perp})$ .

*Proof* If  $\mathcal{M} = \mathcal{N} + \mathcal{N}_{\mathcal{M}}^{\perp}$ , the results are immediate from applying Corollary A.18 and Theorems A.13 and A.14. □

Let  $\mathcal{M} = \mathbf{R}^3$  and let  $\mathcal{N}$  be the space of vectors with the form  $(a, b, b)'$ . It is not difficult to see that the orthogonal complement of  $\mathcal{N}$  consists of vectors of the form  $(0, c, -c)'$ . Any vector  $(x, y, z)'$  can be written uniquely as

$$\begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} x \\ (y+z)/2 \\ (y+z)/2 \end{pmatrix} + \begin{pmatrix} 0 \\ (y-z)/2 \\ -(y-z)/2 \end{pmatrix}.$$

The space of vectors with form  $(a, b, b)'$  has rank 2, and the space  $(0, c, -c)'$  has rank 1.

For additional examples, let

$$X_0 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \quad \text{and} \quad X = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \end{bmatrix}.$$

In this case,

$$C(X_0)^\perp = C\left(\begin{bmatrix} -1 & 1 \\ 0 & -2 \\ 1 & 1 \end{bmatrix}\right), \quad C(X_0)_{C(X)}^\perp = C\left(\begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix}\right),$$

and

$$C(X)^\perp = C\left(\begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix}\right).$$

*Proof of Theorem A.19* It is easily seen that  $\mathcal{N}_{\mathcal{M}}^\perp$  is a subspace by checking Theorem A.3. Since  $\mathcal{N}$  and  $\mathcal{N}_{\mathcal{M}}^\perp$  are subspaces of  $\mathcal{M}$ , we have  $\mathcal{N} + \mathcal{N}_{\mathcal{M}}^\perp \subset \mathcal{M}$ . To show equality it remains to show that  $\mathcal{M} \subset \mathcal{N} + \mathcal{N}_{\mathcal{M}}^\perp$ , i.e., if  $x \in \mathcal{M}$ , then  $x \in \mathcal{N} + \mathcal{N}_{\mathcal{M}}^\perp$ .

Let  $r(\mathcal{M}) = n$  and  $r(\mathcal{N}) = r$ . Let  $v_1, \dots, v_r$  be a basis for  $\mathcal{N}$  and extend this with  $w_1, \dots, w_{n-r}$  to a basis for  $\mathcal{M}$ . (Alternatively, take  $w_1, \dots, w_n$  a basis for  $\mathcal{M}$  but define the spanning set  $\{v_1, \dots, v_r, w_1, \dots, w_n\}$  for  $\mathcal{M}$ .) Apply Gram–Schmidt to get  $v_1^*, \dots, v_r^*, w_1^*, \dots, w_{n-r}^*$  an orthonormal basis for  $\mathcal{M}$  with  $v_1^*, \dots, v_r^*$  an orthonormal basis for  $\mathcal{N}$ . By construction  $\{w_1^*, \dots, w_{n-r}^*\} \subset \mathcal{N}_{\mathcal{M}}^\perp$ .

If  $x \in \mathcal{M}$ , then

$$x = \sum_{i=1}^r \alpha_i v_i^* + \sum_{j=1}^{n-r} \beta_j w_j^*.$$

Let  $x_0 \equiv \sum_{i=1}^r \alpha_i v_i^*$  and  $x_1 \equiv \sum_{j=1}^{n-r} \beta_j w_j^*$ . Then  $x_0 \in \mathcal{N}$ ,  $x_1 \in \mathcal{N}_{\mathcal{M}}^\perp$ , and  $x = x_0 + x_1$ , so  $x \in \mathcal{N} + \mathcal{N}_{\mathcal{M}}^\perp$ .  $\square$

From Corollary A.20 we have  $r(\mathcal{N}_{\mathcal{M}}^\perp) = n - r$ , so  $w_1^*, \dots, w_{n-r}^*$  must be a basis for  $\mathcal{N}_{\mathcal{M}}^\perp$ . This relies on the fact that if the  $w_j^*$ s are all in  $\mathcal{N}_{\mathcal{M}}^\perp$  and if they are linearly independent with the same number of vectors as in a basis for the  $\mathcal{N}_{\mathcal{M}}^\perp$ , then they must be a spanning set for  $\mathcal{N}_{\mathcal{M}}^\perp$ . But we have not shown directly that  $w_1^*, \dots, w_{n-r}^*$  is a spanning set for  $\mathcal{N}_{\mathcal{M}}^\perp$ . We do that now.

If  $x \in \mathcal{N}_{\mathcal{M}}^\perp$ , because  $x \in \mathcal{M}$  write

$$x = \sum_{i=1}^r \alpha_i v_i^* + \sum_{j=1}^{n-r} \beta_j w_j^*.$$

Since  $x \in \mathcal{N}_{\mathcal{M}}^\perp$  and  $v_k^* \in \mathcal{N}$  for  $k = 1, \dots, r$ ,

$$\begin{aligned}
 0 = x'v_k^* &= \left( \sum_{i=1}^r \alpha_i v_i^* + \sum_{j=1}^{n-r} \beta_j w_j^* \right)' v_k^* \\
 &= \sum_{i=1}^r \alpha_i v_i^{*'} v_k^* + \sum_{j=1}^{n-r} \beta_j w_j^{*'} v_k^* \\
 &= \alpha_k v_k^{*'} v_k^* = \alpha_k.
 \end{aligned}$$

Thus  $x = \sum_{j=1}^{n-r} \beta_j w_j^*$ , implying that  $\{w_1^*, \dots, w_{n-r}^*\}$  is a spanning set and a basis for  $\mathcal{N}_{\mathcal{A}}^\perp$ .

**Theorem A.21** Let  $\mathcal{N}_1$  and  $\mathcal{N}_2$  be subspaces, then  $(\mathcal{N}_1 \cap \mathcal{N}_2)^\perp = \mathcal{N}_1^\perp + \mathcal{N}_2^\perp$ .

*Proof* This is a restatement of Proposition 10.4.6. The proof is given there.  $\square$

## Exercises

**Exercise A.1** Give a detailed proof of the Gram–Schmidt theorem.

Questions A.2 through A.13 involve the following matrices:

$$A = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad D = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 2 & 5 \\ 0 & 0 \end{bmatrix}, \quad E = \begin{bmatrix} 1 & 2 \\ 1 & 2 \\ 2 & 7 \\ 0 & 0 \end{bmatrix},$$

$$F = \begin{bmatrix} 1 & 5 & 6 \\ 1 & 5 & 6 \\ 0 & 7 & 2 \\ 0 & 0 & 9 \end{bmatrix}, \quad G = \begin{bmatrix} 1 & 0 & 5 & 2 \\ 1 & 0 & 5 & 2 \\ 2 & 5 & 7 & 9 \\ 0 & 0 & 0 & 3 \end{bmatrix}, \quad H = \begin{bmatrix} 1 & 0 & 2 & 2 & 6 \\ 1 & 0 & 2 & 2 & 6 \\ 7 & 9 & 3 & 9 & -1 \\ 0 & 0 & 0 & 3 & -7 \end{bmatrix},$$

$$K = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}, \quad L = \begin{bmatrix} 2 & 0 & 0 \\ 2 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}, \quad N = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}.$$

**Exercise A.2** Is the space spanned by the columns of  $A$  the same as the space spanned by the columns of  $B$ ? How about the spaces spanned by the columns of  $K$ ,  $L$ ,  $F$ ,  $D$ , and  $G$ ?

**Exercise A.3** Give a matrix whose column space contains  $C(A)$ .



**Exercise A.4** Give two matrices whose column spaces contain  $C(B)$ .

**Exercise A.5** Which of the following equalities are valid:  $C(A) = C(A, D)$ ,  $C(D) = C(A, B)$ ,  $C(A, N) = C(A)$ ,  $C(N) = C(A)$ ,  $C(A) = C(F)$ ,  $C(A) = C(G)$ ,  $C(A) = C(H)$ ,  $C(A) = C(D)$ ?

**Exercise A.6** Which of the following matrices have linearly independent columns:  $A, B, D, N, F, H, G$ ?

**Exercise A.7** Give a basis for the space spanned by the columns of each of the following matrices:  $A, B, D, N, F, H, G$ .

**Exercise A.8** Give the ranks of  $A, B, D, E, F, G, H, K, L, N$ .

**Exercise A.9** Which of the following matrices have columns that are mutually orthogonal:  $B, A, D$ ?

**Exercise A.10** Give an orthogonal basis for the space spanned by the columns of each of the following matrices:  $A, D, N, K, H, G$ .

**Exercise A.11** Find  $C(A)^\perp$  and  $C(B)^\perp$  (with respect to  $\mathbf{R}^4$ ).

**Exercise A.12** Find two linearly independent vectors in the orthogonal complement of  $C(D)$  (with respect to  $\mathbf{R}^4$ ).

**Exercise A.13** Find a vector in the orthogonal complement of  $C(D)$  with respect to  $C(A)$ .

**Exercise A.14** Find an orthogonal basis for the space spanned by the columns of

$$X = \begin{bmatrix} 1 & 1 & 4 \\ 1 & 2 & 1 \\ 1 & 3 & 0 \\ 1 & 4 & 0 \\ 1 & 5 & 1 \\ 1 & 6 & 4 \end{bmatrix}.$$

**Exercise A.15** For  $X$  as above, find two linearly independent vectors in the orthogonal complement of  $C(X)$  (with respect to  $\mathbf{R}^6$ ).

**Exercise A.16** Let  $X$  be an  $n \times p$  matrix. Prove or disprove the following statement: Every vector in  $\mathbf{R}^n$  is in either  $C(X)$  or  $C(X)^\perp$  or both.

**Exercise A.17** For any matrix  $A$ , prove that  $C(A)$  and the null space of  $A'$  are orthogonal complements. Note: The null space is defined in Definition B.11.

# Appendix B

## Matrix Results

**Abstract** This appendix reviews standard ideas in matrix theory with emphasis given to important results that are less commonly taught in a junior/senior level linear algebra course. The appendix begins with basic definitions and results. A section devoted to eigenvalues and their applications follows. This section contains a number of standard definitions, but it also contains a number of very specific results that are unlikely to be familiar to people with only an undergraduate background in linear algebra. The third section is devoted to an intense (brief but detailed) examination of projections and their properties. The appendix closes with some miscellaneous results, some results on Kronecker products and Vec operators, and an introduction to tensors.

### B.1 Basic Ideas

**Definition B.1** Any matrix with the same number of rows and columns is called a *square matrix*.

**Definition B.2** Let  $A = [a_{ij}]$  be a matrix. The *transpose* of  $A$ , written  $A'$ , is the matrix  $A' = [b_{ij}]$ , where  $b_{ij} = a_{ji}$ .

**Definition B.3** If  $A = A'$ , then  $A$  is called *symmetric*. Note that only square matrices can be symmetric.

**Definition B.4** If  $A = [a_{ij}]$  is a square matrix and  $a_{ij} = 0$  for  $i \neq j$ , then  $A$  is a *diagonal matrix*. If  $\lambda_1, \dots, \lambda_n$  are scalars, then  $D(\lambda_j)$  and  $\text{Diag}(\lambda_j)$  are used to indicate an  $n \times n$  matrix  $D = [d_{ij}]$  with  $d_{ij} = 0$ ,  $i \neq j$ , and  $d_{ii} = \lambda_i$ . If  $\lambda \equiv (\lambda_1, \dots, \lambda_n)'$ , then  $D(\lambda) \equiv D(\lambda_j)$ . A diagonal matrix with all 1s on the diagonal is called an *identity matrix* and is denoted  $I$ . Occasionally,  $I_n$  is used to denote an  $n \times n$  identity matrix.

If  $A = [a_{ij}]$  is  $n \times p$  and  $B = [b_{ij}]$  is  $n \times q$ , we can write an  $n \times (p + q)$  matrix  $C = [A, B]$ , where  $c_{ij} = a_{ij}$ ,  $i = 1, \dots, n$ ,  $j = 1, \dots, p$ , and  $c_{ij} = b_{i, j-p}$ ,  $i = 1, \dots, n$ ,  $j = p + 1, \dots, p + q$ . This notation can be extended in obvious ways, e.g.,  $C' = \begin{bmatrix} A' \\ B' \end{bmatrix}$ .

**Definition B.5** Let  $A = [a_{ij}]$  be an  $r \times c$  matrix and  $B = [b_{ij}]$  be an  $s \times d$  matrix. The *Kronecker product* of  $A$  and  $B$ , written  $A \otimes B$ , is an  $r \times c$  matrix of  $s \times d$  matrices. The matrix in the  $i$ th row and  $j$ th column is  $a_{ij}B$ . In total,  $A \otimes B$  is an  $rs \times cd$  matrix.

**Definition B.6** Let  $A$  be an  $r \times c$  matrix. Write  $A = [A_1, A_2, \dots, A_c]$ , where  $A_i$  is the  $i$ th column of  $A$ . The *Vec* operator stacks the columns of  $A$  into an  $rc \times 1$  vector; thus,

$$[\text{Vec}(A)]' = [A'_1, A'_2, \dots, A'_c].$$

*Example B.7*

$$A = \begin{bmatrix} 1 & 4 \\ 2 & 5 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & 3 \\ 0 & 4 \end{bmatrix},$$

$$A \otimes B = \begin{bmatrix} 1 \begin{pmatrix} 1 & 3 \\ 0 & 4 \end{pmatrix} & 4 \begin{pmatrix} 1 & 3 \\ 0 & 4 \end{pmatrix} \\ 2 \begin{pmatrix} 1 & 3 \\ 0 & 4 \end{pmatrix} & 5 \begin{pmatrix} 1 & 3 \\ 0 & 4 \end{pmatrix} \end{bmatrix} = \begin{bmatrix} 1 & 3 & 4 & 12 \\ 0 & 4 & 0 & 16 \\ 2 & 6 & 5 & 15 \\ 0 & 8 & 0 & 20 \end{bmatrix},$$

$$\text{Vec}(A) = [1, 2, 4, 5]'$$

**Definition B.8** Let  $A$  be an  $n \times n$  matrix.  $A$  is *nonsingular* if there exists a matrix  $A^{-1}$  such that  $A^{-1}A = I = AA^{-1}$ . If no such matrix exists, then  $A$  is *singular*. If  $A^{-1}$  exists, it is called the *inverse* of  $A$ .

**Theorem B.9** An  $n \times n$  matrix  $A$  is nonsingular if and only if  $r(A) = n$ , i.e., the columns of  $A$  form a basis for  $\mathbf{R}^n$ .

**Corollary B.10**  $A_{n \times n}$  is singular if and only if there exists  $x \neq 0$  such that  $Ax = 0$ .

For any matrix  $A$ , the set of all  $x$  such that  $Ax = 0$  is easily seen to be a vector space.

**Definition B.11** The set of all  $x$  such that  $Ax = 0$  is called the *null space* of  $A$  and written  $\mathcal{N}(A)$ .

**Theorem B.12** If  $A$  is  $n \times n$  and  $r(A) = r$ , then the null space of  $A$  has rank  $n - r$ .

## B.2 Eigenvalues and Related Results

The material in this section deals with eigenvalues and eigenvectors either in the statements of the results or in their proofs. Again, this is meant to be a brief review of important concepts; but, in addition, there are a number of specific results that may be unfamiliar.

**Definition B.13** The scalar  $\lambda$  is an *eigenvalue* of  $A_{n \times n}$  if  $A - \lambda I$  is singular.  $\lambda$  is an eigenvalue of *multiplicity*  $s$  if the rank of the null space of  $A - \lambda I$  is  $s$ . A nonzero vector  $x$  is an *eigenvector* of  $A$  corresponding to the eigenvalue  $\lambda$  if  $x$  is in the null space of  $A - \lambda I$ , i.e., if  $Ax = \lambda x$ . Eigenvalues are also called *singular values* and *characteristic roots*.

For example,

$$\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = 3 \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

and

$$\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \begin{pmatrix} -1 \\ 1 \end{pmatrix} = 1 \begin{pmatrix} -1 \\ 1 \end{pmatrix}.$$

Combining the two equations gives

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

Note that if  $\lambda \neq 0$  is an eigenvalue of  $A$ , the eigenvectors corresponding to  $\lambda$  (along with the vector 0) form a subspace of  $C(A)$ . For example, if  $Ax_1 = \lambda x_1$  and  $Ax_2 = \lambda x_2$ , then  $A(x_1 + x_2) = \lambda(x_1 + x_2)$ , so the set of eigenvectors is closed under vector addition. Similarly, it is closed under scalar multiplication, so it forms a subspace (except that eigenvectors cannot be 0 and every subspace contains 0). If  $\lambda = 0$ , the subspace is the null space of  $A$ .

If  $A$  is a symmetric matrix, and  $\gamma$  and  $\lambda$  are distinct eigenvalues, then the eigenvectors corresponding to  $\lambda$  and  $\gamma$  are orthogonal. To see this, let  $x$  be an eigenvector for  $\lambda$  and  $y$  an eigenvector for  $\gamma$ . Then  $\lambda x'y = x'Ay = \gamma x'y$ , which can happen only if  $\lambda = \gamma$  or if  $x'y = 0$ . Since  $\lambda$  and  $\gamma$  are distinct, we have  $x'y = 0$ .

Let  $\lambda_1, \dots, \lambda_r$  be the distinct nonzero eigenvalues of a symmetric matrix  $A$  with respective multiplicities  $s(1), \dots, s(r)$ . Let  $v_{i1}, \dots, v_{is(i)}$  be a basis for the space of eigenvectors of  $\lambda_i$ . We want to show that  $v_{11}, v_{12}, \dots, v_{rs(r)}$  is a basis for  $C(A)$ .

Suppose  $v_{11}, v_{12}, \dots, v_{rs(r)}$  is not a basis. Since  $v_{ij} \in C(A)$  and the  $v_{ij}$ s are linearly independent, we can pick  $x \in C(A)$  with  $x \perp v_{ij}$  for all  $i$  and  $j$ . Note that since  $Av_{ij} = \lambda_i v_{ij}$ , we have  $(A)^p v_{ij} = (\lambda_i)^p v_{ij}$ . In particular,  $x'(A)^p v_{ij} = x'(\lambda_i)^p v_{ij} = (\lambda_i)^p x' v_{ij} = 0$ , so  $A^p x \perp v_{ij}$  for any  $i, j$ , and  $p$ . The vectors  $x, Ax, A^2x, \dots$  cannot all be linearly independent, so there exists a smallest value  $k \leq n$  such that

$$A^k x + b_{k-1} A^{k-1} x + \dots + b_0 x = 0.$$

Since there is a solution to this, for some real number  $\mu$  we can write the equation as

$$(A - \mu I)(A^{k-1} x + \gamma_{k-2} A^{k-2} x + \dots + \gamma_0 x) = 0,$$

and  $\mu$  is an eigenvalue. (See Exercise B.1.) An eigenvector for  $\mu$  is  $y = A^{k-1} x + \gamma_{k-2} A^{k-2} x + \dots + \gamma_0 x$ . Clearly,  $y \perp v_{ij}$  for any  $i$  and  $j$ . Since  $k$  was chosen as the smallest value to get linear dependence, we have  $y \neq 0$ . If  $\mu \neq 0$ ,  $y$  is an eigenvector that does not correspond to any of  $\lambda_1, \dots, \lambda_r$ , a contradiction. If  $\mu = 0$ , we have  $Ay = 0$ ; and since  $A$  is symmetric,  $y$  is a vector in  $C(A)$  that is orthogonal to every other vector in  $C(A)$ , i.e.,  $y'y = 0$  but  $y \neq 0$ , a contradiction. We have proven

**Theorem B.14** *If  $A$  is a symmetric matrix, then there exists a basis for  $C(A)$  consisting of eigenvectors of nonzero eigenvalues. If  $\lambda$  is a nonzero eigenvalue of multiplicity  $s$ , then the basis will contain  $s$  eigenvectors for  $\lambda$ .*

If  $\lambda$  is an eigenvalue of  $A$  with multiplicity  $s$ , then we can think of  $\lambda$  as being an eigenvalue  $s$  times. With this convention, the rank of  $A$  is the number of nonzero eigenvalues. The total number of eigenvalues is  $n$  if  $A$  is an  $n \times n$  matrix.

For a symmetric matrix  $A$ , if we use eigenvectors corresponding to the zero eigenvalue, we can get a basis for  $\mathbf{R}^n$  consisting of eigenvectors. We already have a basis for  $C(A)$ , and the eigenvectors of 0 are the null space of  $A$ . For  $A$  symmetric,  $C(A)$  and the null space of  $A$  are orthogonal complements. Let  $\lambda_1, \dots, \lambda_n$  be the eigenvalues of a symmetric matrix  $A$ . Let  $v_1, \dots, v_n$  denote a basis of eigenvectors for  $\mathbf{R}^n$ , with  $v_i$  being an eigenvector for  $\lambda_i$  for any  $i$ .

**Theorem B.15** *If  $A$  is symmetric, there exists an orthonormal basis for  $\mathbf{R}^n$  consisting of eigenvectors of  $A$ .*

*Proof* Assume  $\lambda_{i1} = \dots = \lambda_{ik}$  are all the  $\lambda_i$ s equal to any particular value  $\lambda$ , and let  $v_{i1}, \dots, v_{ik}$  be a basis for the space of eigenvectors for  $\lambda$ . By Gram–Schmidt there exists an orthonormal basis  $w_{i1}, \dots, w_{ik}$  for the space of eigenvectors corresponding to  $\lambda$ . If we do this for each distinct eigenvalue, we get a collection of orthonormal sets that form a basis for  $\mathbf{R}^n$ . Since, as we have seen, for  $\lambda_i \neq \lambda_j$ , any eigenvector for  $\lambda_i$  is orthogonal to any eigenvector for  $\lambda_j$ , the basis is orthonormal.  $\square$

**Definition B.16** A square matrix  $P$  is *orthonormal* (more often called *orthogonal*) if  $P' = P^{-1}$ . Note that if  $P$  is orthonormal, so is  $P'$ .

Some examples of orthonormal matrices are

$$P_1 = \frac{1}{\sqrt{6}} \begin{bmatrix} \sqrt{2} & -\sqrt{3} & 1 \\ \sqrt{2} & 0 & -2 \\ \sqrt{2} & \sqrt{3} & 1 \end{bmatrix}, \quad P_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix},$$

$$P_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

**Theorem B.17**  $P_{n \times n}$  is orthonormal if and only if the columns of  $P$  form an orthonormal basis for  $\mathbf{R}^n$ .

*Proof*  $\Leftarrow$  It is clear that if the columns of  $P$  form an orthonormal basis for  $\mathbf{R}^n$ , then  $P'P = I$ .

$\Rightarrow$  Since  $P$  is nonsingular, the columns of  $P$  form a basis for  $\mathbf{R}^n$ . Since  $P'P = I$ , the basis is orthonormal.  $\square$

**Corollary B.18**  $P_{n \times n}$  is orthonormal if and only if the rows of  $P$  form an orthonormal basis for  $\mathbf{R}^n$ .

*Proof*  $P$  is orthonormal if and only if  $P'$  is orthonormal if and only if the columns of  $P'$  are an orthonormal basis if and only if the rows of  $P$  are an orthonormal basis.  $\square$

**Theorem B.19** If  $A$  is an  $n \times n$  symmetric matrix, then there exists an orthonormal matrix  $P$  such that  $P'AP = \text{Diag}(\lambda_i)$ , where  $\lambda_1, \lambda_2, \dots, \lambda_n$  are the eigenvalues of  $A$ .

*Proof* Let  $v_1, v_2, \dots, v_n$  be an orthonormal set of eigenvectors of  $A$  corresponding, respectively, to  $\lambda_1, \lambda_2, \dots, \lambda_n$ . Let  $P = [v_1, \dots, v_n]$ . Then

$$\begin{aligned} P'AP &= \begin{bmatrix} v'_1 \\ \vdots \\ v'_n \end{bmatrix} [Av_1, \dots, Av_n] \\ &= \begin{bmatrix} v'_1 \\ \vdots \\ v'_n \end{bmatrix} [\lambda_1 v_1, \dots, \lambda_n v_n] \\ &= \begin{bmatrix} \lambda_1 v'_1 v_1 & \dots & \lambda_n v'_1 v_n \\ \vdots & \ddots & \vdots \\ \lambda_1 v'_n v_1 & \dots & \lambda_n v'_n v_n \end{bmatrix} \\ &= \text{Diag}(\lambda_i). \end{aligned}$$

$\square$

The *singular value decomposition* for a symmetric matrix is given by the following corollary.

**Corollary B.20**  $A = PD(\lambda_i)P'$ .

For example, using results illustrated earlier,

$$\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} = \begin{bmatrix} 1/\sqrt{2} & -1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix} \begin{bmatrix} 3 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ -1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}.$$

**Definition B.21** A symmetric matrix  $A$  is *positive (nonnegative) definite* if, for any nonzero vector  $v \in \mathbf{R}^n$ ,  $v'Av$  is positive (nonnegative).

**Theorem B.22**  $A$  is *nonnegative definite* if and only if there exists a square matrix  $Q$  such that  $A = QQ'$ .

*Proof*  $\Rightarrow$  We know that there exists  $P$  orthonormal with  $P'AP = \text{Diag}(\lambda_i)$ . The  $\lambda_i$ s must all be nonnegative, because if  $e'_j = (0, \dots, 0, 1, 0, \dots, 0)$  with the 1 in the  $j$ th place and we let  $v = Pe_j$ , then  $0 \leq v'Av = e'_j \text{Diag}(\lambda_i)e_j = \lambda_j$ . Let  $Q = P\text{Diag}(\sqrt{\lambda_i})$ . Then, since  $P\text{Diag}(\lambda_i)P' = A$ , we have

$$QQ' = P\text{Diag}(\lambda_i)P' = A.$$

$\Leftarrow$  If  $A = QQ'$ , then  $v'Av = (Q'v)'(Q'v) \geq 0$ . □

**Corollary B.23**  $A$  is *positive definite* if and only if  $Q$  is nonsingular for any choice of  $Q$ .

*Proof* There exists  $v \neq 0$  such that  $v'Av = 0$  if and only if there exists  $v \neq 0$  such that  $Q'v = 0$ , which occurs if and only if  $Q'$  is singular. The contrapositive of this is that  $v'Av > 0$  for all  $v \neq 0$  if and only if  $Q'$  is nonsingular. □

In the interest of brevity, I have dropped Theorem B.24, Corollary B.25, and Corollary B.26 that appeared in earlier editions.

**Definition B.27** Let  $A = [a_{ij}]$  be an  $n \times n$  matrix. The *trace* of  $A$  is  $\text{tr}(A) = \sum_{i=1}^n a_{ii}$ .

**Theorem B.28** For matrices  $A_{r \times s}$  and  $B_{s \times r}$ ,  $\text{tr}(AB) = \text{tr}(BA)$ .



*Proof* See Exercise B.8.  $\square$

**Theorem B.29** If  $A_{n \times n}$  is a symmetric matrix,  $\text{tr}(A) = \sum_{i=1}^n \lambda_i$ , where  $\lambda_1, \dots, \lambda_n$  are the eigenvalues of  $A$ .

*Proof*  $A = PD(\lambda_i)P'$  with  $P$  orthonormal

$$\begin{aligned}\text{tr}(A) &= \text{tr}[PD(\lambda_i)P'] = \text{tr}[D(\lambda_i)P'P] \\ &= \text{tr}[D(\lambda_i)] = \sum_{i=1}^n \lambda_i.\end{aligned}$$

$\square$

To illustrate, we saw earlier that the matrix  $\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$  had eigenvalues of 3 and 1. In fact, a stronger result than Theorem B.29 is true. We give it without proof.

**Theorem B.30**  $\text{tr}(A) = \sum_{i=1}^n \lambda_i$ , where  $\lambda_1, \dots, \lambda_n$  are the eigenvalues of  $A$ . Moreover, the determinant of  $A$  is  $\det(A) = \prod_{i=1}^n \lambda_i$ .

### B.3 Projections

This section is devoted primarily to a discussion of perpendicular projection operators. It begins with their definition, some basic properties, and two important characterizations: Theorems B.33 and B.35. A third important characterization, Theorem B.44, involves generalized inverses. Generalized inverses are defined, briefly studied, and applied to projection operators. The section continues with the examination of the relationships between two perpendicular projection operators and closes with discussions of the Gram–Schmidt theorem, eigenvalues of projection operators, and oblique (nonperpendicular) projection operators.

We begin by defining a *perpendicular projection operator* (*ppo*) onto an arbitrary space. To be consistent with later usage, we denote the arbitrary space  $C(X)$  for some matrix  $X$ .

**Definition B.31**  $M$  is a perpendicular projection operator (matrix) onto  $C(X)$  if and only if

- (i)  $v \in C(X)$  implies  $Mv = v$  (projection),
- (ii)  $w \perp C(X)$  implies  $Mw = 0$  (perpendicularity).

For example, consider the subspace of  $\mathbf{R}^2$  determined by vectors of the form  $(2a, a)'$ . It is not difficult to see that the orthogonal complement of this subspace

consists of vectors of the form  $(b, -2b)'$ . The perpendicular projection operator onto the  $(2a, a)'$  subspace is

$$M = \begin{bmatrix} 0.8 & 0.4 \\ 0.4 & 0.2 \end{bmatrix}.$$

To verify this note that

$$M \begin{pmatrix} 2a \\ a \end{pmatrix} = \begin{bmatrix} 0.8 & 0.4 \\ 0.4 & 0.2 \end{bmatrix} \begin{pmatrix} 2a \\ a \end{pmatrix} = \begin{pmatrix} (0.8)2a + 0.4a \\ (0.4)2a + 0.2a \end{pmatrix} = \begin{pmatrix} 2a \\ a \end{pmatrix}$$

and

$$M \begin{pmatrix} b \\ -2b \end{pmatrix} = \begin{bmatrix} 0.8 & 0.4 \\ 0.4 & 0.2 \end{bmatrix} \begin{pmatrix} b \\ -2b \end{pmatrix} = \begin{pmatrix} 0.8b + 0.4(-2b) \\ 0.4b + 0.2(-2b) \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}.$$

Notationally,  $M$  is used to indicate the ppo onto  $C(X)$ . If  $A$  is another matrix,  $M_A$  denotes the ppo onto  $C(A)$ . Thus,  $M \equiv M_X$ . When  $X$  has a subscript we typically write the ppo onto  $C(X_0)$  as  $M_0 \equiv M_{X_0}$  and, similarly,  $M_1 \equiv M_{X_1}$ , but often  $M_2 \neq M_{X_2}$ .

**Proposition B.32** *If  $M$  is a perpendicular projection operator onto  $C(X)$ , then  $C(M) = C(X)$ .*

*Proof* See Exercise B.2. □

Note that both columns of

$$M = \begin{bmatrix} 0.8 & 0.4 \\ 0.4 & 0.2 \end{bmatrix}$$

have the form  $(2a, a)'$ .

**Theorem B.33**  *$M$  is a perpendicular projection operator on  $C(M)$  if and only if  $MM = M$  and  $M' = M$ .*

*Proof*  $\Rightarrow$  Write  $v = v_1 + v_2$ , where  $v_1 \in C(M)$  and  $v_2 \perp C(M)$ , and let  $w = w_1 + w_2$  with  $w_1 \in C(M)$  and  $w_2 \perp C(M)$ . Since  $(I - M)v = (I - M)v_2 = v_2$  and  $Mw = Mw_1 = w_1$ , we get

$$w'M'(I - M)v = w'_1M'(I - M)v_2 = w'_1v_2 = 0.$$

This is true for any  $v$  and  $w$ , so we have  $M'(I - M) = 0$  or  $M' = M'M$ . Since  $M'M$  is symmetric,  $M'$  must also be symmetric, and this implies that  $M = MM$ .

$\Leftarrow$  If  $M^2 = M$  and  $v \in C(M)$ , then since  $v = Mb$  we have  $Mv = MMb = Mb = v$ . If  $M' = M$  and  $w \perp C(M)$ , then  $Mw = M'w = 0$  because the columns of  $M$  are in  $C(M)$ . □

In our example,

$$MM = \begin{bmatrix} 0.8 & 0.4 \\ 0.4 & 0.2 \end{bmatrix} \begin{bmatrix} 0.8 & 0.4 \\ 0.4 & 0.2 \end{bmatrix} = \begin{bmatrix} 0.8 & 0.4 \\ 0.4 & 0.2 \end{bmatrix} = M$$

and

$$M = \begin{bmatrix} 0.8 & 0.4 \\ 0.4 & 0.2 \end{bmatrix} = M'.$$

**Proposition B.34** *Perpendicular projection operators are unique.*

*Proof* Let  $M$  and  $P$  be perpendicular projection operators onto some space  $\mathcal{M}$ . Let  $v \in \mathbf{R}^n$  and write  $v = v_1 + v_2$ ,  $v_1 \in \mathcal{M}$ ,  $v_2 \perp \mathcal{M}$ . Since  $v$  is arbitrary and  $Mv = v_1 = Pv$ , we have  $M = P$ .  $\square$

For any matrix  $X$ , we will now find two ways to characterize the perpendicular projection operator onto  $C(X)$ . The first method depends on the Gram–Schmidt theorem; the second depends on the concept of a generalized inverse.

**Theorem B.35** *Let  $o_1, \dots, o_r$  be an orthonormal basis for  $C(X)$ , and let  $O = [o_1, \dots, o_r]$ . Then  $OO' = \sum_{i=1}^r o_i o_i'$  is the perpendicular projection operator onto  $C(X)$ .*

*Proof*  $OO'$  is symmetric and  $OO'OO' = OO'I_r O' = OO'$ ; so, by Theorem B.33, it only remains to show that  $C(OO') = C(X)$ . Clearly  $C(OO') \subset C(O) = C(X)$ . On the other hand, if  $v \in C(O)$ , then  $v = Ob$  for some vector  $b \in \mathbf{R}^r$  and  $v = Ob = OO'r b = OO'Ob$ ; so clearly  $v \in C(OO')$ .  $\square$

For example, to find the perpendicular projection operator for vectors of the form  $(2a, a)'$ , we can find an orthonormal basis. The space has rank 1 and to normalize  $(2a, a)'$ , we must have

$$1 = (2a, a) \begin{pmatrix} 2a \\ a \end{pmatrix} = 4a^2 + a^2 = 5a^2;$$

so  $a^2 = 1/5$  and  $a = \pm 1/\sqrt{5}$ . If we take  $(2/\sqrt{5}, 1/\sqrt{5})'$  as our orthonormal basis, then

$$M = \begin{pmatrix} 2/\sqrt{5} \\ 1/\sqrt{5} \end{pmatrix} (2/\sqrt{5}, 1/\sqrt{5}) = \begin{bmatrix} 0.8 & 0.4 \\ 0.4 & 0.2 \end{bmatrix},$$

as was demonstrated earlier.

One use of Theorem B.35 is that, given a matrix  $X$ , one can use the Gram–Schmidt theorem to get an orthonormal basis for  $C(X)$  and thus obtain the perpendicular projection operator.

We now examine properties of generalized inverses. Generalized inverses are a generalization on the concept of the inverse of a matrix. Although the most common use of generalized inverses is in solving systems of linear equations, our interest lies

primarily in their relationship to projection operators. The discussion below is given for an arbitrary matrix  $A$ .

**Definition B.36** A *generalized inverse* of a matrix  $A$  is any matrix  $G$  such that  $AGA = A$ . The notation  $A^-$  is used to indicate a generalized inverse of  $A$ .

**Theorem B.37** If  $A$  is nonsingular, the unique generalized inverse of  $A$  is  $A^{-1}$ .

*Proof*  $AA^{-1}A = IA = A$ , so  $A^{-1}$  is a generalized inverse. If  $AA^-A = A$ , then  $AA^- = AA^-AA^{-1} = AA^{-1} = I$ ; so  $A^-$  is the inverse of  $A$ .  $\square$

**Theorem B.38** For any symmetric matrix  $A$ , there exists a generalized inverse of  $A$ .

*Proof* There exists  $P$  orthonormal so that  $P'AP = D(\lambda_i)$  and  $A = PD(\lambda_i)P'$ . Let

$$\gamma_i = \begin{cases} 1/\lambda_i, & \text{if } \lambda_i \neq 0 \\ 0, & \text{if } \lambda_i = 0, \end{cases}$$

and  $G = PD(\gamma_i)P'$ . We now show that  $G$  is a generalized inverse of  $A$ .  $P$  is orthonormal, so  $P'P = I$  and

$$\begin{aligned} AGA &= PD(\lambda_i)P'PD(\gamma_i)P'PD(\lambda_i)P' \\ &= PD(\lambda_i)D(\gamma_i)D(\lambda_i)P' \\ &= PD(\lambda_i)P' \\ &= A. \end{aligned} \quad \square$$

Although this is the only existence result we really need, later we will show that generalized inverses exist for arbitrary matrices.

**Theorem B.39** If  $G_1$  and  $G_2$  are generalized inverses of  $A$ , then so is  $G_1AG_2$ .

*Proof*  $A(G_1AG_2)A = (AG_1A)G_2A = AG_2A = A$ .  $\square$

For  $A$  symmetric,  $A^-$  need not be symmetric.

*Example B.40* Consider the matrix

$$\begin{bmatrix} a & b \\ b & b^2/a \end{bmatrix}.$$

It has a generalized inverse

$$\begin{bmatrix} 1/a & -1 \\ 1 & 0 \end{bmatrix},$$

and in fact, by considering the equation

$$\begin{bmatrix} a & b \\ b & b^2/a \end{bmatrix} \begin{bmatrix} r & s \\ t & u \end{bmatrix} \begin{bmatrix} a & b \\ b & b^2/a \end{bmatrix} = \begin{bmatrix} a & b \\ b & b^2/a \end{bmatrix},$$

it can be shown that if  $r = 1/a$ , then any solution of  $at + as + bu = 0$  gives a generalized inverse.

**Corollary B.41** For a symmetric matrix  $A$ , there exists  $A^-$  such that  $A^-AA^- = A^-$  and  $(A^-)' = A^-$ .

*Proof* Take  $A^-$  as the generalized inverse in the proof of Theorem B.38. Clearly,  $A^- = PD(\gamma_i)P'$  is symmetric and

$$A^-AA^- = PD(\gamma_i)P'PD(\lambda_i)P'PD(\gamma_i)P' = PD(\gamma_i)D(\lambda_i)D(\gamma_i)P' = PD(\gamma_i)P' = A^-.$$

□

**Definition B.42** A generalized inverse  $A^-$  for a matrix  $A$  that has the property  $A^-AA^- = A^-$  is said to be *reflexive*.

Corollary B.41 establishes the existence of a reflexive generalized inverse for any symmetric matrix.

Generalized inverses are of interest in that they provide an alternative to the characterization of perpendicular projection matrices given in Theorem B.35. The two results immediately below characterize the perpendicular projection matrix onto  $C(X)$ .

**Lemma B.43** If  $G$  and  $H$  are generalized inverses of  $(X'X)$ , then

- (i)  $XGX'X = XHX'X = X$ ,
- (ii)  $XGX' = XHX'$ .

*Proof* For  $v \in \mathbf{R}^n$ , let  $v = v_1 + v_2$  with  $v_1 \in C(X)$  and  $v_2 \perp C(X)$ . Also let  $v_1 = Xb$  for some vector  $b$ . Then

$$v'XGX'X = v_1'XGX'X = b'(X'X)G(X'X) = b'(X'X) = v_1'X.$$

Since  $v$  and  $G$  are arbitrary, we have shown (i).

To see (ii), observe that for the arbitrary vector  $v$  above,

$$XGX'v = XGX'Xb = XHX'Xb = XHX'v.$$

□

Since  $X'X$  is symmetric, there exists a generalized inverse  $(X'X)^-$  that is symmetric. For this generalized inverse,  $X(X'X)^-X'$  is symmetric; so, by the above lemma,  $X(X'X)^-X'$  must be symmetric for any choice of  $(X'X)^-$ .

**Theorem B.44**  $X(X'X)^-X'$  is the perpendicular projection operator onto  $C(X)$ .

*Proof* We need to establish conditions (i) and (ii) of Definition B.31. (i) For  $v \in C(X)$ , write  $v = Xb$ , so by Lemma B.43,  $X(X'X)^-X'v = X(X'X)^-X'Xb = Xb = v$ . (ii) If  $w \perp C(X)$ ,  $X(X'X)^-X'w = 0$ .  $\square$

For example, one spanning set for the subspace of vectors with the form  $(2a, a)'$  is  $(2, 1)'$ . It follows that

$$M = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \left[ (2, 1) \begin{pmatrix} 2 \\ 1 \end{pmatrix} \right]^{-1} (2, 1) = \begin{bmatrix} 0.8 & 0.4 \\ 0.4 & 0.2 \end{bmatrix},$$

as was shown earlier.

The next five results examine the relationships between two perpendicular projection matrices.

**Theorem B.45** Let  $M_1$  and  $M_2$  be perpendicular projection matrices on  $\mathbf{R}^n$ .  $(M_1 + M_2)$  is the perpendicular projection matrix onto  $C(M_1, M_2)$  if and only if  $C(M_1) \perp C(M_2)$ .

*Proof*  $\Leftarrow$  If  $C(M_1) \perp C(M_2)$ , then  $M_1M_2 = M_2M_1 = 0$ . Because

$$(M_1 + M_2)^2 = M_1^2 + M_2^2 + M_1M_2 + M_2M_1 = M_1^2 + M_2^2 = M_1 + M_2$$

and

$$(M_1 + M_2)' = M_1' + M_2' = M_1 + M_2,$$

$M_1 + M_2$  is the perpendicular projection matrix onto  $C(M_1 + M_2)$ . Clearly  $C(M_1 + M_2) \subset C(M_1, M_2)$ . To see that  $C(M_1, M_2) \subset C(M_1 + M_2)$ , write  $v = M_1b_1 + M_2b_2$ . Then, because  $M_1M_2 = M_2M_1 = 0$ ,  $(M_1 + M_2)v = v$ . Thus,  $C(M_1, M_2) = C(M_1 + M_2)$ .

$\Rightarrow$  If  $M_1 + M_2$  is a perpendicular projection matrix, then

$$\begin{aligned} (M_1 + M_2) &= (M_1 + M_2)^2 = M_1^2 + M_2^2 + M_1M_2 + M_2M_1 \\ &= M_1 + M_2 + M_1M_2 + M_2M_1. \end{aligned}$$

Thus,  $M_1M_2 + M_2M_1 = 0$ .

Multiplying by  $M_1$  gives  $0 = M_1^2M_2 + M_1M_2M_1 = M_1M_2 + M_1M_2M_1$  and thus  $-M_1M_2M_1 = M_1M_2$ . Since  $-M_1M_2M_1$  is symmetric, so is  $M_1M_2$ . This gives  $M_1M_2 = (M_1M_2)' = M_2M_1$ , so the condition  $M_1M_2 + M_2M_1 = 0$  becomes

$2(M_1M_2) = 0$  or  $M_1M_2 = 0$ . By symmetry, this says that the columns of  $M_1$  are orthogonal to the columns of  $M_2$ .  $\square$

**Theorem B.46** *If  $M_1$  and  $M_2$  are symmetric,  $C(M_1) \perp C(M_2)$ , and  $(M_1 + M_2)$  is a perpendicular projection matrix, then  $M_1$  and  $M_2$  are perpendicular projection matrices.*

*Proof*

$$(M_1 + M_2) = (M_1 + M_2)^2 = M_1^2 + M_2^2 + M_1M_2 + M_2M_1.$$

Since  $M_1$  and  $M_2$  are symmetric with  $C(M_1) \perp C(M_2)$ , we have  $M_1M_2 + M_2M_1 = 0$  and  $M_1 + M_2 = M_1^2 + M_2^2$ . Rearranging gives  $M_2 - M_2^2 = M_1^2 - M_1$ , so  $C(M_2 - M_2^2) = C(M_1^2 - M_1)$ . Now  $C(M_2 - M_2^2) \subset C(M_2)$  and  $C(M_1^2 - M_1) \subset C(M_1)$ , so  $C(M_2 - M_2^2) \perp C(M_1^2 - M_1)$ . The only way a vector space can be orthogonal to itself is if it consists only of the zero vector. Thus,  $M_2 - M_2^2 = M_1^2 - M_1 = 0$ , and  $M_2 = M_2^2$  and  $M_1 = M_1^2$ .  $\square$

**Theorem B.47** *Let  $M$  and  $M_0$  be perpendicular projection matrices with  $C(M_0) \subset C(M)$ . Then  $M - M_0$  is the perpendicular projection matrix onto  $C(M_0)_{C(M)}^\perp$ .*

*Proof* Since  $C(M_0) \subset C(M)$ ,  $MM_0 = M_0$  and, by symmetry,  $M_0M = M_0$ . Checking the conditions of Theorem B.33, we see that  $(M - M_0)^2 = M^2 - MM_0 - M_0M + M_0^2 = M - M_0 - M_0 + M_0 = M - M_0$ , and  $(M - M_0)' = M - M_0$ , so  $M - M_0$  is a ppo onto  $C(M - M_0)$ .

To see that  $C(M - M_0) = C(M_0)_{C(M)}^\perp$  note that  $C(M - M_0) \perp C(M_0)$ , because  $(M - M_0)M_0 = MM_0 - M_0^2 = M_0 - M_0 = 0$ . Thus,  $C(M - M_0) \subset C(M_0)_{C(M)}^\perp$ . If  $x \in C(M)$  and  $x \perp C(M_0)$ , then  $x = Mx = (M - M_0)x + M_0x = (M - M_0)x$ . Thus,  $x \in C(M - M_0)$  and  $C(M_0)_{C(M)}^\perp \subset C(M - M_0)$ .  $\square$

**Corollary B.48**  $C(M - M_0) = C(M_0)_{C(M)}^\perp$ .

**Corollary B.49**  $r(M) = r(M_0) + r(M - M_0)$ .

One particular application of these results involves  $I$ , the perpendicular projection operator onto  $\mathbf{R}^n$ . For any other perpendicular projection operator  $M$ ,  $I - M$  is the perpendicular projection operator onto the orthogonal complement of  $C(M)$  with respect to  $\mathbf{R}^n$ . For example, the subspace of vectors with the form  $(2a, a)'$  has an orthogonal complement consisting of vectors with the form  $(b, -2b)'$ . With  $M$  as given earlier,

$$I - M = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - \begin{bmatrix} 0.8 & 0.4 \\ 0.4 & 0.2 \end{bmatrix} = \begin{bmatrix} 0.2 & -0.4 \\ -0.4 & 0.8 \end{bmatrix}.$$

Note that

$$(I - M) \begin{pmatrix} b \\ -2b \end{pmatrix} = \begin{pmatrix} b \\ -2b \end{pmatrix} \quad \text{and} \quad (I - M) \begin{pmatrix} 2a \\ a \end{pmatrix} = 0;$$

so by definition  $I - M$  is the perpendicular projection operator onto the space of vectors with the form  $(b, -2b)'$ .

At this point, we examine the relationship between perpendicular projection operations and the Gram–Schmidt theorem (Theorem A.16). Recall that in the Gram–Schmidt theorem,  $x_1, \dots, x_r$  denotes the original basis and  $y_1, \dots, y_r$  denotes the orthonormal basis. Let

$$M_s = \sum_{i=1}^s y_i y_i'.$$

Applying Theorem B.35,  $M_s$  is the ppo onto  $C(x_1, \dots, x_s)$ . Now define

$$w_{s+1} = (I - M_s)x_{s+1}.$$

Thus,  $w_{s+1}$  is the perpendicular projection of  $x_{s+1}$  onto the orthogonal complement of  $C(x_1, \dots, x_s)$ . Finally,  $y_{s+1}$  is just  $w_{s+1}$  normalized.

Consider the eigenvalues of a perpendicular projection operator  $M$ . Let  $v_1, \dots, v_r$  be a basis for  $C(M)$ . Then  $Mv_i = v_i$ , so  $v_i$  is an eigenvector of  $M$  with eigenvalue 1. In fact, 1 is an eigenvalue of  $M$  with multiplicity  $r$ . Now, let  $w_1, \dots, w_{n-r}$  be a basis for  $C(M)^\perp$ .  $Mw_j = 0$ , so 0 is an eigenvalue of  $M$  with multiplicity  $n - r$ . We have completely characterized the  $n$  eigenvalues of  $M$ . Since  $\text{tr}(M)$  equals the sum of the eigenvalues, we have  $\text{tr}(M) = r(M)$ .

In fact, if  $A$  is an  $n \times n$  matrix with  $A^2 = A$ , any basis for  $C(A)$  is a basis for the space of eigenvectors for the eigenvalue 1. The null space of  $A$  is the space of eigenvectors for the eigenvalue 0. The rank of  $A$  and the rank of the null space of  $A$  add to  $n$ , and  $A$  has  $n$  eigenvalues, so all the eigenvalues are accounted for. Again,  $\text{tr}(A) = r(A)$ .

### Definition B.50

- (a) If  $A$  is a square matrix with  $A^2 = A$ , then  $A$  is called *idempotent*.
- (b) Let  $\mathcal{N}$  and  $\mathcal{M}$  be two spaces with  $\mathcal{N} \cap \mathcal{M} = \{0\}$  and  $r(\mathcal{N}) + r(\mathcal{M}) = n$ . The square matrix  $A$  is a *projection operator* onto  $\mathcal{N}$  along  $\mathcal{M}$  if 1)  $Av = v$  for any  $v \in \mathcal{N}$ , and 2)  $Aw = 0$  for any  $w \in \mathcal{M}$ .

If the square matrix  $A$  has the property that  $Av = v$  for any  $v \in C(A)$ , then  $A$  is the projection operator (matrix) onto  $C(A)$  along  $C(A)^\perp$ . (Note that  $C(A)^\perp$  is the null space of  $A$ .) It follows immediately that if  $A$  is idempotent, then  $A$  is a projection operator onto  $C(A)$  along  $\mathcal{N}(A) = C(A)^\perp = C(I - A)$ , see Exercise B.22.

The uniqueness of projection operators can be established like it was for perpendicular projection operators. Note that  $x \in \mathbf{R}^n$  can be written uniquely as  $x = v + w$



for  $v \in \mathcal{N}$  and  $w \in \mathcal{M}$ , i.e.,  $\mathbf{R}^n = \mathcal{N} + \mathcal{M}$ . To see this, take basis matrices for the two spaces, say  $N$  and  $M$ , respectively. The result follows from observing that  $[N, M]$  is a basis matrix for  $\mathbf{R}^n$ . Because of the rank conditions,  $[N, M]$  is an  $n \times n$  matrix. It is enough to show that the columns of  $[N, M]$  must be linearly independent.

$$0 = [N, M] \begin{bmatrix} b \\ c \end{bmatrix} = Nb + Mc$$

implies  $Nb = M(-c)$  which, since  $\mathcal{N} \cap \mathcal{M} = \{0\}$ , can only happen when  $Nb = 0 = M(-c)$ , which, because they are basis matrices, can only happen when  $b = 0 = (-c)$ , which implies that  $\begin{bmatrix} b \\ c \end{bmatrix} = 0$ , and we are done.

Any projection operator that is not a perpendicular projection is referred to as an *oblique projection operator*.

To show that a matrix  $A$  is a projection operator onto an arbitrary space, say  $C(X)$ , it is necessary to show that  $C(A) = C(X)$  and that for  $x \in C(X)$ ,  $Ax = x$ . A typical proof runs in the following pattern. First, show that  $Ax = x$  for any  $x \in C(X)$ . This also establishes that  $C(X) \subset C(A)$ . To finish the proof, it suffices to show that  $Av \in C(X)$  for any  $v \in \mathbf{R}^n$  because this implies that  $C(A) \subset C(X)$ .

In this book, our use of the word “perpendicular” is based on the standard inner product that defines Euclidean distance. In other words, for two vectors  $x$  and  $y$ , their inner product is  $x'y$ . By definition, the vectors  $x$  and  $y$  are orthogonal if their inner product is 0. In fact, for any two vectors  $x$  and  $y$ , let  $\theta$  be the angle between  $x$  and  $y$ . Then  $x'y = \sqrt{x'x}\sqrt{y'y} \cos \theta$ . The length of a vector  $x$  is defined as the square root of the inner product of  $x$  with itself, i.e.,  $\|x\| \equiv \sqrt{x'x}$ . The distance between two vectors  $x$  and  $y$  is the length of their difference, i.e.,  $\|x - y\|$ .

These concepts can be generalized. For a positive definite matrix  $B$ , we can define an inner product between  $x$  and  $y$  as  $x'By$ . As before,  $x$  and  $y$  are orthogonal if their inner product is 0 and the length of  $x$  is the square root of its inner product with itself (now  $\|x\|_B \equiv \sqrt{x'Bx}$ ). As argued above, any idempotent matrix is always a projection operator, but which one is the perpendicular projection operator depends on the inner product. As can be seen from Proposition 2.7.2 and Exercise 2.5, the matrix  $X(X'BX)^{-1}X'B$  is an oblique projection onto  $C(X)$  for the standard inner product; but it is the perpendicular projection operator onto  $C(X)$  with the inner product defined using the matrix  $B$ .

## B.4 Miscellaneous Results

**Proposition B.51** For any matrix  $X$ ,  $C(XX') = C(X)$ .

*Proof* Clearly  $C(XX') \subset C(X)$ , so we need to show that  $C(X) \subset C(XX')$ . Let  $x \in C(X)$ . Then  $x = Xb$  for some  $b$ . Write  $b = b_0 + b_1$ , where  $b_0 \in C(X')$  and  $b_1 \perp C(X')$ . Clearly,  $Xb_1 = 0$ , so we have  $x = Xb_0$ . But  $b_0 = X'd$  for some  $d$ ; so  $x = Xb_0 = XX'd$  and  $x \in C(XX')$ . □

**Corollary B.52** For any matrix  $X$ ,  $r(XX') = r(X)$ .

*Proof* See Exercise B.4. □

**Corollary B.53** If  $X_{n \times p}$  has  $r(X) = p$ , then the  $p \times p$  matrix  $X'X$  is nonsingular.

*Proof* See Exercise B.5. □

**Proposition B.54**

- (a) If  $C(U_1) \subset C(U_2)$ , then  $C(XU_1) \subset C(XU_2)$ .
- (b) If  $C(U_1) = C(U_2)$ , then  $C(XU_1) = C(XU_2)$ .
- (c)  $C(XB) \subset C(X)$
- (d) If  $B$  is nonsingular,  $C(XB) = C(X)$ .

*Proof* (a) Take  $v \in C(XU_1)$ . For some  $\gamma_1$ ,  $v = XU_1\gamma_1$ . Because,  $C(U_1) \subset C(U_2)$ , there exists  $\gamma_2$  so that  $U_1\gamma_1 = U_2\gamma_2$ . Clearly,  $v = XU_1\gamma_1 = XU_2\gamma_2 \in C(XU_2)$ . (b) Use (a) as is and with the roles of  $U_1$  and  $U_2$  reversed. (c) This is immediate from the definition of a column space but also, in a) take  $U_1 = B$  and  $U_2 = I$ . (d) If  $B$  is nonsingular,  $C(B) = C(I)$  and use (b). □

It follows immediately from Proposition B.54 that, for  $B$  nonsingular, the perpendicular projection operators onto  $C(XB)$  and  $C(X)$  are identical.

We now show that generalized inverses always exist.

**Theorem B.55** For any matrix  $X$ , there exists a generalized inverse  $X^-$ .

*Proof* We know that  $(X'X)^-$  exists. Set  $X^- = (X'X)^-X'$ . Then  $XX^-X = X(X'X)^-X'X = X$  because  $X(X'X)^-X'$  is a projection matrix onto  $C(X)$ . □

Note that for any  $X^-$ , the matrix  $XX^-$  is idempotent and hence a projection operator.

**Proposition B.56** When all inverses exist,

$$[A + BCD]^{-1} = A^{-1} - A^{-1}B [C^{-1} + DA^{-1}B]^{-1} DA^{-1}.$$

*Proof* If all inverses exist

$$\begin{aligned} & [A + BCD] \left[ A^{-1} - A^{-1}B [C^{-1} + DA^{-1}B]^{-1} DA^{-1} \right] \\ &= I - B [C^{-1} + DA^{-1}B]^{-1} DA^{-1} + BCDA^{-1} \\ &\quad - BCDA^{-1}B [C^{-1} + DA^{-1}B]^{-1} DA^{-1} \\ &= I - B [I + CDA^{-1}B] [C^{-1} + DA^{-1}B]^{-1} DA^{-1} + BCDA^{-1} \\ &= I - BC [C^{-1} + DA^{-1}B] [C^{-1} + DA^{-1}B]^{-1} DA^{-1} + BCDA^{-1} \\ &= I - BCDA^{-1} + BCDA^{-1} = I. \end{aligned}$$
□

**Proposition B.57** *Let  $P$  be a projection operator (idempotent), and let  $a$  and  $b$  be real numbers. Then*

$$[aI + bP]^{-1} = \frac{1}{a} \left[ I - \frac{b}{a+b} P \right].$$

*Proof*

$$\frac{1}{a} \left[ I - \frac{b}{a+b} P \right] [aI + bP] = \frac{1}{a} \left[ aI + bP - \frac{ab}{a+b} P - \frac{b^2}{a+b} P \right] = I. \quad \square$$

When we study linear models, we frequently need to refer to matrices and vectors that consist entirely of 1s. Such matrices are denoted by the letter  $J$  with various subscripts and superscripts to specify their dimensions.  $J_r^c$  is an  $r \times c$  matrix of 1s. The subscript indicates the number of rows and the superscript indicates the number of columns. If there is only one column, the superscript may be suppressed, e.g.,  $J_r \equiv J_r^1$ . In a context where we are dealing with vectors in  $\mathbf{R}^n$ , the subscript may also be suppressed, e.g.,  $J \equiv J_n \equiv J_n^1$ .

A matrix of 0s is always denoted by 0.

## B.5 Properties of Kronecker Products and Vec Operators

Kronecker products and Vec operators are extremely useful in multivariate analysis and some approaches to variance component estimation. (Both are discussed in *ALM-III*.) They are also often used in writing balanced ANOVA models. We now present their basic algebraic properties.

1. If the matrices are of conformable sizes,  $[A \otimes (B + C)] = [A \otimes B] + [A \otimes C]$ .
2. If the matrices are of conformable sizes,  $[(A + B) \otimes C] = [A \otimes C] + [B \otimes C]$ .
3. If  $a$  and  $b$  are scalars,  $ab[A \otimes B] = [aA \otimes bB]$ .
4. If the matrices are of conformable sizes,  $[A \otimes B][C \otimes D] = [AC \otimes BD]$ .
5. The transpose of a Kronecker product matrix is  $[A \otimes B]' = [A' \otimes B']$ .
6. The generalized inverse of a Kronecker product matrix is  $[A \otimes B]^- = [A^- \otimes B^-]$ .
7. For two vectors  $v$  and  $w$ ,  $\text{Vec}(vw')$  =  $w \otimes v$ .
8. For a matrix  $W$  and conformable matrices  $A$  and  $B$ ,  $\text{Vec}(AWB')$  =  $[B \otimes A]\text{Vec}(W)$ .
9. For conformable matrices  $A$  and  $B$ ,  $\text{Vec}(A)'\text{Vec}(B)$  =  $\text{tr}(A'B)$ .
10. The Vec operator commutes with any matrix operation that is performed elementwise. For example,  $E\{\text{Vec}(W)\}$  =  $\text{Vec}\{E(W)\}$  when  $W$  is a random matrix. Similarly, for conformable matrices  $A$  and  $B$  and scalar  $\phi$ ,  $\text{Vec}(A + B)$  =  $\text{Vec}(A) + \text{Vec}(B)$  and  $\text{Vec}(\phi A)$  =  $\phi \text{Vec}(A)$ .
11. If  $A$  and  $B$  are positive definite, then  $A \otimes B$  is positive definite.

Most of these are well-known facts and easy to establish. Two of them are somewhat more unusual. Proofs for Items 8 and 11 are given in *ALM-III*, Appendix A.2 and in earlier editions of this book.

## B.6 Tensors

Tensors are simply an alternative notation for writing vectors. This notation has substantial advantages when dealing with quadratic forms and when dealing with more general concepts than quadratic forms. Our main purpose in discussing them here is simply to illustrate how flexibly subscripts can be used in writing vectors.

Consider a vector  $Y = (y_1, \dots, y_n)'$ . The tensor notation for this is simply  $y_i$ . We can write another vector  $a = (a_1, \dots, a_n)'$  as  $a_i$ . When written individually, the subscript is not important. In other words,  $a_i$  is the same vector as  $a_j$ . Note that the length of these vectors needs to be understood from the context. Just as when we write  $Y$  and  $a$  in conventional vector notation, there is nothing in the notation  $y_i$  or  $a_i$  to tell us how many elements are in the vector.

If we want the inner product  $a'Y$ , in tensor notation we write  $a_i y_i$ . Here we are using something called the *summation convention*. Because the subscripts on  $a_i$  and  $y_i$  are the same,  $a_i y_i$  is taken to mean  $\sum_{i=1}^n a_i y_i$ . If, on the other hand, we wrote  $a_i y_j$ , this means something completely different.  $a_i y_j$  is an alternative notation for the Kronecker product  $[a \otimes Y] = (a_1 y_1, \dots, a_1 y_n, a_2 y_1, \dots, a_n y_n)'$ . In  $[a \otimes Y] \equiv a_i y_j$ , we have two subscripts identifying the rows of the vector.

Now, suppose we want to look at a quadratic form  $Y'AY$ , where  $Y$  is an  $n$  vector and  $A$  is  $n \times n$ . One way to rewrite this is

$$Y'AY = \sum_{i=1}^n \sum_{j=1}^n y_i a_{ij} y_j = \sum_{i=1}^n \sum_{j=1}^n a_{ij} y_i y_j = \text{Vec}(A)'[Y \otimes Y].$$

(From Property B.5.8 we also have  $Y'AY = [Y' \otimes Y']\text{Vec}(A)$ .) Here we have rewritten the quadratic form as a linear combination of the elements in the vector  $[Y \otimes Y]$ . The linear combination is determined by the elements of the vector  $\text{Vec}(A)$ . In tensor notation, this becomes quite simple. Using the summation convention in which objects with the same subscript are summed over,

$$Y'AY = y_i a_{ij} y_j = a_{ij} y_i y_j.$$

The second term just has the summation signs removed, but the third term, which obviously gives the same sum as the second, is actually the tensor notation for  $\text{Vec}(A)'[Y \otimes Y]$ . Again,  $\text{Vec}(A) = (a_{11}, a_{21}, a_{31}, \dots, a_{nn})'$  uses two subscripts to identify rows of the vector. Obviously, if you had a need to consider things like

$$\sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n a_{ijk} y_i y_j y_k \equiv a_{ijk} y_i y_j y_k,$$

the tensor version  $a_{ijk} y_i y_j y_k$  saves some work.

There is one slight complication in how we have been writing things. Suppose  $A$  is not symmetric and we have another  $n$  vector  $W$ . Then we might want to consider

$$W'AY = \sum_{i=1}^n \sum_{j=1}^n w_i a_{ij} y_j.$$

From item 8 in the previous subsection,

$$W'AY = \text{Vec}(W'AY) = [Y' \otimes W'] \text{Vec}(A).$$

Alternatively,

$$W'AY = \sum_{i=1}^n \sum_{j=1}^n w_i a_{ij} y_j = \sum_{i=1}^n \sum_{j=1}^n a_{ij} y_j w_i = \text{Vec}(A)' [Y \otimes W]$$

or  $W'AY = Y'A'W = \text{Vec}(A)' [W \otimes Y]$ . However, with  $A$  nonsymmetric,  $W'A'Y = \text{Vec}(A)' [Y \otimes W]$  is typically different from  $W'AY$ . The Kronecker notation requires that care be taken in specifying the order of the vectors in the Kronecker product, and whether or not to transpose  $A$  before using the  $\text{Vec}$  operator. In tensor notation,  $W'AY$  is simply  $w_i a_{ij} y_j$ . In fact, the orders of the vectors can be permuted in any way; so, for example,  $a_{ij} y_j w_i$  means the same thing.  $W'A'Y$  is simply  $w_i a_{ji} y_j$ . The tensor notation and the matrix notation require less effort than the Kronecker notation.

For our purposes, the real moral here is simply that the subscripting of an individual vector does not matter. We can write a vector  $Y = (y_1, \dots, y_n)'$  as  $Y = [y_k]$  (in tensor notation as simply  $y_k$ ), or we can write the same  $n$  vector as  $Y = [y_{ij}]$  (in tensor notation, simply  $y_{ij}$ ), where, as long as we know the possible values that  $i$  and  $j$  can take on, the actual order in which we list the elements is not of much importance. Thus, if  $i = 1, \dots, t$  and  $j = 1, \dots, N_i$ , with  $n = \sum_{i=1}^t N_i$ , it really does not matter if we write a vector  $Y$  as  $(y_1, \dots, y_n)$ , or  $(y_{11}, \dots, y_{1N_1}, y_{21}, \dots, y_{tN_t})'$  or  $(y_{t1}, \dots, y_{tN_t}, y_{t-1,1}, \dots, y_{1N_1})'$  or in any other fashion we may choose, as long as we keep straight which row of the vector is which. Thus, a linear combination  $a'Y$  can be written  $\sum_{k=1}^n a_k y_k$  or  $\sum_{i=1}^t \sum_{j=1}^{N_i} a_{ij} y_{ij}$ . In tensor notation, the first of these is simply  $a_k y_k$  and the second is  $a_{ij} y_{ij}$ . These ideas become very handy in examining analysis of variance models, where the standard approach is to use multiple subscripts to identify the various observations. The subscripting has no intrinsic importance; the only thing that matters is knowing which row is which in the vectors. The subscripts are an aid in this identification, but they do not create any problems. We can still put all of the observations into a vector and use standard operations on them.

## B.7 Exercises

### Exercise B.0

(a) Let  $X = [X_0, X_1]$  with  $M$  and  $M_0$  the ppos onto  $C(X)$  and  $C(X_0)$ , respectively. Show that  $(I - M_0)X_1[X'_0(I - M_0)X_1]^{-1}X'_1(I - M_0)$  is the ppo onto  $C(X_0)^\perp_{C(X)}$ .

(b) Let  $r$  and  $s$  be two  $n$  vectors. Let  $M_r$  be the ppo onto  $C(r)$ , then  $s'(I - M_r)s \geq 0$ . Use this fact to prove the Cauchy–Schwarz inequality,

$$(s'r)^2 \leq s's r'r.$$

### Exercise B.1

(a) Show that

$$A^k x + b_{k-1}A^{k-1}x + \cdots + b_0x = (A - \mu I) \left( A^{k-1}x + \tau_{k-2}A^{k-2}x + \cdots + \tau_0x \right) = 0,$$

where  $\mu$  is any nonzero solution of  $b_0 + b_1w + \cdots + b_kw^k = 0$  with  $b_k = 1$  and  $\tau_j = -(b_0 + b_1\mu + \cdots + b_j\mu^j)/\mu^{j+1}$ ,  $j = 0, \dots, k$ .

(b) Show that if the only root of  $b_0 + b_1w + \cdots + b_kw^k$  is zero, then the factorization in (a) still holds.

(c) The solution  $\mu$  used in (a) need not be a real number, in which case  $\mu$  is a complex eigenvalue and the  $\tau_i$ s are complex; so the eigenvector is complex. Show that with  $A$  symmetric,  $\mu$  must be real because the eigenvalues of  $A$  must be real. In particular, assume that

$$A(y + iz) = (\lambda + i\gamma)(y + iz),$$

for  $y, z, \lambda$ , and  $\gamma$  real vectors and scalars, respectively, set  $Ay = \lambda y - \gamma z$ ,  $Az = \gamma y + \lambda z$ , and examine  $z'Ay = y'Az$ .

**Exercise B.2** Prove Proposition B.32.

**Exercise B.3** Show that any nonzero symmetric matrix  $A$  can be written as  $A = PDP'$ , where  $C(A) = C(P)$ ,  $P'P = I$ , and  $D$  is nonsingular.

**Exercise B.4** Prove Corollary B.52.

**Exercise B.5** Prove Corollary B.53.

**Exercise B.6** Show  $\text{tr}(cI_n) = nc$ .

**Exercise B.7** Let  $a, b, c,$  and  $d$  be real numbers. If  $ad - bc \neq 0$ , find the inverse of

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix}.$$

**Exercise B.8** Prove Theorem B.28, i.e., let  $A$  be an  $r \times s$  matrix, let  $B$  be an  $s \times r$  matrix, and show that  $\text{tr}(AB) = \text{tr}(BA)$ .

**Exercise B.9** Determine whether the matrices given below are positive definite, nonnegative definite, or neither.

$$\begin{bmatrix} 3 & 2 & -2 \\ 2 & 2 & -2 \\ -2 & -2 & 10 \end{bmatrix}, \quad \begin{bmatrix} 26 & -2 & -7 \\ -2 & 4 & -6 \\ -7 & -6 & 13 \end{bmatrix}, \quad \begin{bmatrix} 26 & 2 & 13 \\ 2 & 4 & 6 \\ 13 & 6 & 13 \end{bmatrix}, \quad \begin{bmatrix} 3 & 2 & -2 \\ 2 & -2 & -2 \\ -2 & -2 & 10 \end{bmatrix}.$$

**Exercise B.10** Show that the matrix  $B$  given below is positive definite, and find a matrix  $Q$  such that  $B = QQ'$ . (Hint: The first row of  $Q$  can be taken as  $(1, -1, 0)$ .)

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

**Exercise B.11** Let

$$A = \begin{bmatrix} 2 & 0 & 4 \\ 1 & 5 & 7 \\ 1 & -5 & -3 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}, \quad C = \begin{bmatrix} 1 & 4 & 1 \\ 2 & 5 & 1 \\ -3 & 0 & 1 \end{bmatrix}.$$

Use Theorem B.35 to find the perpendicular projection operator onto the column space of each matrix.

**Exercise B.12** Show that for a perpendicular projection matrix  $M$ ,

$$\sum_i \sum_j m_{ij}^2 = r(M).$$

**Exercise B.13** Prove that if  $M = M'M$ , then  $M = M'$  and  $M = M^2$ .

**Exercise B.14** Let  $M_1$  and  $M_2$  be perpendicular projection matrices, and let  $M_0$  be a perpendicular projection operator onto  $C(M_1) \cap C(M_2)$ . Show that the following are equivalent:

(a)  $M_1M_2 = M_2M_1$ .

(b)  $M_1M_2 = M_0$ .

(c)  $\{C(M_1) \cap [C(M_1) \cap C(M_2)]^\perp\} \perp \{C(M_2) \cap [C(M_1) \cap C(M_2)]^\perp\}$ .

Hints: (i) Show that  $M_1M_2$  is a projection operator. (ii) Show that  $M_1M_2$  is symmetric. (iii) Note that  $C(M_1) \cap [C(M_1) \cap C(M_2)]^\perp = C(M_1 - M_0)$ .

**Exercise B.15** Let  $M_1$  and  $M_2$  be perpendicular projection matrices. Show that

(a) the eigenvalues of  $M_1M_2$  are no greater than 1 in absolute value (they may be complex);

(b)  $\text{tr}(M_1M_2) \leq r(M_1M_2)$ .

Hints: For part (a) show that with  $x'Mx \equiv \|Mx\|^2$ ,  $\|Mx\| \leq \|x\|$  for any perpendicular projection operator  $M$ . Use this to show that if  $M_1M_2x = \lambda x$ , then  $\|M_1M_2x\| \geq |\lambda| \|M_1M_2x\|$ .

**Exercise B.16** For vectors  $x$  and  $y$ , let  $M_x = x(x'x)^{-1}x'$  and  $M_y = y(y'y)^{-1}y'$ . Show that  $M_xM_y = M_yM_x$  if and only if  $C(x) = C(y)$  or  $x \perp y$ .

**Exercise B.17** Consider the matrix

$$A = \begin{bmatrix} 0 & 1 \\ 0 & 1 \end{bmatrix}.$$

(a) Show that  $A$  is a projection matrix.

(b) Is  $A$  a perpendicular projection matrix? Why or why not?

(c) Describe the space that  $A$  projects onto and the space that  $A$  projects along. Sketch these spaces.

(d) Find another projection operator onto the space that  $A$  projects onto.

**Exercise B.18** Let  $A$  be an arbitrary projection matrix. Show that  $C(I - A) = C(A')^\perp$ .

Hints: Recall that  $C(A')^\perp$  is the null space of  $A$ . Show that  $(I - A)$  is a projection matrix.

**Exercise B.19** Show that if  $A^-$  is a generalized inverse of  $A$ , then so is

$$G = A^-AA^- + (I - A^-A)B_1 + B_2(I - AA^-)$$

for any choices of  $B_1$  and  $B_2$  with conformable dimensions.

**Exercise B.20** Let  $A$  be positive definite with eigenvalues  $\lambda_1, \dots, \lambda_n$ . Show that  $A^{-1}$  has eigenvalues  $1/\lambda_1, \dots, 1/\lambda_n$  and the same eigenvectors as  $A$ .



**Exercise B.21** For  $A$  nonsingular, let

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix},$$

and let  $A_{1.2} = A_{11} - A_{12}A_{22}^{-1}A_{21}$ . Show that if all inverses exist,

$$A^{-1} = \begin{bmatrix} A_{1.2}^{-1} & -A_{1.2}^{-1}A_{12}A_{22}^{-1} \\ -A_{22}^{-1}A_{21}A_{1.2}^{-1} & A_{22}^{-1} + A_{22}^{-1}A_{21}A_{1.2}^{-1}A_{12}A_{22}^{-1} \end{bmatrix}$$

and that

$$A_{22}^{-1} + A_{22}^{-1}A_{21}A_{1.2}^{-1}A_{12}A_{22}^{-1} = [A_{22} - A_{21}A_{11}^{-1}A_{12}]^{-1}.$$

**Exercise B.22** Show that if  $A$  is idempotent, then  $\mathcal{N}(A) = C(I - A)$ . Hint: Show that each set is contained in the other.

**Exercise B.23** Consider the vectors that are the columns of a matrix  $X_{n \times p}$  with  $r(X) = r$ . A *rotation* of these vectors keeps their lengths the same and the angles between them the same. In other words, it keeps all of the inner products between them the same. If  $P$  is an orthonormal matrix, then  $Z = PX$  is a rotation of the columns of  $X$  because  $Z'Z = X'P'PX = X'X$ .

Rotating vectors within a subspace is more difficult. Let the  $n \times r$  matrix  $Q$  have columns that are an orthonormal basis for  $C(X)$ . Write  $X = QB$ . Show that a rotation of the columns of  $X$  that remains within  $C(X)$  is obtained from  $Z = QPB$  where  $P$  is an  $r \times r$  orthonormal matrix.

**Exercise B.24** Show that  $r(X) = r(X')$ . Hints: Let  $X$  be  $n \times p$  with  $r(X') = r$ . Let  $\tilde{X}$  be an  $r \times p$  matrix with the rows of  $\tilde{X}$  forming a basis for  $C(X')$ . Write  $X = B\tilde{X}$  and argue that  $r(X) \leq r = r(X')$ . Reverse the roles of  $X$  and  $X'$ .

## Appendix C

### Some Univariate Distributions

**Abstract** The tests and confidence intervals presented in this book rely almost exclusively on the  $\chi^2$ ,  $t$ , and  $F$  distributions. This appendix defines each of the distributions.

**Definition C.1** Let  $Z_1, \dots, Z_n$  be independent with  $Z_i \sim N(\mu_i, 1)$ . Then if

$$W \sim \sum_{i=1}^n Z_i^2$$

we say  $W$  has a *noncentral chi-squared distribution* with  $n$  degrees of freedom and *noncentrality parameter*  $\gamma = \sum_{i=1}^n \mu_i^2/2$ . Write  $W \sim \chi^2(n, \gamma)$ .

See Rao (1973, Section 3b.2) for a proof that the distribution of  $W$  depends only on  $n$  and  $\gamma$ .

It is evident from the definition that if  $X \sim \chi^2(r, \gamma)$  and  $Y \sim \chi^2(s, \delta)$  with  $X$  and  $Y$  independent, then  $(X + Y) \sim \chi^2(r + s, \gamma + \delta)$ . A central  $\chi^2$  distribution is a distribution with a noncentrality parameter of zero, i.e.,  $\chi^2(r, 0)$ . We will use  $\chi^2(r)$  to denote a  $\chi^2(r, 0)$  distribution. The  $100\alpha$ th percentile of a  $\chi^2(r)$  distribution is the point  $\chi^2(\alpha, r)$  that satisfies the equation

$$\Pr[\chi^2(r) \leq \chi^2(\alpha, r)] = \alpha.$$

Note that if  $0 \leq a < 1$ , the  $100a$  percentile of a central  $\chi^2(b)$  is denoted  $\chi^2(a, b)$ . However, if  $a$  is a positive integer,  $\chi^2(a, b)$  denotes a noncentral chi-squared distribution.

**Definition C.2** Let  $X \sim N(\mu, 1)$  and  $Y \sim \chi^2(n)$  with  $X$  and  $Y$  independent. Then

$$W = \frac{X}{\sqrt{Y/n}}$$

has a *noncentral  $t$  distribution* with  $n$  degrees of freedom and noncentrality parameter  $\mu$ . Write  $W \sim t(n, \mu)$ . If  $\mu = 0$ , we say that the distribution is a central  $t$  distribution and write  $W \sim t(n)$ . The  $100\alpha$ th percentile of a  $t(n)$  distribution is denoted  $t(\alpha, n)$ .

**Definition C.3** Let  $X \sim \chi^2(r, \gamma)$  and  $Y \sim \chi^2(s, 0)$  with  $X$  and  $Y$  independent. Then

$$W = \frac{X/r}{Y/s}$$

has a *noncentral  $F$  distribution* with  $r$  numerator and  $s$  denominator degrees of freedom and noncentrality parameter  $\gamma$ . Write  $W \sim F(r, s, \gamma)$ . If  $\gamma = 0$ , write  $W \sim F(r, s)$  for the central  $F$  distribution. The  $100\alpha$ th percentile of  $F(r, s)$  is denoted  $F(\alpha, r, s)$ .

As indicated, if the noncentrality parameter of any of these distributions is zero, the distribution is referred to as a *central distribution* (e.g., central  $F$  distribution). The central distributions are those commonly used in statistical methods courses. If any of these distributions is not specifically identified as a noncentral distribution, it should be assumed to be a central distribution.

It is easily seen from Definition C.1 that any noncentral chi-squared distribution *tends* to be larger than the central chi-squared distribution with the same number of degrees of freedom. Similarly, from Definition C.3, a noncentral  $F$  tends to be larger than the corresponding central  $F$  distribution. (These ideas are made rigorous in Exercise C.1.) The fact that the noncentral  $F$  distribution tends to be larger than the corresponding central  $F$  distribution is the basis for many of the tests used in linear models. Typically, test statistics are used that have a central  $F$  distribution if the reduced (null) model is true and a noncentral  $F$  distribution if the full model is true but the null model is not. Since the noncentral  $F$  distribution tends to be larger, large values of the test statistic are more consistent with the full model than with the null. Thus, the form of an appropriate rejection region when the full model is true is to reject the null hypothesis for large values of the test statistic.

The power of these  $F$  tests is simply a function of the noncentrality parameter. Given a value for the noncentrality parameter, there is no theoretical difficulty in finding the power of an  $F$  test. The power simply involves computing the probability of the rejection region when the probability distribution is a noncentral  $F$ . Davies (1980) gives an algorithm for making these and more general computations.

We now prove a theorem about central  $F$  distributions that will be useful in Chapter 5.

**Theorem C.4** If  $s > t$ , then  $sF(1 - \alpha, s, v) \geq tF(1 - \alpha, t, v)$ .

*Proof* Let  $X \sim \chi^2(s)$ ,  $Y \sim \chi^2(t)$ , and  $Z \sim \chi^2(v)$ . Let  $Z$  be independent of  $X$  and  $Y$ . Note that  $(X/s)/(Z/v)$  has an  $F(s, v)$  distribution; so  $sF(1 - \alpha, s, v)$  is the  $100(1 - \alpha)$  percentile of the distribution of  $X/(Z/v)$ . Similarly,  $tF(1 - \alpha, t, v)$  is the  $100(1 - \alpha)$  percentile of the distribution of  $Y/(Z/v)$ .

We will first argue that to prove the theorem it is enough to show that

$$\Pr[X \leq d] \leq \Pr[Y \leq d] \quad (1)$$

for all real numbers  $d$ . We will then show that (1) is true.

If (1) is true, if  $c$  is any real number, and if  $Z = z$ , by independence we have

$$\Pr[X \leq cz/v] = \Pr[X \leq cz/v | Z = z] \leq \Pr[Y \leq cz/v | Z = z] = \Pr[Y \leq cz/v].$$

Taking expectations with respect to  $Z$ ,

$$\begin{aligned} \Pr[X/(Z/v) \leq c] &= E(\Pr[X \leq cz/v | Z = z]) \\ &\leq E(\Pr[Y \leq cz/v | Z = z]) \\ &= \Pr[Y/(Z/v) \leq c]. \end{aligned}$$

Since the cumulative distribution function (cdf) for  $X/(Z/v)$  is always no greater than the cdf for  $Y/(Z/v)$ , the point at which a probability of  $1 - \alpha$  is attained for  $X/(Z/v)$  must be no less than the similar point for  $Y/(Z/v)$ . Therefore,

$$sF(1 - \alpha, s, v) \geq tF(1 - \alpha, t, v).$$

To see that (1) holds, let  $Q$  be independent of  $Y$  and  $Q \sim \chi^2(s - t)$ . Then, because  $Q$  is nonnegative,

$$\Pr[X \leq d] = \Pr[Y + Q \leq d] \leq \Pr[Y \leq d]. \quad \square$$

## Exercise

**Definition C.5** Consider two random variables  $W_1$  and  $W_2$ .  $W_2$  is said to be *stochastically larger* than  $W_1$  if for every real number  $w$

$$\Pr[W_1 > w] \leq \Pr[W_2 > w].$$

If for some random variables  $W_1$  and  $W_2$ ,  $W_2$  is stochastically larger than  $W_1$ , then we also say that the distribution of  $W_2$  is stochastically larger than the distribution of  $W_1$ .

**Exercise C.1** Show that a noncentral chi-squared distribution is stochastically larger than the central chi-squared distribution with the same degrees of freedom. Show that a noncentral  $F$  distribution is stochastically larger than the corresponding central  $F$  distribution.

## Appendix D

# Multivariate Distributions

**Abstract** This appendix reviews properties of multivariate distributions. It also examines the concept of identifiable parameters.

Let  $(x_1, \dots, x_n)'$  be a random vector. The joint cumulative distribution function (cdf) of  $(x_1, \dots, x_n)'$  is

$$F(u_1, \dots, u_n) \equiv \Pr [x_1 \leq u_1, \dots, x_n \leq u_n].$$

If  $F(u_1, \dots, u_n)$  is the cdf of a discrete random variable, we can define a (joint) probability mass function

$$f(u_1, \dots, u_n) \equiv \Pr [x_1 = u_1, \dots, x_n = u_n].$$

If  $F(u_1, \dots, u_n)$  admits the  $n$ th order mixed partial derivative, then we can define a (joint) density function

$$f(u_1, \dots, u_n) \equiv \frac{\partial^n}{\partial u_1 \cdots \partial u_n} F(u_1, \dots, u_n).$$

The cdf can be recovered from the density as

$$F(u_1, \dots, u_n) = \int_{-\infty}^{u_1} \cdots \int_{-\infty}^{u_n} f(w_1, \dots, w_n) dw_1 \cdots dw_n.$$

For a function  $g(\cdot)$  of  $(x_1, \dots, x_n)'$  into  $\mathbf{R}$ , the expected value is defined as

$$E[g(x_1, \dots, x_n)] = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} g(u_1, \dots, u_n) f(u_1, \dots, u_n) du_1 \cdots du_n.$$

We might also write this as  $E_x [g(x)]$ .

We now consider relationships between two random vectors, say  $x = (x_1, \dots, x_n)'$  and  $y = (y_1, \dots, y_m)'$ . Assume that the joint vector  $(x', y')' = (x_1, \dots, x_n, y_1, \dots, y_m)'$  has a density function

$$f_{x,y}(u, v) \equiv f_{x,y}(u_1, \dots, u_n, v_1, \dots, v_m).$$

Similar definitions and results hold if  $(x', y')'$  has a probability mass function.

The distribution of one random vector, say  $x$ , ignoring the other vector,  $y$ , is called the *marginal distribution* of  $x$ . The marginal cdf of  $x$  can be obtained by substituting the value  $+\infty$  into the joint cdf for all of the  $y$  variables:

$$F_x(u) = F_{x,y}(u_1, \dots, u_n, +\infty, \dots, +\infty).$$

The marginal density can be obtained either by partial differentiation of  $F_x(u)$  or by integrating the joint density over the  $y$  variables:

$$f_x(u) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f_{x,y}(u_1, \dots, u_n, v_1, \dots, v_m) dv_1 \cdots dv_m.$$

The conditional density of a vector, say  $x$ , given the value of the other vector, say  $y = v$ , is obtained by dividing the density of  $(x', y')'$  by the density of  $y$  evaluated at  $v$ , i.e.,

$$f_{x|y}(u|v) \equiv f_{x,y}(u, v) / f_y(v).$$

The conditional density is a well-defined density, so expectations with respect to it are well defined. Let  $g$  be a function from  $\mathbf{R}^n$  into  $\mathbf{R}$ ,

$$E[g(x)|y = v] = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} g(u) f_{x|y}(u|v) du,$$

where  $du \equiv du_1 du_2 \cdots du_n$ . Sometimes we write

$$E_{x|y=v}[g(x)] \equiv E[g(x)|y = v].$$

The standard properties of expectations hold for conditional expectations. For example, with  $a$  and  $b$  real,

$$E[ag_1(x) + bg_2(x)|y = v] = aE[g_1(x)|y = v] + bE[g_2(x)|y = v].$$

The conditional expectation of  $E[g(x)|y = v]$  is a function of the value  $v$ . Since  $y$  is random, we can consider  $E[g(x)|y = v]$  as a random variable. In this context we write  $E[g(x)|y]$  or  $E_{x|y}[g(x)]$ . An important property of conditional expectations is

$$E[g(x)] = E[E[g(x)|y]].$$

To see this, note that  $f_{x|y}(u|v) f_y(v) = f_{x,y}(u, v)$  and

$$\begin{aligned}
 E[E[g(x)|y]] &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} E[g(x)|y = v] f_y(v) dv \\
 &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \left[ \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} g(u) f_{x|y}(u|v) du \right] f_y(v) dv \\
 &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} g(u) f_{x|y}(u|v) f_y(v) du dv \\
 &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} g(u) f_{x,y}(u, v) du dv \\
 &= E[g(x)].
 \end{aligned}$$

In fact, both the notion of conditional expectation and this result can be generalized. Consider a function  $g(x, y)$  from  $\mathbf{R}^{n+m}$  into  $\mathbf{R}$ . If  $y = v$ , we can define  $E[g(x, y)|y = v]$  in a natural manner. If we consider  $y$  as random, we write  $E[g(x, y)|y]$ . It can be easily shown that

$$E[g(x, y)] = E[E[g(x, y)|y]].$$

A function of  $x$  or  $y$  alone can also be considered as a function from  $\mathbf{R}^{n+m}$  into  $\mathbf{R}$ .

A second important property of conditional expectations is that if  $h(y)$  is a function from  $\mathbf{R}^m$  into  $\mathbf{R}$ , we have

$$E[h(y)g(x, y)|y] = h(y)E[g(x, y)|y]. \tag{1}$$

This follows because if  $y = v$ ,

$$\begin{aligned}
 E[h(y)g(x, y)|y = v] &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} h(v)g(u, v) f_{x|y}(u|v) du \\
 &= h(v) \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} g(u, v) f_{x|y}(u|v) du \\
 &= h(v)E[g(x, y)|y = v].
 \end{aligned}$$

This is true for all  $v$ , so (1) holds. In particular, if  $g(x, y) \equiv 1$ , we get

$$E[h(y)|y] = h(y).$$

Finally, we can extend the idea of conditional expectation to a function  $g(x, y)$  from  $\mathbf{R}^{n+m}$  into  $\mathbf{R}^s$ . Write  $g(x, y) = [g_1(x, y), \dots, g_s(x, y)]'$ . Then define

$$E[g(x, y)|y] = (E[g_1(x, y)|y], \dots, E[g_s(x, y)|y])'$$

If their densities exist, two random vectors are *independent* if and only if their joint density is equal to the product of their marginal densities, i.e.,  $x$  and  $y$  are independent if and only if



$$f_{x,y}(u, v) = f_x(u)f_y(v).$$

Note that if  $x$  and  $y$  are independent,

$$f_{x|y}(u|v) = f_x(u).$$

If the random vectors  $x$  and  $y$  are independent, then any (reasonable) vector-valued functions of them, say  $g(x)$  and  $h(y)$ , are also independent. This follows easily from a more general definition of the independence of two random vectors: The random vectors  $x$  and  $y$  are independent if for any two (reasonable) sets  $A$  and  $B$ ,

$$\Pr[x \in A, y \in B] = \Pr[x \in A]\Pr[y \in B].$$

To prove that functions of random variables are independent, recall that the set inverse of a function  $g(u)$  on a set  $A_0$  is  $g^{-1}(A_0) \equiv \{u | g(u) \in A_0\}$ . That  $g(x)$  and  $h(y)$  are independent follows from the fact that for any (reasonable) sets  $A_0$  and  $B_0$ ,

$$\begin{aligned} \Pr[g(x) \in A_0, h(y) \in B_0] &= \Pr[x \in g^{-1}(A_0), y \in h^{-1}(B_0)] \\ &= \Pr[x \in g^{-1}(A_0)]\Pr[y \in h^{-1}(B_0)] \\ &= \Pr[g(x) \in A_0]\Pr[h(y) \in B_0]. \end{aligned}$$

By “reasonable” I mean things that satisfy the mathematical definitions of being measurable.

The *characteristic function* of a random vector  $x = (x_1, \dots, x_n)'$  is a function from  $\mathbf{R}^n$  to  $\mathbf{C}$ , the complex numbers. It is defined by

$$\varphi_x(t_1, \dots, t_n) = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \exp\left[i \sum_{j=1}^n t_j u_j\right] f_x(u_1, \dots, u_n) du_1 \cdots du_n.$$

We are interested in characteristic functions because if  $x = (x_1, \dots, x_n)'$  and  $y = (y_1, \dots, y_n)'$  are random vectors and if

$$\varphi_x(t_1, \dots, t_n) = \varphi_y(t_1, \dots, t_n)$$

for all  $(t_1, \dots, t_n)$ , then  $x$  and  $y$  have the same distribution.

For convenience, we have assumed the existence of densities. With minor modifications, the definitions and results of this appendix hold for any probability defined on  $\mathbf{R}^n$ .

## D.1 Identifiability

For better or worse (usually worse) much of statistical practice focuses on estimating and testing parameters. Identifiability is a property that ensures that this process is a sensible one.

Consider a collection of probability distributions  $Y \sim P_\theta, \theta \in \Theta$ . The parameter  $\theta$  merely provides the name (index) for each distribution in the collection. Identifiability ensures that each distribution has a unique name/index.

**Definition D.1** The parameterization  $\theta \in \Theta$  is *identifiable* if  $Y_1 \sim P_{\theta_1}, Y_2 \sim P_{\theta_2}$ , and  $Y_1 \sim Y_2$  imply that  $\theta_1 = \theta_2$ .

Being identifiable is easily confused with the concept of being well defined.

**Definition D.2** The parameterization  $\theta \in \Theta$  is *well defined* if  $Y_1 \sim P_{\theta_1}, Y_2 \sim P_{\theta_2}$ , and  $\theta_1 = \theta_2$  imply that  $Y_1 \sim Y_2$ .

The problem with not being identifiable is that some distributions have more than one name. Observed data give you information about the correct distribution and thus about the correct name. Typically, the more data you have, the more information you have about the correct name. Estimation is about getting close to the correct name and testing hypotheses is about deciding which of two lists contains the correct name. If a distribution has more than one name, it could be in both lists. (Significance testing is about whether it seems plausible that a name is on a list, so identifiability seems less of an issue.) If a distribution has more than one name, does getting close to one of those names really help? In applications to linear models, typically distributions have only one name or they have an infinite number of names.

The ideas are roughly this. If the distributions are well defined and I know that Wesley O. Johnson ( $\theta_1$ ) and O. Wesley Johnson ( $\theta_2$ ) are the same person ( $\theta_1 = \theta_2$ ), then, say, any collection of blood pressure readings on Wesley O. should look pretty much the same as comparable readings on O. Wesley. They would be two samples from the same distribution. Identifiability is the following: if all the samples I have taken or ever could take on Wesley O. look pretty much the same as samples on O. Wesley, then Wesley O. would have to be the same person as O. Wesley. (The reader might consider whether personhood is actually an identifiable parameter for blood pressure.)

For the multivariate normal distributions of Section 1.2, being well defined is the requirement that if  $Y_1 \sim N(\mu_1, V_1), Y_2 \sim N(\mu_2, V_2)$ , and  $\mu_1 = \mu_2$  and  $V_1 = V_2$ , then  $Y_1 \sim Y_2$ . Theorem 1.2.2 establishes that the mean and covariance of a multivariate normal determine the distribution. Being identifiable is that if  $Y_1 \sim N(\mu_1, V_1), Y_2 \sim N(\mu_2, V_2)$ , and  $Y_1 \sim Y_2$ , then  $\mu_1 = \mu_2$  and  $V_1 = V_2$ . Obviously, two random vectors with the same distribution have to have the same mean vector and covariance matrix. But life gets more complicated.

The more interesting problem for multivariate normality is a model

$$Y \sim N [F(\beta), V(\phi)]$$

where  $F$  and  $V$  are known functions of parameter vectors  $\beta$  and  $\phi$ . To show that  $\beta$  and  $\phi$  are identifiable we need to consider

$$Y_1 \sim N [F(\beta_1), V(\phi_1)], \quad Y_2 \sim N [F(\beta_2), V(\phi_2)]$$

and show that if  $Y_1 \sim Y_2$  then  $\beta_1 = \beta_2$  and  $\phi_1 = \phi_2$ . From our earlier discussion, if  $Y_1 \sim Y_2$  then  $F(\beta_1) = F(\beta_2)$  and  $V(\phi_1) = V(\phi_2)$ . We need to check that  $F(\beta_1) = F(\beta_2)$  implies  $\beta_1 = \beta_2$  and that  $V(\phi_1) = V(\phi_2)$  implies  $\phi_1 = \phi_2$ .

Section 2.1 gives an extensive discussion of when the mean parameterization is identifiable, i.e., when  $F(\beta_1) = F(\beta_2)$  implies  $\beta_1 = \beta_2$ . There we defined identifiable functions of  $\beta$  as those that are functions of  $F(\beta)$ .

In this book, we mostly consider simple models for the covariance parameterization  $V(\phi)$ ; models that are clearly identifiable because they involve at most one scalar parameter. We consider  $V(\phi) \equiv \sigma^2 I$ , and  $V(\phi) \equiv \sigma^2 V$  where  $V$  is a known nonnegative definite matrix, and, again with known  $V$ ,  $V(\phi) \equiv V$ , which involves no parameterization. For example, if  $V(\phi_1) \equiv \sigma_1^2 I = V(\phi_2) \equiv \sigma_2^2 I$ , we must have  $\phi_1 \equiv \sigma_1^2 = \phi_2 \equiv \sigma_2^2$ . As long as  $V$  is not the zero matrix, the covariance parameterizations in this book are identifiable.

*ALM-III* examines many commonly used models for  $V(\phi)$  using similar notation to that used here. In particular, linear covariance parameterizations of the form  $\sum_{r=0}^s \phi_r V_r$  for nonnegative  $\phi_r$ s and nonnegative definite known  $V_r$ s are identifiable if and only if the  $V_r$ s are linearly independent, i.e., the  $\text{Vec}(V_r)$ s are linearly independent. The covariance matrices of Chapter 11 fall into this category.

## Exercise

**Exercise D.1** Let  $x$  and  $y$  be independent. Show that

- (a)  $E[g(x)|y] = E[g(x)]$ ;
- (b)  $E[g(x)h(y)] = E[g(x)]E[h(y)]$ .

## Appendix E

# Inference for One Parameter

**Abstract** Since the third edition of this book, I have thought hard about the philosophy of testing as a basis for non-Bayesian statistical inference, cf. Christensen (2005, 2008). This appendix has been modified accordingly. The approach taken is one I call Fisherian, as opposed to the Neyman–Pearson approach. The theory presented here has no formal role for alternative hypotheses. A more extensive discussion of these ideas appears in Chapter 3 of Christensen (2015).

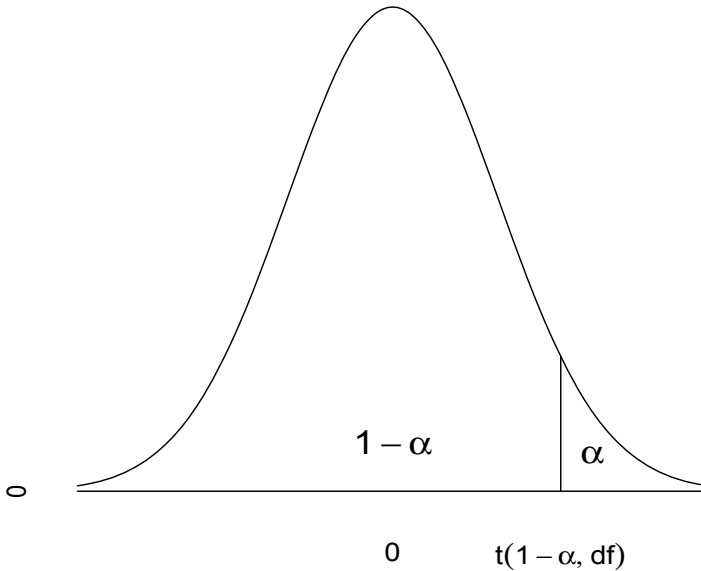
A significance testing problem is essentially a form of proof by contradiction. We have a *null model* for the data and we determine whether the observed data seem to contradict that null model or whether they are consistent with it. If the data contradict the null model, something must be wrong with the null model. Having data consistent with the null model certainly does not suggest that the null model is correct but may suggest that the model is tentatively adequate. The catch is that we rarely get an absolute contradiction to the null model, so we use probability to determine the extent to which the data seem inconsistent with the null model.

In the current discussion, *it is convenient to break the null model into two parts: a general model for the data and a particular statement about a single parameter of interest, called the null hypothesis ( $H_0$ ).*

Many statistical tests and confidence intervals for a single parameter are applications of the same theory. (Tests and confidence intervals for variances are an exception.) To use this theory we need to know four things: [1] The unobservable *parameter* of interest (*Par*). [2] The *estimate* of the parameter (*Est*). [3] The *standard error* of the estimate ( $SE(Est)$ ), wherein  $SE(Est)$  is typically an estimate of the standard deviation of *Est*, but if we happened to know the actual standard deviation, we would be happy to use it. And [4] an appropriate *reference distribution*. Specifically, we need the distribution of

$$\frac{Est - Par}{SE(Est)}.$$

If the  $SE(Est)$  is estimated, the reference distribution is usually the *t* distribution with



**Fig. E.1** Percentiles of  $t(df)$  distributions

some known number of degrees of freedom  $df$ , say,  $t(df)$ . If the  $SE(Est)$  is known, then the distribution is usually the standard normal distribution, i.e., a  $t(\infty)$ . In some problems (e.g., problems involving the binomial distribution) large sample results are used to get an approximate distribution and then the technique proceeds as if the approximate distribution were correct. When appealing to large sample results, the known distribution of part [4] is the standard normal (although I suspect that a  $t(df)$  distribution with a reasonable, finite number of degrees of freedom would give more realistic results).

*These four required items are derived from the model for the data* (although sometimes the standard error incorporates the null hypothesis). For convenience, we may refer to these four items as “the model.”

The  $1 - \alpha$  percentile of a distribution is the point that cuts off the top  $\alpha$  of the distribution. For a  $t$  distribution, denote this  $t(1 - \alpha, df)$  as seen in Figure E.1. Formally, we can write

$$\Pr \left[ \frac{Est - Par}{SE(Est)} \geq t(1 - \alpha, df) \right] = \alpha.$$

By symmetry about zero, we also have

$$\Pr \left[ \frac{Est - Par}{SE(Est)} \leq -t(1 - \alpha, df) \right] = \alpha.$$

To keep the discussion as simple as possible, numerical examples have been restricted to one-sample normal theory. However, the results also apply to inferences

on each individual mean and the difference between the means in two-sample problems, contrasts in analysis of variance, coefficients in regression, and, in general, to one-dimension estimable parametric functions in arbitrary linear models.

## E.1 Testing

We want to test the null hypothesis

$$H_0 : Par = m,$$

where  $m$  is some known number. In *significance (Fisherian) testing*, we cannot do that. *What we can do* is test the null model, which is the combination of the model and the null hypothesis. The test is based on the assumption that both the model and  $H_0$  are true. As mentioned earlier, it is rare that data contradict the null model absolutely, so we check to see if the data seem inconsistent with the null model.

What kind of data are inconsistent with the null model? Consider the *test statistic*

$$\frac{Est - m}{SE(Est)}.$$

With  $m$  known, the test statistic is an observable random variable. If the null model is true, the test statistic has a known  $t(df)$  distribution as illustrated in Figure E.1. The  $t(df)$  distribution is likely to give values near 0 and is increasingly less likely to give values far from 0. Therefore, weird data, i.e., those that are most inconsistent with the null model, are large positive and large negative values of  $[Est - m]/SE(Est)$ . The density (shape) of the  $t(df)$  distribution allows us to order the possible values of the test statistic in terms of how weird they are relative to the null model.

To decide on a formal test, we need to decide which values of the test statistic will cause us to reject the null model and which will not. In other words, “How weird must data be before we question the null model?” We solve this problem by picking a small probability  $\alpha$  that determines a *rejection region*, sometimes called a *critical region*. The rejection region consists of the weirdest test statistic values under the null model, but is restricted to have a probability of only  $\alpha$  under the null model. Since a  $t(df)$  distribution is symmetric about 0 and the density decreases as we go away from 0, the  $\alpha$  critical region consists of points less than  $-t(1 - \alpha/2, df)$  and points larger than  $t(1 - \alpha/2, df)$ . In other words, the  $\alpha$  level test for the model with  $H_0 : Par = m$  is to reject the null model if

$$\frac{Est - m}{SE(Est)} \geq t\left(1 - \frac{\alpha}{2}, df\right)$$

or if

$$\frac{Est - m}{SE(Est)} \leq -t\left(1 - \frac{\alpha}{2}, df\right).$$

This is equivalent to rejecting the null model if

$$\frac{|Est - m|}{SE(Est)} \geq t\left(1 - \frac{\alpha}{2}, df\right).$$

What causes us to reject the null model? Either having a true model that is so different from the null that the data look “weird,” or having the null model true and getting unlucky with the data.

Observing weird data, i.e., data that are inconsistent with the null model, gives us cause to question the validity of the null model. Specifying a small  $\alpha$  level merely ensures that everything in the rejection region really constitutes weird data. More properly, specifying a small  $\alpha$  level is our means of determining what constitutes weird data. Although  $\alpha$  can be viewed as a probability, it is better viewed as a measure of how weird the data must be relative to the null model before we will reject. We want  $\alpha$  small so that we only reject the null model for data that are truly weird, but we do not want  $\alpha$  so small that we fail to reject the null model even when very strange data occur.

Rejecting the null model means that *either* the null hypothesis *or* the model is deemed incorrect. Only if we are confident that the model is correct can we conclude that the null hypothesis is wrong. If we want to make conclusions about the null hypothesis, it is important to do everything possible to assure ourselves that the model is reasonable.

If we do not reject the null model, we merely have data that are consistent with the null model. That in no way implies that the null model is true. Many other models will also be consistent with the data. Typically,  $Par = m + 0.00001$  fits the data about as well as the null model. Not rejecting the test does not imply that the null model is true any more than rejecting the null model implies that the underlying model is true.

*Example E.1* Suppose that 16 independent observations are taken from a normal population. Test  $H_0 : \mu = 20$  with  $\alpha$  level 0.01. The observed values of  $\bar{y}$  and  $s^2$  were 19.78 and 0.25, respectively.

[1]  $Par = \mu,$

[2]  $Est = \bar{y},$

[3]  $SE(Est) = \sqrt{s^2/16}$ . In this case, the  $SE(Est)$  is estimated.

[4]  $[Est - Par]/SE(Est) = [\bar{y} - \mu]/\sqrt{s^2/16}$  has a  $t(15)$  distribution.

With  $m = 20$ , the  $\alpha = 0.01$  test is to reject the  $H_0$  model if

$$|\bar{y} - 20|/\sqrt{s/4} \geq 2.947 = t(0.995, 15).$$

Having  $\bar{y} = 19.78$  and  $s^2 = 0.25$ , we reject if

$$\frac{|19.78 - 20|}{\sqrt{.25/16}} \geq 2.947.$$

Since  $|19.78 - 20|/\sqrt{.25/16} = |-1.76|$  is less than 2.947, we do not reject the null model at the  $\alpha = 0.01$  level.

*Nobody actually does this!* Or at least, nobody should do it. Although this procedure provides a philosophical basis for our statistical inferences, there are two other procedures, both based on this, that give uniformly more information. This procedure requires us to specify the model, the null hypothesis parameter value  $m$ , and the  $\alpha$  level. For a fixed model and a fixed null parameter  $m$ ,  $P$  values are more informative because they allow us to report test results for all  $\alpha$  levels. Alternatively, for a fixed model and a fixed  $\alpha$  level, confidence intervals report the values of all parameters that are consistent with the model and the data. (Parameter values that are inconsistent with the model and the data are those that would be rejected, assuming the model is true.) We now discuss these other procedures.

## E.2 $P$ Values

*The  $P$  value of a test is the probability under the null model of seeing data as weird or weirder than we actually saw.* Weirdness is determined by the distribution of the test statistic. If the observed value of the test statistic from Section 1 is  $t_{obs}$ , then the  $P$  value is the probability of seeing data as far or farther from 0 than  $t_{obs}$ . In general, we do not know if  $t_{obs}$  will be positive or negative, but its distance from 0 is  $|t_{obs}|$ . The  $P$  value is the probability that a  $t(df)$  distribution is less than or equal to  $-|t_{obs}|$  or greater than or equal to  $|t_{obs}|$ .

In Example E.1, the value of the test statistic is  $-1.76$ . Since  $t(0.95, 15) = 1.75$ , the  $P$  value of the test is approximately (just smaller than) 0.10. An  $\alpha = 0.10$  test would use the  $t(0.95, 15)$  value.

It is not difficult to see that the  $P$  value is the  $\alpha$  level at which the test would just barely be rejected. So if  $P \leq \alpha$ , the null model is rejected, and if  $P > \alpha$ , the data are deemed consistent with the null model. Knowing the  $P$  value lets us do all  $\alpha$  level tests of the null model. In fact, historically and philosophically,  $P$  values come before  $\alpha$  level tests. Rather than noticing that the  $\alpha$  level test has this relationship with  $P$  values, it is more general to define the  $\alpha$  level test as rejecting precisely when  $P \leq \alpha$ . We can then observe that, for our setup, the  $\alpha$  level test has the form given in Section 1.

While an  $\alpha$  level constitutes a particular choice about how weird the data must be before we decide to reject the null model, the  $P$  value measures the evidence against the null hypothesis. The smaller the  $P$  value, the more evidence against the null model.



### E.3 Confidence Intervals

A  $(1 - \alpha)100\%$  confidence interval (CI) for  $Par$  is defined to be the set of all parameter values  $m$  that would not be rejected by an  $\alpha$  level test. In Section 1 we gave the rule for when an  $\alpha$  level test of  $H_0 : Par = m$  rejects. Conversely, the null model will not be rejected if

$$-t\left(1 - \frac{\alpha}{2}, df\right) < \frac{Est - m}{SE(Est)} < t\left(1 - \frac{\alpha}{2}, df\right). \quad (1)$$

Some algebra, given later, establishes that we do not reject the null model if and only if

$$Est - t\left(1 - \frac{\alpha}{2}, df\right) SE(Est) < m < Est + t\left(1 - \frac{\alpha}{2}, df\right) SE(Est). \quad (2)$$

This interval consists of all the parameter values  $m$  that are consistent with the data and the model as determined by an  $\alpha$  level test. The endpoints of the CI can be written

$$Est \pm t\left(1 - \frac{\alpha}{2}, df\right) SE(Est).$$

On occasion (as with binomial data), when doing an  $\alpha$  level test or a  $P$  value, we may let the standard error depend on the null hypothesis. To obtain a confidence interval using this approach, we need a standard error that does not depend on  $m$ .

*Example E.2* We have 10 independent observations from a normal population with variance 6.  $\bar{y}$  is observed to be 17. We find a 95% CI for  $\mu$ , the mean of the population.

[1]  $Par = \mu,$

[3]  $Est = \bar{y},$

[3]  $SE(Est) = \sqrt{6/10}$ . In this case,  $SE(Est)$  is known and not estimated.

[4]  $[Est - Par]/SE(Est) = [\bar{y} - \mu]/\sqrt{6/10} \sim N(0, 1) = t(\infty).$

The confidence coefficient is  $95\% = (1 - \alpha)100\%$ , so  $1 - \alpha = 0.95$  and  $\alpha = 0.05$ . The percentage point from the normal distribution that we require is  $t(1 - \frac{\alpha}{2}, \infty) = t(0.975, \infty) = 1.96$ . The limits of the 95% CI are, in general,

$$\bar{y} \pm 1.96\sqrt{6/10}$$

or, since  $\bar{y} = 17$ ,

$$17 \pm 1.96\sqrt{6/10}.$$

The  $\mu$  values in the interval (15.48, 18.52) are consistent with the data and the normal random sampling model as determined by an  $\alpha = 0.05$  test.

To see that statements (1) and (2) are algebraically equivalent, the argument runs as follows:

$$-t\left(1 - \frac{\alpha}{2}, df\right) < \frac{Est - m}{SE(Est)} < t\left(1 - \frac{\alpha}{2}, df\right)$$

if and only if  $-t\left(1 - \frac{\alpha}{2}, df\right) SE(Est) < Est - m < t\left(1 - \frac{\alpha}{2}, df\right) SE(Est)$ ;

if and only if  $t\left(1 - \frac{\alpha}{2}, df\right) SE(Est) > -Est + m > -t\left(1 - \frac{\alpha}{2}, df\right) SE(Est)$ ;

if and only if  $Est + t\left(1 - \frac{\alpha}{2}, df\right) SE(Est) > m > Est - t\left(1 - \frac{\alpha}{2}, df\right) SE(Est)$ ;

if and only if  $Est - t\left(1 - \frac{\alpha}{2}, df\right) SE(Est) < m < Est + t\left(1 - \frac{\alpha}{2}, df\right) SE(Est)$ .

## E.4 Final Comments on Significance Testing

The most arbitrary element in Fisherian testing is the choice of a test statistic. Although alternative hypotheses do not play a formal role in significance testing, interesting possible alternative hypotheses do inform the choice of test statistic.

For example, in linear models we often test a full model  $Y = X\beta + e$  against a reduced model  $Y = X_0\gamma + e$ , with  $e \sim N(0, \sigma^2 I)$  and  $C(X_0) \subset C(X)$ . Although we choose a test statistic based on comparing these models, the significance test is only a test of whether the data are consistent with the reduced model (and is a two-sided  $F$  test when  $r(X) - r(X_0) \geq 3$ ). Rejecting the  $F$  test does not suggest that the full model is correct, it only suggests that the reduced model is wrong. Nonetheless, it is of interest to see how the test behaves if the full model is correct. But models other than the full model can also cause the test to reject, see Appendix F, especially Section F.2. For example, it is of interest to examine the *power* of a test. The power of an  $\alpha$  level test at some alternative model is the probability of rejecting the null model when the alternative model is true. But in significance testing, there is no thought of accepting any alternative model. Any number of things can cause the rejection of the null model. Similar comments hold for testing generalized linear models.

When testing a null model based on a single parameter hypothesis  $H_0 : Par = m$ , interesting possible alternatives include  $Par \neq m$ . Our test statistic is designed to be sensitive to these alternatives, but problems with the null model other than  $Par \neq m$  can cause us to reject the null model.

In general, a test statistic can be any function of the data for which the distribution under the null model is known (or can be approximated). But finding a usable test statistic can be difficult. Having to choose between alternative test statistics for the same null model is something of a luxury. For example, to test the null model with equal means in a balanced one-way ANOVA, we can use either the  $F$  test of Chapter 4 or the Studentized range test of Section 5.5.

## Appendix F

# Significantly Insignificant Tests

**Abstract** Computer programs for fitting linear models typical focus on the significance of large  $F$  statistics. This appendix discusses why one should always be concerned about observing  $F$  statistics very near 0 (when the numerator degrees of freedom are 3 or more).

Philosophically, the test of a null model occurs almost in a vacuum. Either the data contradict the null model or they are consistent with it. The discussion of model testing in Section 3.2 largely assumes that the full model is true. While it is interesting to explore the behavior of the  $F$  test statistic when the full model is true, and indeed it is reasonable and appropriate to choose a test statistic that will work well when the full model is true, the act of rejecting the null model in no way implies that the full model is true. It is perfectly reasonable that the null (reduced) model can be rejected when the full model is false.

Throughout this book we have examined standard approaches to testing in which  $F$  tests are rejected only for large values. The rationale for this is based on the full model being true. We now examine the significance of small  $F$  statistics. Small  $F$  statistics can be caused by an unsuspected lack of fit or, when the mean structure of the reduced model is correct, they can be caused by not accounting for negatively correlated data or not accounting for heteroscedasticity. We also demonstrate that large  $F$  statistics can be generated by not accounting for positively correlated data or heteroscedasticity, even when the mean structure of the reduced model is correct.

Christensen (1995, 2005, 2008) argues that (non-Bayesian) testing should be viewed as an exercise in examining whether or not the data are consistent with a particular (predictive) model. While possible alternative hypotheses may drive the choice of a test statistic, any unusual values of the test statistic should be considered important. By this standard, perhaps the only general way to decide which values of the test statistic are unusual is to identify as unusual those values that have small probabilities or small densities under the model being tested.

The  $F$  test statistic is driven by the idea of testing the reduced model against the full model. However, given the test statistic, any unusual values of that statistic

should be recognized as indicating data that are inconsistent with the model being tested. If the full model is true, values of  $F$  much larger than 1 are inconsistent with the reduced model. Values of  $F$  much larger than 1 are consistent with the full model but, as we shall see, they are consistent with other models as well. Similarly, (when the numerator degrees of freedom are 3 or more) values of  $F$  much smaller than 1 are also inconsistent with the reduced model and we will examine models that can generate small  $F$  statistics.

I have been hesitant to discuss what I think of as a Fisherian  $F$  test, since nobody actually performs them. (That includes me, because it is so much easier to use the reported  $P$  values provided by standard computer programs.) Although the test statistic comes from considering both the reduced (null) model and the full model, once the test statistic is chosen, the full model no longer plays a role. From Theorem 3.2.1(ii), if the reduced model is true,

$$F \equiv \frac{Y'(M - M_0)Y/r(M - M_0)}{Y'(I - M)Y/r(I - M)} \sim F(r(M - M_0), r(I - M), 0).$$

We use the density to define “weird” values of the  $F$  distribution. The smaller the density, the weirder the observation. Write  $r_1 \equiv r(M - M_0)$  and  $r_2 \equiv r(I - M)$ , denote the density  $g(f|r_1, r_2)$ , and let  $F_{obs}$  denote the observed value of the  $F$  statistic. Since the  $P$  value of a test is the probability under the null model of seeing data as weird or weirder than we actually saw, and weirdness is defined by the density, the  $P$  value of the test is

$$P = \Pr[g(F|r_1, r_2) \leq g(F_{obs}|r_1, r_2)],$$

wherein  $F_{obs}$  is treated as fixed and known. This is computed under the only distribution we have, the  $F(r_1, r_2)$  distribution. An  $\alpha$  level test is defined as rejecting the null model precisely when  $P \leq \alpha$ .

If  $r_1 > 2$ , the  $F(r_1, r_2)$  density has the familiar shape that starts at 0, rises to a maximum in the vicinity of 1, and drops back down to zero for large values. Unless  $F_{obs}$  happens to be the mode, there are two values  $f_1 < f_2$  that have

$$g(F_{obs}|r_1, r_2) = g(f_1|r_1, r_2) = g(f_2|r_1, r_2).$$

(One of  $f_1$  and  $f_2$  will be  $F_{obs}$ .) In this case, the  $P$  value reduces to

$$P = \Pr[F \leq f_1] + \Pr[F \geq f_2].$$

In other words, the Fisherian  $F$  test is a two-sided  $F$  test, rejecting both for very small and very large values of  $F_{obs}$ . For  $r_1 = 1, 2$ , the Fisherian test agrees with the usual test because then the  $F(r_1, r_2)$  density starts high and decreases as  $f$  gets larger.

I should also admit that there remain open questions about the appropriateness of using densities, rather than actual probabilities, to define the weirdness of observations. In fact, I have long speculated whether Fisher’s “ $z$ ” distribution, i.e.,  $z \equiv \log(F)/2$  might not provide a more appropriate density for significance testing than the  $F$  density. The remainder of this appendix is closely related to Christensen (2003). See also Högfeldt (1979).

## F.1 Lack of Fit and Small $F$ Statistics

The standard assumption in testing models is that there is a full model  $Y = X\beta + e$ ,  $E(e) = 0$ ,  $\text{Cov}(e) = \sigma^2 I$  that fits the data. We then test the adequacy of a reduced model  $Y = X_0\gamma + e$ ,  $E(e) = 0$ ,  $\text{Cov}(e) = \sigma^2 I$  in which  $C(X_0) \subset C(X)$ , cf. Section 3.2. Based on second moment arguments, the test statistic is a ratio of variance estimates. We construct an unbiased estimate of  $\sigma^2$ ,  $Y'(I - M)Y/r(I - M)$ , and another statistic  $Y'(M - M_0)Y/r(M - M_0)$  that has  $E[Y'(M - M_0)Y/r(M - M_0)] = \sigma^2 + \beta'X'(M - M_0)X\beta/r(M - M_0)$ . Under the assumed covariance structure, this second statistic is an unbiased estimate of  $\sigma^2$  if and only if the reduced model is correct. The test statistic

$$F = \frac{Y'(M - M_0)Y/r(M - M_0)}{Y'(I - M)Y/r(I - M)}$$

is a (biased) estimate of

$$\frac{\sigma^2 + \beta'X'(M - M_0)X\beta/r(M - M_0)}{\sigma^2} = 1 + \frac{\beta'X'(M - M_0)X\beta}{\sigma^2 r(M - M_0)}.$$

Under the null model,  $F$  is an estimate of the number 1. When the full model is true, values of  $F$  much larger than 1 suggest that  $F$  is estimating something larger than 1, which suggests that  $\beta'X'(M - M_0)X\beta/\sigma^2 r(M - M_0) > 0$ , something that occurs if and only if the reduced model is false. The standard normality assumption leads to an exact central  $F$  distribution for the test statistic under the null model, so we are able to quantify how unusual it is to observe any  $F$  statistic greater than 1. Although the test is based on second moment considerations, under the normality assumption it is also the generalized likelihood ratio test, see Exercise 3.1, and a uniformly most powerful invariant test, see Lehmann (1986, Section 7.1).

In testing lack of fit, the same basic ideas apply except that we start with the (reduced) model  $Y = X\beta + e$ . The ideal situation would be to know that if  $Y = X\beta + e$  has the wrong mean structure, then a model of the form

$$Y = X\beta + W\delta + e, \quad C(W) \perp C(X) \tag{1}$$

fits the data where assuming  $C(W) \perp C(X)$  creates no loss of generality. Unfortunately, there is rarely anyone to tell us the true matrix  $W$ . Lack-or-fit testing is largely about constructing a full model, say,  $Y = X_*\beta_* + e$  with  $C(X) \subset C(X_*)$  based on reasonable assumptions about the nature of any lack of fit. The test for lack of fit is simply the test of  $Y = X\beta + e$  against the constructed model  $Y = X_*\beta_* + e$ . Typically, the constructed full model involves somehow generalizing the structure already observed in  $Y = X\beta + e$ . Section 6.7 discusses the rationale for several choices of constructed full models. For example, the traditional lack-or-fit test for simple linear regression begins with the replication model  $y_{ij} = \beta_0 + \beta_1 x_i + e_{ij}$ ,  $i = 1, \dots, a$ ,  $j = 1, \dots, N_i$ . It then assumes  $E(y_{ij}) = f(x_i)$  for some function  $f(\cdot)$ , in other

words, it assumes that the several observations associated with  $x_i$  have the same expected value. Making no additional assumptions leads to fitting the full model  $y_{ij} = \mu_i + e_{ij}$  and the traditional lack-of-fit test. Another way to think of this traditional test views the reduced model relative to the one-way ANOVA as having only the linear contrast important. The traditional lack-of-fit test statistic becomes

$$F = \frac{SSTrts - SS(lin)}{a - 2} \bigg/ MSE, \quad (2)$$

where  $SS(lin)$  is the sum of squares for the linear contrast. If there is no lack of fit in the reduced model,  $F$  should be near 1. If lack of fit exists because the more general mean structure of the one-way ANOVA fits the data better than the simple linear regression model, the  $F$  statistic tends to be larger than 1.

Unfortunately, if the lack of fit exists because of features that are not part of the original model, generalizing the structure observed in  $Y = X\beta + e$  is often inappropriate. Suppose that the simple linear regression model is balanced, i.e., all  $N_i = N$ , that for each  $i$  the data are taken in time order  $t_1 < t_2 < \dots < t_N$ , and that the lack of fit is due to the true model being

$$y_{ij} = \beta_0 + \beta_1 x_i + \delta t_j + e_{ij}, \quad \delta \neq 0. \quad (3)$$

Thus, depending on the sign of  $\delta$ , the observations within each group are subject to an increasing or decreasing trend. Note that in this model, for fixed  $i$ , the  $E(y_{ij})$ s are *not* the same for all  $j$ , thus invalidating the assumption of the traditional test. In fact, this causes the traditional lack of fit test to have a *small*  $F$  statistic. One way to see this is to view the problem in terms of a balanced two-way ANOVA. The true model (3) is a special case of the two-way ANOVA model  $y_{ij} = \mu + \alpha_i + \eta_j + e_{ij}$  in which the only nonzero terms are the linear contrast in the  $\alpha_i$ s and the linear contrast in the  $\eta_j$ s. Under model (3), the numerator of the statistic (2) gives an unbiased estimate of  $\sigma^2$  because  $SSTrts$  in (2) is  $SS(\alpha)$  for the two-way model and the only nonzero  $\alpha$  effect is being eliminated from the treatments. However, the mean squared error in the denominator of (2) is a weighted average of the error mean square from the two-way model and the mean square for the  $\eta_j$ s in the two-way model. The sum of squares for the significant linear contrast in the  $\eta_j$ s from model (3) is included in the error term of the lack-of-fit test (2), thus biasing the error term to estimate something larger than  $\sigma^2$ . In particular, the denominator has an expected value of  $\sigma^2 + \delta^2 a \sum_{j=1}^N (t_j - \bar{t})^2 / a(N - 1)$ . Thus, if the appropriate model is (3), the statistic in (2) estimates  $\sigma^2 / [\sigma^2 + \delta^2 a \sum_{j=1}^N (t_j - \bar{t})^2 / a(N - 1)]$  which is a number that is less than 1. Values of  $F$  much smaller than 1, i.e., very near 0, are consistent with a lack of fit that exists within the groups of the one-way ANOVA. Note that in this balanced case, true models involving interaction terms, e.g., models like

$$y_{ij} = \beta_0 + \beta_1 x_i + \delta t_j + \gamma x_i t_j + e_{ij},$$

also tend to make the  $F$  statistic small if either  $\delta \neq 0$  or  $\gamma \neq 0$ . Finally, if there exists lack of fit both between the groups of observations and within the groups, it can be very difficult to identify. For example, if  $\beta_2 \neq 0$  and either  $\delta \neq 0$  or  $\gamma \neq 0$  in the true model

$$y_{ij} = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \delta t_j + \gamma x_i t_j + e_{ij},$$

there is both a traditional lack of fit between the groups (the significant  $\beta_2 x_i^2$  term) and lack of fit within the groups ( $\delta t_j + \gamma x_i t_j$ ). In this case, neither the numerator nor the denominator in (2) is an estimate of  $\sigma^2$ .

More generally, start with a model  $Y = X\beta + e$ . This is tested against a larger model  $Y = X_*\beta_* + e$  with  $C(X) \subset C(X_*)$ , regardless of where the larger model comes from. The  $F$  statistic is

$$F = \frac{Y'(M_* - M)Y/r(M_* - M)}{Y'(I - M_*)Y/r(I - M_*)}.$$

We assume that the true model is (1). The  $F$  statistic estimates 1 if the original model  $Y = X\beta + e$  is correct. It estimates something greater than 1 if the larger model  $Y = X_*\beta_* + e$  is correct, i.e., if  $W\delta \in C(X)^\perp_{C(X_*)}$ .  $F$  estimates something less than 1 if  $W\delta \in C(X_*)^\perp$ , i.e., if  $W\delta$  is actually in the error space of the larger model, because then the numerator estimates  $\sigma^2$  but the denominator estimates

$$\sigma^2 + \delta'W'(I - M_*)W\delta/r(I - M_*) = \sigma^2 + \delta'W'W\delta/r(I - M_*).$$

If  $W\delta$  is in neither of  $C(X)^\perp_{C(X_*)}$  nor  $C(X_*)^\perp$ , it is not clear how the test will behave because neither the numerator nor the denominator estimates  $\sigma^2$ . Christensen (1989, 1991) contains related discussion of these concepts.

The main point is that, when testing a full model  $Y = X\beta + e, E(e) = 0, \text{Cov}(e) = \sigma^2 I$  against a reduced model  $Y = X_0\gamma + e, C(X_0) \subset C(X)$ , if the  $F$  statistic is small, it suggests that  $Y = X_0\gamma + e$  may suffer from lack of fit in which the lack of fit exists in the error space of  $Y = X\beta + e$ . We will see in the next section that other possible explanations for a small  $F$  statistic are the existence of “negative correlation” in the data or heteroscedasticity.

## F.2 The Effect of Correlation and Heteroscedasticity on $F$ Statistics

The test of a reduced model assumes that the full model  $Y = X\beta + e, E(e) = 0, \text{Cov}(e) = \sigma^2 I$  holds and tests the adequacy of a reduced model  $Y = X_0\gamma + e, E(e) = 0, \text{Cov}(e) = \sigma^2 I, C(X_0) \subset C(X)$ . Rejecting the reduced model does not imply that the full model is correct. The mean structure of the reduced model may

be perfectly valid, but the  $F$  statistic can become large or small because the assumed covariance structure is incorrect.

We begin with a concrete example, one-way ANOVA. Let  $i = 1, \dots, a$ ,  $j = 1, \dots, N$ , and  $n \equiv aN$ . Consider a reduced model  $y_{ij} = \mu + e_{ij}$  which in matrix terms we write  $Y = J\mu + e$ , and a full model  $y_{ij} = \mu_i + e_{ij}$ , which we write  $Y = Z\gamma + e$ . In matrix terms the usual one-way ANOVA  $F$  statistic is

$$F = \frac{Y'[M_Z - (1/n)J_n^n]Y/(a-1)}{Y'(I - M_Z)Y/a(N-1)}. \quad (1)$$

We now assume that the true model is  $Y = J\mu + e$ ,  $E(e) = 0$ ,  $\text{Cov}(e) = \sigma^2 V$  and examine the behavior of the  $F$  statistic (1).

For a homoscedastic balanced one-way ANOVA we want to characterize the concepts of overall positive correlation, positive correlation within groups, and positive correlation for evaluating differences between groups. Consider first a simple example with  $a = 2$ ,  $N = 2$ . The first two observations are a group and the last two are a group. Consider a covariance structure

$$V_1 = \begin{bmatrix} 1 & 0.9 & 0.1 & 0.09 \\ 0.9 & 1 & 0.09 & 0.1 \\ 0.1 & 0.09 & 1 & 0.9 \\ 0.09 & 0.1 & 0.9 & 1 \end{bmatrix}.$$

There is an overall positive correlation, high positive correlation between the two observations in each group, and weak positive correlation between the groups. A second example,

$$V_2 = \begin{bmatrix} 1 & 0.1 & 0.9 & 0.09 \\ 0.1 & 1 & 0.09 & 0.9 \\ 0.9 & 0.09 & 1 & 0.1 \\ 0.09 & 0.9 & 0.1 & 1 \end{bmatrix},$$

has an overall positive correlation but weak positive correlation between the two observations in each group, with high positive correlation between some observations in different groups.

We now make a series of definitions for homoscedastic balanced one-way ANOVA based on the projection operators in (1) and  $V$ . Overall positive correlation is characterized by  $\text{Var}(\bar{y}_{..}) > \sigma^2/n$ , which in matrix terms is written

$$n \frac{\text{Var}(\bar{y}_{..})}{\sigma^2} = \text{tr}[(1/n)JJ'V] > \frac{1}{n} \text{tr}(V) \text{tr}[(1/n)JJ'] = \frac{1}{n} \text{tr}(V). \quad (2)$$

Overall negative correlation is characterized by the reverse inequality. For homoscedastic models the term  $\text{tr}(V)/n$  is 1. For heteroscedastic models the term on the right is the average variance of the observations divided by  $\sigma^2$ .



Positive correlation within groups is characterized by  $\sum_{i=1}^a \text{Var}(\bar{y}_{i.})/a > \sigma^2/N$ , which in matrix terms is written

$$\sum_{i=1}^a N \frac{\text{Var}(\bar{y}_{i.})}{\sigma^2} = \text{tr}[M_Z V] > \frac{1}{n} \text{tr}(V) \text{tr}[M_Z] = \frac{a}{n} \text{tr}(V). \tag{3}$$

Negative correlation within groups is characterized by the reverse inequality.

Positive correlation for evaluating differences between groups is characterized by

$$\frac{\sum_{i=1}^a \text{Var}(\bar{y}_{i.} - \bar{y}_{..})}{a} > \frac{a-1}{a} \frac{\sigma^2}{N}.$$

Note that equality obtains if  $V = I$ . In matrix terms, this is written

$$\begin{aligned} \frac{N}{\sigma^2} \sum_{i=1}^a \text{Var}(\bar{y}_{i.} - \bar{y}_{..}) &= \text{tr}([M_Z - (1/n)J J']V) \\ &> \frac{1}{n} \text{tr}(V) \text{tr}[M_Z - (1/n)J J'] = \frac{a-1}{n} \text{tr}(V) \end{aligned} \tag{4}$$

and negative correlation for evaluating differences between groups is characterized by the reverse inequality. If all the observations in different groups are uncorrelated, there will be positive correlation for evaluating differences between groups if and only if there is positive correlation within groups. This follows because having a block diagonal covariance matrix  $\sigma^2 V$  implies that  $\text{tr}(M_Z V) = \text{tr}[(1/N)Z' V Z] = \text{atr}[(1/n)J' V J] = \text{atr}[(1/n)J J' V]$ .

For our example  $V_1$ ,

$$2.09 = (1/4)[4(2.09)] = \text{tr}[(1/n)J_n^n V_1] > \frac{1}{n} \text{tr}(V_1) = 4/4 = 1,$$

so there is an overall positive correlation,

$$3.8 = 2(1/2)[3.8] = \text{tr}[M_Z V_1] > \frac{a}{n} \text{tr}(V_1) = (2/4)4 = 2,$$

so there is positive correlation within groups, and

$$1.71 = 3.8 - 2.09 = \text{tr}([M_Z - (1/n)J_n^n]V_1) > \frac{a-1}{n} \text{tr}(V_1) = (1/4)4 = 1,$$

so there is positive correlation for evaluating differences between groups.

For the second example  $V_2$ ,

$$2.09 = (1/4)[4(2.09)] = \text{tr}[(1/n)J_n^n V_2] > \frac{1}{n} \text{tr}(V_2) = 4/4 = 1,$$

so there is an overall positive correlation,

$$2.2 = 2(1/2)[2.2] = \text{tr}[M_Z V_2] > \frac{a}{n} \text{tr}(V_2) = (2/4)4 = 2,$$

so there is positive correlation within groups, but

$$0.11 = 2.2 - 2.09 = \text{tr}([M_Z - (1/n)J_n^n]V_2) < \frac{a-1}{n} \text{tr}(V_2) = (1/4)4 = 1,$$

so positive correlation for evaluating differences between groups does not exist.

The existence of positive correlation within groups and positive correlation for evaluating differences between groups causes the one-way ANOVA  $F$  statistic in (1) to get large even when there are no differences in the group means. Assuming that the correct model is  $Y = J\mu + e$ ,  $E(e) = 0$ ,  $\text{Cov}(e) = \sigma^2 V$ , by Theorem 1.3.1, the numerator of the  $F$  statistic estimates

$$\begin{aligned} E\{Y'[M_Z - (1/n)J_n^n]Y/(a-1)\} &= \text{tr}\{[M_Z - (1/n)J_n^n]V\}/(a-1) \\ &> \frac{a-1}{n} \text{tr}(V)/(a-1) = \text{tr}(V)/n \end{aligned}$$

and the denominator of the  $F$  statistic estimates

$$\begin{aligned} E\{Y'(I - M_Z)Y/a(N-1)\} &= \text{tr}\{[I - M_Z]V\}/a(N-1) \\ &= (\text{tr}\{V\} - \text{tr}\{[M_Z]V\})/a(N-1) \\ &< \left(\text{tr}\{V\} - \frac{a}{n} \text{tr}(V)\right)/a(N-1) \\ &= \frac{n-a}{n} \text{tr}(V)/a(N-1) = \text{tr}(V)/n. \end{aligned}$$

In (1),  $F$  is an estimate of

$$\frac{E\{Y'[M_Z - (1/n)J_n^n]Y/(a-1)\}}{E\{Y'(I - M_Z)Y/a(N-1)\}} = \frac{\text{tr}\{[M_Z - (1/n)J_n^n]V\}/(a-1)}{\text{tr}\{[I - M_Z]V\}/a(N-1)} > \frac{\text{tr}(V)/n}{\text{tr}(V)/n} = 1,$$

so having both positive correlation within groups and positive correlation for evaluating differences between groups tends to make  $F$  statistics large. Exactly analogous computations show that having both negative correlation within groups and negative correlation for evaluating differences between groups tends to make  $F$  statistics less than 1.

Another example elucidates some additional points. Suppose the observations have the AR(1) correlation structure discussed in Subsection 12.3.1:

$$V_3 = \begin{bmatrix} 1 & \rho & \rho^2 & \rho^3 \\ \rho & 1 & \rho & \rho^2 \\ \rho^2 & \rho & 1 & \rho \\ \rho^3 & \rho^2 & \rho & 1 \end{bmatrix}.$$

Using the same grouping structure as before, when  $0 < \rho < 1$ , we have overall positive correlation because

$$1 + \frac{\rho}{2}(3 + 2\rho + \rho^2) = \text{tr}[(1/n)JJ'V_3] > 1,$$

and we have positive correlation within groups because

$$2(1 + \rho) = \text{tr}[M_Z V_3] > 2.$$

If  $-1 < \rho < 0$ , the inequalities are reversed. Similarly, for  $-1 < \rho < 0$  we have negative correlation for evaluating differences between groups because

$$1 + \frac{\rho}{2}(1 - 2\rho - \rho^2)^2 = \text{tr}([M_Z - (1/n)JJ']V_3) < 1.$$

However, we only get positive correlation for evaluating differences between groups when  $0 < \rho < \sqrt{2} - 1$ . Thus, for negative  $\rho$  we tend to get small  $F$  statistics, for  $0 < \rho < \sqrt{2} - 1$  we tend to get large  $F$  statistics, and for  $\sqrt{2} - 1 < \rho < 1$  the result is not clear.

To illustrate, suppose  $\rho = 1$  and the observations all have the same mean, then with probability 1, all the observations are equal and, in particular,  $\bar{y}_i = \bar{y}_..$  with probability 1. It follows that

$$0 = \frac{\sum_{i=1}^a \text{Var}(\bar{y}_i - \bar{y}_..)}{a} < \frac{a-1}{a} \frac{\sigma^2}{N}$$

and no positive correlation exists for evaluating differences between groups. More generally, for very strong positive correlations, both the numerator and the denominator of the  $F$  statistic estimate numbers close to 0 and both are smaller than they would be under  $V = I$ . On the other hand, it is not difficult to see that, for  $\rho = -1$ , the  $F$  statistic is 0.

In the balanced heteroscedastic one-way ANOVA,  $V$  is diagonal. This generates equality between the left sides and right sides of (2), (3), and (4), so under heteroscedasticity  $F$  still estimates the number 1. We now generalize the ideas of within group correlation and correlation for evaluating differences between groups, and see that heteroscedasticity can affect unbalanced one-way ANOVA.

In general, we test a full model  $Y = X\beta + e$ ,  $E(e) = 0$ ,  $\text{Cov}(e) = \sigma^2 I$  against a reduced model  $Y = X_0\gamma + e$ , in which  $C(X_0) \subset C(X)$ . We examine the  $F$  statistic when the true model is  $Y = X_0\gamma + e$ ,  $E(e) = 0$ ,  $\text{Cov}(e) = \sigma^2 V$ . Using arguments similar to those for balanced one-way ANOVA, having

$$\text{tr}[MV] > \frac{1}{n} \text{tr}(V) \text{tr}[M] = \frac{r(X)}{n} \text{tr}(V)$$

and

$$\text{tr}[(M - M_0)V] > \frac{1}{n} \text{tr}(V) \text{tr}[M - M_0] = \frac{r(X) - r(X_0)}{n} \text{tr}(V)$$

causes large  $F$  statistics even when the mean structure of the reduced model is true, and reversing the inequalities causes small  $F$  statistics. These are merely sufficient conditions so that the tests intuitively behave certain ways. The actual behavior of the tests under normal distributions can be determined numerically, cf. Christensen and Bedrick (1997).

These covariance conditions can be caused by patterns of positive and negative correlations as discussed earlier, but they can also be caused by heteroscedasticity. For example, consider the behavior of the unbalanced one-way ANOVA  $F$  test when the observations are uncorrelated but heteroscedastic. For concreteness, assume that  $\text{Var}(y_{ij}) = \sigma_i^2$ . Because the observations are uncorrelated, we need only check the condition

$$\text{tr}[MV] \equiv \text{tr}[M_Z V] > \frac{1}{n} \text{tr}(V) \text{tr}[M_Z] = \frac{a}{n} \text{tr}(V),$$

which amounts to

$$\sum_{i=1}^a \sigma_i^2 / a > \sum_{i=1}^a \frac{N_i}{n} \sigma_i^2.$$

Thus, when the groups' means are equal,  $F$  statistics will get large if many observations are taken in groups with small variances and few observations are taken on groups with large variances.  $F$  statistics will get small if the reverse relationship holds.

The general condition

$$\text{tr}[MV] > \frac{1}{n} \text{tr}(V) \text{tr}[M] = \frac{r(X)}{n} \text{tr}(V)$$

is equivalent to

$$\frac{\sum_{i=1}^n \text{Var}(x'_i \hat{\beta})}{r(X)} > \frac{\sum_{i=1}^n \text{Var}(y_i)}{n}.$$

So, under homoscedasticity, positive correlation in the full model amounts to having an average variance for the predicted values (averaging over the rank of the covariance matrix of the predicted values) that is larger than the common variance of the observations. Negative correlation in the full model involves reversing the inequality. Similarly, having positive correlation for distinguishing the full model from the reduced model means

$$\frac{\sum_{i=1}^n \text{Var}(x'_i \hat{\beta} - x'_{0i} \hat{\gamma})}{r(X) - r(X_0)} = \frac{\text{tr}[(M - M_0)V]}{r(M - M_0)} > \frac{\text{tr}(V)}{n} = \frac{\sum_{i=1}^n \text{Var}(y_i)}{n}.$$

# Appendix G

## Randomization Theory Models

**Abstract** This appendix introduces randomization theory models in which, rather than assuming the existence of an error term with certain properties, the random variability in the data is constructed either by random sampling from a population or by randomly assigning treatments to experimental units.

The division of labor in statistics has traditionally designated randomization theory as an area of nonparametric statistics. Randomization theory is also of special interest in the theory of experimental design because randomization has been used to justify the analysis of designed experiments.

It can be argued that the linear models given in Chapter 8 are merely good approximations to more appropriate models based on randomization theory. One aspect of this argument is that the  $F$  tests based on the theory of normal errors are a good approximation to randomization (permutation) tests. Investigating this is beyond the scope of a linear models book, cf. Hinkelmann and Kempthorne (2005) and Puri and Sen (1971). Another aspect of the approximation argument is that the BLUEs under randomization theory are precisely the least squares estimates. By Theorem 10.4.5, to establish this we need to show that  $C(VX) \subset C(X)$  for the model

$$Y = X\beta + e, \quad E(e) = 0, \quad \text{Cov}(e) = V,$$

where  $V$  is the covariance matrix under randomization theory. This argument will be examined here for two experimental design models: the model for a completely randomized design and the model for a randomized complete block design. First, we introduce the subject with a discussion of simple random sampling.

### G.1 Simple Random Sampling

Randomization theory for a simple random sample assumes that observations  $y_i$  are picked at random (without replacement) from a larger finite population. Suppose the elements of the population are  $s_1, s_2, \dots, s_N$ . We can define elementary sampling

random variables for  $i = 1, \dots, n$  and  $j = 1, \dots, N$ ,

$$\delta_j^i = \begin{cases} 1, & \text{if } y_i = s_j \\ 0, & \text{otherwise.} \end{cases}$$

Under simple random sampling without replacement

$$E[\delta_j^i] = \Pr[\delta_j^i = 1] = \frac{1}{N}.$$

$$E[\delta_j^i \delta_{j'}^{i'}] = \Pr[\delta_j^i \delta_{j'}^{i'} = 1] = \begin{cases} 1/N, & \text{if } (i, j) = (i', j') \\ 1/N(N-1), & \text{if } i \neq i' \text{ and } j \neq j' \\ 0, & \text{otherwise.} \end{cases}$$

If we write  $\mu = \sum_{j=1}^N s_j/N$  and  $\sigma^2 = \sum_{j=1}^N (s_j - \mu)^2/N$ , then

$$y_i = \sum_{j=1}^N \delta_j^i s_j = \mu + \sum_{j=1}^N \delta_j^i (s_j - \mu).$$

Letting  $e_i = \sum_{j=1}^N \delta_j^i (s_j - \mu)$  gives the linear model

$$y_i = \mu + e_i.$$

The population mean  $\mu$  is a fixed unknown constant. The  $e_i$ s have the properties

$$E[e_i] = E\left[\sum_{j=1}^N \delta_j^i (s_j - \mu)\right] = \sum_{j=1}^N E[\delta_j^i] (s_j - \mu) = \sum_{j=1}^N (s_j - \mu)/N = 0,$$

$$\text{Var}(e_i) = E[e_i^2] = \sum_{j=1}^N \sum_{j'=1}^N (s_j - \mu)(s_{j'} - \mu) E[\delta_j^i \delta_{j'}^i] = \sum_{j=1}^N (s_j - \mu)^2/N = \sigma^2.$$

For  $i \neq i'$ ,

$$\begin{aligned} \text{Cov}(e_i, e_{i'}) &= E[e_i e_{i'}] = \sum_{j=1}^N \sum_{j'=1}^N (s_j - \mu)(s_{j'} - \mu) E[\delta_j^i \delta_{j'}^{i'}] \\ &= [N(N-1)]^{-1} \sum_{j \neq j'} (s_j - \mu)(s_{j'} - \mu) \\ &= [N(N-1)]^{-1} \left( \left[ \sum_{j=1}^N (s_j - \mu) \right]^2 - \sum_{j=1}^N (s_j - \mu)^2 \right) \\ &= -\sigma^2/(N-1). \end{aligned}$$

In matrix terms, the linear model can be written

$$Y = J\mu + e, \quad E(e) = 0, \quad \text{Cov}(e) = \sigma^2 V,$$

where

$$V = \begin{bmatrix} 1 & -(N-1)^{-1} & -(N-1)^{-1} & \cdots & -(N-1)^{-1} \\ -(N-1)^{-1} & 1 & -(N-1)^{-1} & \cdots & -(N-1)^{-1} \\ -(N-1)^{-1} & -(N-1)^{-1} & 1 & \cdots & -(N-1)^{-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ -(N-1)^{-1} & -(N-1)^{-1} & -(N-1)^{-1} & \cdots & 1 \end{bmatrix}.$$

Clearly  $VJ = [(N-n)/(N-1)]J$ , so the BLUE of  $\mu$  is  $\bar{y}$ .

## G.2 Completely Randomized Designs

Suppose that there are  $t$  treatments, each to be randomly assigned to  $N$  units out of a collection of  $n = tN$  experimental units. A one-way ANOVA model for this design is

$$y_{ij} = \mu_i + e_{ij}, \quad (1)$$

$i = 1, \dots, t, j = 1, \dots, N$ . Suppose further that the  $i$ th treatment has an effect  $\tau_i$  and that the experimental units without treatment effects would have readings  $s_1, \dots, s_n$ . The elementary sampling random variables are

$$\delta_k^{ij} = \begin{cases} 1, & \text{if replication } j \text{ of treatment } i \text{ is assigned to unit } k \\ 0, & \text{otherwise.} \end{cases}$$

With this restricted random sampling,

$$E[\delta_k^{ij}] = \Pr[\delta_k^{ij} = 1] = \frac{1}{n}$$

$$E[\delta_k^{ij} \delta_{k'}^{i'j'}] = \Pr[\delta_k^{ij} \delta_{k'}^{i'j'} = 1] = \begin{cases} 1/n, & \text{if } (i, j, k) = (i', j', k') \\ 1/n(n-1), & \text{if } k \neq k' \text{ and } (i, j) \neq (i', j') \\ 0, & \text{otherwise.} \end{cases}$$

We can write

$$y_{ij} = \tau_i + \sum_{k=1}^n \delta_k^{ij} s_k.$$

Taking  $\mu = \sum_{k=1}^n s_k/n$  and  $\mu_i = \mu + \tau_i$  gives

$$y_{ij} = \mu_i + \sum_{k=1}^n \delta_k^{ij} (s_k - \mu).$$

To obtain the linear model (1), let  $\text{sp } e_{ij} = \sum_{k=1}^n \delta_k^{ij} (s_k - \mu)$ . Write  $\sigma^2 = \sum_{k=1}^n (s_k - \mu)^2/n$ . Then

$$\text{E}[e_{ij}] = \text{E}\left[\sum_{k=1}^n \delta_k^{ij} (s_k - \mu)\right] = \sum_{k=1}^n \text{E}[\delta_k^{ij}] (s_k - \mu) = \sum_{k=1}^n (s_k - \mu)/n = 0,$$

$$\text{Var}(e_{ij}) = \text{E}[e_{ij}^2] = \sum_{k=1}^n \sum_{k'=1}^n (s_k - \mu)(s_{k'} - \mu) \text{E}[\delta_k^{ij} \delta_{k'}^{ij}] = \sum_{k=1}^n (s_k - \mu)^2/n = \sigma^2.$$

For  $(i, j) \neq (i', j')$ ,

$$\begin{aligned} \text{Cov}(e_{ij}, e_{i'j'}) &= \text{E}[e_{ij} e_{i'j'}] = \sum_{k=1}^n \sum_{k'=1}^n (s_k - \mu)(s_{k'} - \mu) \text{E}[\delta_k^{ij} \delta_{k'}^{i'j'}] \\ &= [n(n-1)]^{-1} \sum_{k \neq k'} (s_k - \mu)(s_{k'} - \mu) \\ &= [n(n-1)]^{-1} \left( \left[ \sum_{k=1}^n (s_k - \mu) \right]^2 - \sum_{k=1}^n (s_k - \mu)^2 \right) \\ &= -\sigma^2/(n-1). \end{aligned}$$

In matrix terms, writing  $Y = (y_{11}, y_{12}, \dots, y_{tN})'$ , we get

$$Y = X \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_t \end{bmatrix} + e, \quad \text{E}(e) = 0, \quad \text{Cov}(e) = \sigma^2 V,$$

where

$$\begin{aligned} V &= \begin{bmatrix} 1 & -1/(n-1) & -1/(n-1) & \cdots & -1/(n-1) \\ -1/(n-1) & 1 & -1/(n-1) & \cdots & -1/(n-1) \\ -1/(n-1) & -1/(n-1) & 1 & \cdots & -1/(n-1) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ -1/(n-1) & -1/(n-1) & -1/(n-1) & \cdots & 1 \end{bmatrix} \\ &= \frac{n}{n-1} I - \frac{1}{n-1} J_n^n. \end{aligned}$$



It follows that

$$VX = \frac{n}{n-1}X - \frac{1}{n-1}J_n^n X.$$

Since  $J \in C(X)$ ,  $C(VX) \subset C(X)$ , and least squares estimates are BLUEs. Standard errors for estimable functions can be found as in Section 11.1 using the fact that this model involves only one cluster.

**Exercise G.1** Establish whether least squares estimates are BLUEs in a completely randomized design with unequal numbers of observations on the treatments.

### G.3 Randomized Complete Block Designs

Suppose there are  $a$  treatments and  $b$  blocks. The experimental units must be grouped into  $b$  blocks, each of  $a$  units. Let the experimental unit effects be  $s_{kj}$ ,  $k = 1, \dots, a$ ,  $j = 1, \dots, b$ . Treatments are assigned at random to the  $a$  units in each block. The elementary sampling random variables are

$$\delta_{kj}^i = \begin{cases} 1, & \text{if treatment } i \text{ is assigned to unit } k \text{ in block } j \\ 0, & \text{otherwise.} \end{cases}$$

$$E[\delta_{kj}^i] = \Pr[\delta_{kj}^i = 1] = \frac{1}{a}.$$

$$E[\delta_{kj}^i \delta_{k'j'}^{i'}] = \Pr[\delta_{kj}^i \delta_{k'j'}^{i'} = 1] = \begin{cases} 1/a, & \text{if } (i, j, k) = (i', j', k') \\ 1/a^2, & \text{if } j \neq j' \\ 1/a(a-1), & \text{if } j = j', k \neq k', i \neq i' \\ 0, & \text{otherwise.} \end{cases}$$

If  $\alpha_i$  is the additive effect of the  $i$ th treatment and  $\beta_j \equiv \bar{s}_{.j}$ , then

$$y_{ij} = \alpha_i + \beta_j + \sum_{k=1}^a \delta_{kj}^i (s_{kj} - \beta_j).$$

Letting  $e_{ij} = \sum_{k=1}^a \delta_{kj}^i (s_{kj} - \beta_j)$  gives the linear model

$$y_{ij} = \alpha_i + \beta_j + e_{ij}. \quad (1)$$

The column space of the design matrix for this model is precisely that of the model considered in Section 8.3. Let  $\sigma_j^2 = \sum_{k=1}^a (s_{kj} - \beta_j)^2/a$ . Then

$$\begin{aligned} E[e_{ij}] &= \sum_{k=1}^a (s_{kj} - \beta_j) / a = 0, \\ \text{Var}(e_{ij}) &= \sum_{k=1}^a \sum_{k'=1}^a (s_{kj} - \beta_j)(s_{k'j} - \beta_j) E[\delta_{kj}^i \delta_{k'j}^i] \\ &= \sum_{k=1}^a (s_{kj} - \beta_j)^2 / a = \sigma_j^2. \end{aligned}$$

For  $j \neq j'$ ,

$$\begin{aligned} \text{Cov}(e_{ij}, e_{i'j'}) &= \sum_{k=1}^a \sum_{k'=1}^a (s_{kj} - \beta_j)(s_{k'j'} - \beta_{j'}) E[\delta_{kj}^i \delta_{k'j'}^{i'}] \\ &= a^{-2} \sum_{k=1}^a (s_{kj} - \beta_j) \sum_{k'=1}^a (s_{k'j'} - \beta_{j'}) \\ &= 0. \end{aligned}$$

For  $j = j', i \neq i'$ ,

$$\begin{aligned} \text{Cov}(e_{ij}, e_{i'j'}) &= \sum_{k=1}^a \sum_{k'=1}^a (s_{kj} - \beta_j)(s_{k'j} - \beta_j) E[\delta_{kj}^i \delta_{k'j}^{i'}] \\ &= \sum_{k \neq k'} (s_{kj} - \beta_j)(s_{k'j} - \beta_j) / a(a-1) \\ &= [a(a-1)]^{-1} \left( \left[ \sum_{k=1}^a (s_{kj} - \beta_j) \right]^2 - \sum_{k=1}^a (s_{kj} - \beta_j)^2 \right) \\ &= -\sigma_j^2 / (a-1). \end{aligned}$$

Before proceeding, we show that although the terms  $\beta_j$  are not known, the differences among these are known constants under randomization theory. For any unit  $k$  in block  $j$ , some treatment is assigned, so  $\sum_{i=1}^a \delta_{kj}^i = 1$ .

$$\begin{aligned} \bar{y}_{\cdot j} &= \frac{1}{a} \left[ \sum_{i=1}^a \left( \alpha_i + \beta_j + \sum_{k=1}^a \delta_{kj}^i (s_{kj} - \beta_j) \right) \right] \\ &= \frac{1}{a} \left[ \sum_{i=1}^a \alpha_i + a\beta_j + \sum_{k=1}^a (s_{kj} - \beta_j) \sum_{i=1}^a \delta_{kj}^i \right] \\ &= \bar{\alpha}_{\cdot} + \beta_j + \sum_{k=1}^a (s_{kj} - \beta_j) \\ &= \bar{\alpha}_{\cdot} + \beta_j. \end{aligned}$$

Therefore,  $\bar{y}_{.j} - \bar{y}_{.j'} = \beta_j - \beta_{j'} = \bar{s}_{.j} - \bar{s}_{.j'}$ . Since these differences are fixed and known, there is no basis for a test of  $H_0 : \beta_1 = \dots = \beta_b$ . In fact, the linear model is not just model (1) but model (1) subject to these estimable constraints on the  $\beta$ s.

To get best linear unbiased estimates we need to assume that  $\sigma_1^2 = \sigma_2^2 = \dots = \sigma_b^2 = \sigma^2$ . We can now write the linear model in matrix form and establish that least squares estimates of treatment means and contrasts in the  $\alpha_i$ s are BLUEs. In the discussion that follows, we use notation from Section 7.1. Model (1) can be rewritten

$$Y = X\eta + e, \quad E(e) = 0, \quad \text{Cov}(e) = V, \quad (2)$$

where  $\eta = [\mu, \alpha_1, \dots, \alpha_a, \beta_1, \dots, \beta_b]'$ . If we let  $X_2$  be the columns of  $X$  corresponding to  $\beta_1, \dots, \beta_b$ , then (cf. Section 11.1)

$$V = \sigma^2 [a/(a-1)] [I - (1/a)X_2X_2'] = \sigma^2 [a/(a-1)] [I - M_\mu - M_\beta].$$

If model (2) were the appropriate model, checking that  $C(VX) \subset C(X)$  would be trivial based on the fact that  $C(X_2) \subset C(X)$ . However, we must account for the estimable constraints on the model discussed above. In particular, consider

$$M_\beta X\eta = [t_{ij}],$$

where

$$t_{ij} = \beta_j - \bar{\beta}_{.} = \bar{y}_{.j} - \bar{y}_{..} = \bar{s}_{.j} - \bar{s}_{..}$$

This is a fixed known quantity. Proceeding as in Section 3.3, the model is subject to the estimable constraint

$$M_\beta X\eta = M_\beta Y.$$

Normally a constraint has the form  $A'\beta = d$ , where  $d$  is known. Here  $d = M_\beta Y$ , which appears to be random but, as discussed,  $M_\beta Y$  is not random; it is fixed and upon observing  $Y$  it is known.

The equivalent reduced model involves  $X_0 = (I - M_{MP})X = (I - M_\beta)X$  and a known vector  $Xb = M_\beta Y$ . Thus, the constrained model is equivalent to

$$(Y - M_\beta Y) = (I - M_\beta)X\gamma + e. \quad (3)$$

We want to show that least squares estimates of contrasts in the  $\alpha$ s based on  $Y$  are BLUEs with respect to this model. First we show that least squares estimates from model (3) based on  $(Y - M_\beta Y) = (I - M_\beta)Y$  are BLUEs. We need to show that

$$C(V(I - M_\beta)X) = C[(I - M_\mu - M_\beta)(I - M_\beta)X] \subset C[(I - M_\beta)X].$$

Because  $(I - M_\mu - M_\beta)(I - M_\beta) = (I - M_\mu - M_\beta)$ , we have

$$C(V(I - M_\beta)X) = C[(I - M_\mu - M_\beta)X],$$

and because  $C(I - M_\mu - M_\beta) \subset C(I - M_\beta)$  we have

$$C[(I - M_\mu - M_\beta)X] \subset C[(I - M_\beta)X].$$

To finish the proof that least squares estimates based on  $Y$  are BLUEs, note that the estimation space for model (3) is  $C[(I - M_\beta)X] = C(M_\mu + M_\alpha)$ . BLUEs are based on

$$(M_\mu + M_\alpha)(I - M_\beta)Y = (M_\mu + M_\alpha)Y.$$

Thus, any linear parametric function in model (2) that generates a constraint on  $C(M_\mu + M_\alpha)$  has a BLUE based on  $(M_\mu + M_\alpha)Y$  (cf. Exercise 3.9.5). In particular, this is true for contrasts in the  $\alpha$ s. Standard errors for estimable functions are found in a manner analogous to Section 11.1. This is true even though model (3) is not the form considered in Section 11.1 and is a result of the orthogonality relationships that are present.

The assumption that  $\sigma_1^2 = \sigma_2^2 = \dots = \sigma_b^2$  is a substantial one. Least squares estimates without this assumption are unbiased, but may be far from optimal. It is important to choose blocks so that their variances are approximately equal.

**Exercise G.2** Find the standard error for a contrast in the  $\alpha_i$ s of model (1).

## References

- Christensen, R. (1989). Lack of fit tests based on near or exact replicates. *The Annals of Statistics*, 17, 673–683.
- Christensen, R. (1991). Small sample characterizations of near replicate lack of fit tests. *Journal of the American Statistical Association*, 86, 752–756.
- Christensen, R. (1995). Comment on Inman (1994). *The American Statistician*, 49, 400.
- Christensen, R. (2003). Significantly insignificant  $F$  tests. *The American Statistician*, 57, 27–32.
- Christensen, R. (2005). Testing Fisher, Neyman, Pearson, and Bayes. *The American Statistician*, 59, 121–126.
- Christensen, R. (2008). Review of *Principals of statistical inference* by D. R. Cox. *Journal of the American Statistical Association*, 103, 1719–1723.
- Christensen, R. (2015). *Analysis of variance, design, and regression: Linear modeling for unbalanced data* (2nd ed.). Boca Raton, FL: Chapman and Hall/CRC Pres.
- Christensen, R., & Bedrick, E. J. (1997). Testing the independence assumption in linear models. *Journal of the American Statistical Association*, 92, 1006–1016.
- Davies, R. B. (1980). The distribution of linear combinations of  $\chi^2$  random variables. *Applied Statistics*, 29, 323–333.
- Hinkelmann, K., & Kempthorne, O. (2005). *Design and analysis of experiments: Volume 2, Advanced experimental design*. Hoboken, NJ: Wiley.
- Högfeldt, P. (1979). On low  $F$ -test values in linearmodels. *Scandinavian Journal of Statistics*, 6, 175–178.
- Lehmann, E. L. (1986). *Testing statistical hypotheses* (2nd ed.). New York: Wiley.
- Puri, M. L., & Sen, P. K. (1971). *Nonparametric methods in multivariate analysis*. New York: Wiley.
- Rao, C. R. (1973). *Linear statistical inference and its applications* (2nd ed.). New York: Wiley.

# Author Index

## A

Aitchison, J., 57, 192  
Anderson, T. W., 160  
Andrews, D.F., 389  
Arnold, S. F., 22, 356  
Atkinson, A. C., 343, 361, 362, 388  
Atwood, C. L., 184

## B

Bailey, D. W., 224  
Bedrick, E. J., 16, 49, 51, 87, 230, 233, 237,  
252, 276, 362, 363, 389, 421, 426,  
508  
Belsley, D. A., 396–398  
Benedetti, J. K., 237  
Berger, J. O., 48, 73  
Berger, R. L., 73  
Berk, K. N., 383  
Berry, D. A., 48  
Blom, G., 355  
Blouin, D. C., 314  
Box, G. E. P., 48, 192, 372, 388  
Branscum, A., 48, 49, 364, 427  
Bretz, F., 124  
Brockwell, P. J., 160  
Brownlee, K. A., 389  
Brown, M. B., 237

## C

Carlin, J. B., 48  
Casella, G., 73, 242  
Cavanaugh, J. E., 426  
Chen, G., 387  
Christensen, R., 16, 31, 48, 49, 87, 88, 124,  
127, 136, 145, 153, 154, 156, 178,

180–182, 192, 224, 225, 230, 233,  
237, 242, 251, 252, 276, 304, 314,  
328, 341, 356, 359, 362–364, 380,  
387–389, 396, 404, 420, 421, 427,  
437, 491, 499, 500, 503, 508  
Cochran, W. G., 136, 189, 242, 256, 276,  
371–373  
Cook, R. D., 146, 343, 372, 383, 385–388,  
420  
Cornell, J. A., 337, 338  
Cox, D. R., 35, 242, 256, 388  
Cox, G. M., 242, 276  
Cressie, N., 160

## D

Daniel, C., 146, 362, 389  
Davies, R. B., 482  
Davis, R. A., 160  
de Finetti, B., 48  
DeGroot, M. H., 48, 53  
deLaubenfels, R., 157  
Dey, D., 22  
Doob, J. L., 160  
Draper, N., 146, 150, 342, 343, 389  
Draper, N. R., 192, 406  
Duan, N., 355  
Dunsmore, I. R., 57, 192  
Dunson, D. B., 48  
Durbin, J., 367

## E

Eaton, M. L., 22  
Efron, B., 438

**F**

Feldt, L. S., 101–103  
 Ferguson, T. S., 287  
 Fienberg, S. E., 230, 233, 421  
 Fisher, R. A., 107, 125, 142, 145, 178, 242,  
 332, 500  
 Forzani, L., 420  
 Francia, R. S., 356  
 Fraser, D. A. S., 442  
 Freedman, D. A., 375, 376  
 Friedman, J., 438, 441, 443  
 Furnival, G. M., 421

**G**

Geisser, S., 48, 57, 156, 436  
 Gelman, A., 48  
 Gnanadesikan, R., 180  
 Goldstein, M., 410  
 Graybill, F. A., 22, 255, 276  
 Grizzle, J. E., 388  
 Groß, J., 282  
 Guttman, I., 57

**H**

Haberman, S. J., xvi  
 Hamada, M. S., 242  
 Hanson, T. E., 48, 49, 364, 427  
 Hartigan, J., 160  
 Harville, D. A., 22  
 Haslett, J., 387  
 Hastie, T., 428, 438, 441, 443  
 Hayes, K., 387  
 Hinkelmann, K., 242, 509  
 Hinkley, D. V., 35, 193  
 Hochberg, Y., 124  
 Hodges, J. S., 160  
 Hoerl, A. E., 405  
 Högfeldt, P., 500  
 Holt, D., 314  
 Hothorn, T., 124  
 Hsu, J. C., 124  
 Huber, P. J., 410  
 Hunter, J. S., 192  
 Hunter, W. G., 192  
 Hurvich, C. M., 426  
 Huynh, H., 101–103

**J**

James, G., 428, 438  
 Jeffreys, H., 48  
 John, P. W. M., 242

Johnson, D. E., 180  
 Johnson, R. A., 160, 404  
 Johnson, W., 48, 49, 362–364, 387, 427, 508

**K**

Kempthorne, O., 242, 509  
 Kennard, R., 405  
 Koch, G. G., 388  
 Kuh, E., 397, 398  
 Kutner, M. H., 136

**L**

LaMotte, L. R., 264  
 Lehmann, E. L., 35, 36, 70, 501  
 Lenth, R. V., 362  
 Lin, Y., 304  
 Lindley, D. V., 48  
 Li, W., 136

**M**

Mandansky, A., 372  
 Mandel, J., 276  
 Marquardt, D. W., 404, 405  
 Martin, R. J., 387  
 Mathew, T., 314  
 McCullagh, P., 16, 250, 252  
 McCulloch, C. E., 160  
 Miller, F. R., 180  
 Miller, R. G., Jr., 124  
 Milliken, G. A., 255, 276  
 Mitra, S. K., 303  
 Moguerza, J. M., 416  
 Monlezun, C. J., 314  
 Morrison, Donald F., 161  
 Mosteller, F., 407  
 Muñoz, A., 416

**N**

Nachtsheim, C. J., 136  
 Neill, J. W., 180  
 Nelder, J. A., 16  
 Neter, J., 136  
 Neuhaus, J. M., 160

**O**

Oehlert, G. W., 242  
 Ogden, R., 420

**P**

Pearson, L. M., 387  
 Peixoto, J. L., 63  
 Petkova, E., 420  
 Picard, R. R., 383  
 Pukelsheim, F., 286  
 Puri, M. L., 509

**R**

Raiffa, H., 48  
 Rao, C. R., 22, 35, 43, 101, 168, 276, 303,  
 388, 481  
 Ravishanker, N., 22  
 Reid, N., 242  
 Rencher, A. C., 22  
 Ripley, B. D., 160  
 Robert, C. P., 48  
 Ronchetti, E. M., 410  
 Rothman, A. J., 420  
 Rubin, D. B., 48  
 Ryan, T. A., Jr., 184

**S**

Savage, L. J., 48  
 Schaalje, G. B., 22  
 Schafer, D. W., 51  
 Schatzoff, M., 421  
 Scheffé, H., 22, 123, 126–135, 139–142,  
 148, 213, 224, 226, 232, 236  
 Schlaifer, R., 48  
 Schwarz, G., 166, 168, 427, 476  
 Searle, S. R., 22, 160, 286, 380  
 Seber, G. A. F., 22  
 Sen, P. K., 509  
 Shapiro, S. S., 356  
 Sherfey, B. W., 180  
 Shewhart, W. A., 127, 363  
 Shi, L., 387  
 Shillington, E. R., 180–182, 363  
 Shumway, R. H., 160  
 Sinha, B. K., 314  
 Skinner, C. J., 314  
 Smith, A. F. M., 155, 410  
 Smith, H., 146, 150, 342, 343, 389  
 Smith, T. M. F., 314  
 Snedecor, G. W., 136, 189, 256  
 Starmer, C. F., 388  
 Stefanski, L. A., 380

Stern, H. S., 48  
 St. Laurent, R. T., 276  
 Stoffer, D. S., 160  
 Sugiura, N., 426  
 Sulzberger, P. H., 277

**T**

Tamhane, A., 124  
 Tarpey, T., 420  
 Tiao, G. C., 48  
 Tibshirani, R., 428, 438, 441, 443  
 Tsai, C.-L., 426  
 Tsao, R., 421  
 Tukey, J. W., 123, 127, 135–139, 141–143,  
 276, 388, 407, 413, 435

**U**

Utts, J., 184, 185, 266, 276, 386

**V**

Van Nostrand, R. C., 406  
 Vehtari, A., 48  
 Velilla, S., 396

**W**

Watson, G. S., 367  
 Weisberg, S., 146, 343, 372, 373, 380, 385,  
 387, 388  
 Welsch, R. E., 397, 398  
 Wermuth, N., 237  
 Westfall, P., 124  
 Wichern, D. W., 160, 404  
 Wichura, M. J., 22  
 Wikle, C. K., 160  
 Wilk, M. B., 356  
 Williams, E. J., 277  
 Wilson, R. W., 421  
 Witten, D., 428, 438  
 Wood, F. S., 146, 389  
 Wu, C. F. J., 242

**Z**

Zellner, A., 48  
 Zhu, M., 416

# Subject Index

## Symbols

*C*

correction factor, 114

$C(A)$ , 448

$CN$ , 398

$C_p$ , 222, 424

$F$  distribution, 69, 98, 482

doubly noncentral, 183

$J$ , 473

$J_n$ , 473

$J_r^c$ , 473

$M$ , 464

$MSE$ , 32

$MSGrps$ , 115

$M_0$ , 464

$M_A$ , 464

$M_\alpha$ , 113

$P$  value, 495, 500

$R(\cdot)$ , 95

$RMS$ , 394, 422

$RSS$ , 394

$R^2$ , 165, 421

$SSE$ , 32

$SSGrps$ , 115

$SSLF$ , 177

$SSPE$ , 177

$SSR(\cdot)$ , 90

$SSReg$ , 150

$SSTot$ , 114

$SSTot - C$ , 114

$\underline{\perp}$ , 15

$\alpha$  level test, 70, 493, 500

$\chi^2$ , 481

$L^p$  distance, 412

$\varepsilon$  ill-defined, 395

$dfE$ , 32

$r(A)$ , 449

$t$  distribution, 42, 482

$t$  residual, 384

$ALM$ , vii

$PA$ , vii

$ppo$ , 463

## A

ACOVA, 255

ACOVA table, 262

Added variable plot, 263, 377

Adjusted  $R^2$ , 422

Adjusted treatment means, 272

Affine transformation, 27

AIC, 425

AIC corrected, 427

ALM-III, vii

Almost sure, 12

Almost surely consistent, 288

Alternating additive effects, 251

Analysis of covariance, 255

estimation, 256

for BIBs, 267

missing observations, 265

nonlinear model, 276

testing, 262

Analysis of covariance table, 262

Analysis of means, 124

Analysis of variance, 2

balanced incomplete blocks (BIBs), 267

definition, 2

multifactor, 197

one-way, 107

three-factor

balanced, 230

unbalanced, 232



- two-factor
  - balanced with interaction, 204
  - balanced with quantitative factors, 214
  - balanced without interaction, 197
  - proportional numbers, 217
  - unbalanced, 219
- Analysis of variance table, 115, 150, 212
- Angle between vectors, 471
- ANOVA, 2
- ANOVA table, 202, 212
- Assumptions, 341, 493, 501
- Asymptotic consistency, 355
- Asymptotic results, viii
  
- B**
- Backwards elimination, 431
- Balanced ANOVA, 197
- Balanced incomplete block design, 267
- Balanced three-way ANOVA, 230
- Balanced two-way ANOVA, 197, 204, 214
- Basis, 448
- Basis functions, 153
- Bayesian estimation, 405, 406
- Bayesian statistics, 48, 155, 491
- Bayes's theorem, 49
- Best predictor, 156, 160, 173
- Best subset regression, 420
  - variable limits, 420
- Best subset selection, 421
- BIB, 267
- BIC, 427
- Binomial distribution, 16, 388
- Biweight, 413
- BLP, 156, 160, 173
- Blocking, 241
- BLUE, 33, 39, 42, 174, 282, 300, 315, 319
- BLUP, 163, 172, 174, 343
- BSD, 135
- Bonferroni, 123, 135, 385
- Bootstrap, 438
- Box-Cox transformations, 388
  
- C**
- Calibration, 191
- Candidate model, 419
- Canonical form, regression in, 401, 405
- Cauchy-Schwarz inequality, 166
  - proof, 476
- Cell means model, 205, 223
- Centering data, 146, 152, 396
- Central
  - chi-squared distribution, 481
  - $F$  distribution, 482
  - $t$  distribution, 482
- Central distribution, 482
- Central limit theorem, viii
- Change of scale, 397
- Characteristic function, 7, 488
- Characteristic root, 459
- Chi-squared distribution, 11, 37, 69, 100, 292, 293, 316, 323, 324, 481
- Classification models, 2
- Cluster error, 324
- Cluster sampling, 314
- Coefficient of
  - determination, 165, 421
  - partial determination, 171
  - variation, 388
- Collinearity, 393
- Column space, 448
- Comparisons, 116
- Completely randomized design, 242
- Complete statistic, 35, 36
- Compound symmetry, 102
- Concomitant variable, 261
- Condition number, 398
- Confidence
  - bands, 148
  - ellipsoid, 97
  - interval, 37, 496
  - simultaneous, 132
  - region, 97
- Confirmatory data analysis, 435
- Consistent, 288, 355
- Constrained estimation, 82, 84, 104
- Constraint on, 81
- Constraints
  - estimable, 77, 84, 259
  - estimation under, 82, 84, 104
  - imposed by hypothesis, 81, 84
  - linear, 74
  - nonestimable, 24
  - nonidentifiable, 75, 76, 83, 84, 114, 120, 206
- Constructed variable tests, 388
- Contrasts
  - balanced incomplete blocks (BIBs), 271
  - one-way, 116
  - orthogonal, 94, 119
  - polynomial, 188, 215
  - two-way with interaction, 207, 214
  - two-way without interaction, 203
- Cook's distance, 386
- Corrected AIC, 427

Correction factor, 114, 150  
 Correlation coefficient, 166, 193  
   multiple, 166  
   partial, 170, 193, 429  
   serial, 363  
 Cost complexity pruning, 427  
 Counts, 388  
 Covariance, 5, 160  
   analysis of, 255  
 Covariance matrix, 5  
 Covariance parameterization, 489  
   linear, 490  
 Covariate, 261  
 CRD, 242  
 Critical region, 493  
 Cross-validation, 427

**D**

Degrees of freedom, 481  
   for error, 32  
 Deleted residual, 380  
 Design matrix, 1  
 Design space, 62  
 Determinant, 8, 463  
 Diagnostics, 341  
 Diagonal matrix, 457  
 Dispersion matrix, 5  
 Distance measures, 165, 345, 450, 471  
 Distributions  
   chi-squared, 11, 481  
   *F*, 482  
   doubly noncentral, 183  
   gamma, 49, 53  
   multivariate normal, 7  
   *t*, 482  
 Duncan's multiple range test, 124, 139  
 Dunnett's method, 124  
 Durbin-Watson test, 367

**E**

EDA, 435  
 Eigenvalue, 459  
 Eigenvector, 459  
 Empirical estimate, 375  
 Error degrees of freedom, 32  
 Error mean square, 32  
 Error rate, 123  
 Error space, 62  
 Estimable, 22, 24, 39  
 Estimable constraints, 77, 84, 259  
 Estimable part, 82

**Estimation**

Bayesian, 48, 405, 406  
 best linear unbiased (BLUE), 33, 39, 174  
 consistent linear unbiased (CLUE), 290  
 general Gauss-Markov models, 281  
 generalized least squares (GLS), 39, 281  
 generalized split plot, 322, 325  
 least squares, 28, 42, 300, 315, 319  
 maximum likelihood, 34, 39  
 minimum variance unbiased, 35, 39  
 ordinary least squares (OLS), 28, 42,  
   300, 315, 319  
 simple least squares, 28, 42, 300, 315,  
   319  
 unbiased, 27  
   for variance, 31  
 uniformly minimum variance unbiased  
   (UMVU), 35, 39  
 variance, unbiased, 31  
 weighted least squares, 42  
 with constraints, 82, 84, 104

**Estimation space, 62**

Expected mean squares, 115, 201, 206

Expected squared error, 156, 394

Expected values, 485

  quadratic forms, 11  
 random vectors and matrices, 4

Experimental unit, 241

Experimentwise error rate, 123, 124, 142

Exploratory data analysis, 435

Exponential regression, 16

**F**

Factor, 247

Factorial design, 247, 250

Factorial experiment, 247

Factorial treatment structure, 247, 250

Fieller's method, 97, 191

Fisherian testing, 70, 142, 491, 493, 499, 500

Fisher significant difference (FSD), 133

Fisher's "z" distribution, 500

Fitted values, 31, 369

  definition, 2

Forward selection, 429

Full model, 64, 68, 500

Fundamental theorem of least squares esti-  
 mation, 28

**G**

Gamma distribution, 49, 53

Gamma regression, 16

Gaussian distribution, 7

Gauss–Markov Theorem, 33  
 General Gauss–Markov  
   estimation, 281  
   testing, 291  
 Generalized additive models, 154, 186  
 Generalized inverse, 466  
   general form, 478  
 Generalized inverse regression, 404  
 Generalized least squares  
   estimation, 39  
   testing, 98  
 Generalized likelihood ratio test, 70, 73, 104  
 Generalized linear model, 16, 22, 27, 153,  
   157, 413  
 Generalized split plot models, 319  
 General linear model, 16  
 Graeco–Latin squares, 246  
 Gram–Schmidt orthogonalization, 91, 187,  
   199, 451, 460, 465, 470  
 Grand mean, 114  
 Greedy algorithm, 419

**H**

Hamming loss, 159  
 Heterogeneous variances, 342  
 Heteroscedastic, 342, 507  
 High leverage, 347  
 Homologous factors, 251  
 Homoscedastic, 504  
 HSD, 124, 135  
 Huber–White estimator, 376  
 Huynh–Feldt condition, 102

**I**

Idempotent matrix, 470  
 Identifiable, 22, 23, 489  
 Identity matrix, 457  
 Ill-conditioned, 398  
 Ill-conditioned model matrix, 395  
 i.i.d., 7  
 Ill-defined, 395  
 Incomplete blocks, 241, 267  
 Independence, 9, 13, 230  
   contingency table, 240  
   linear, 448  
   random vectors, 487  
 Independent identically distributed, 7  
 Influential observation, 342  
 Information criteria, 425  
 Inner product, 43, 164, 450, 471  
 Interaction  
   BIB designs, 275

  contrasts, 207  
   factorial treatment structure, 247  
   plot, 213  
   split plot designs, 329  
   test, 205  
   three-way ANOVA, 230, 232  
   two-way ANOVA, 207, 215, 222  
 Interval estimation, 37, 496  
 Intraclass correlation, 102, 315  
 Invariance, 70  
 Inverse matrix, 458

**J**

Joint distribution, 485

**K**

Kernel trick, 409  
 Kriging, 176  
 Kronecker product, 240, 250, 458, 473, 474

**L**

Lack of fit, 177, 204, 369, 499  
   near replicate tests, 180  
   partitioning tests, 182  
   residual analysis, 157, 377  
   traditional test, 178  
 Lasso, 421  
 Latin square design, 243  
 Least squares  
   consistent estimate, 296  
   estimate, 28, 319  
   generalized, 39, 98, 281  
   ordinary, 39, 42, 300, 315, 328  
   simple, 39, 42, 300, 315, 328  
   weighted, 42  
 Legendre polynomials, 190  
 Length, 165, 450, 471  
 Leverage, 344, 347  
 Likelihood function, 34, 40  
 Likelihood ratio test, 70, 73, 104  
 Linear combination, 62, 448  
 Linear constraint, 74  
 Linear covariance parameterization, 490  
 Linear dependence, 448  
 Linear estimable function, 24  
 Linear estimate, 27  
 Linear independence, 448  
 Linear model  
   standard, 1  
 Locally weighted scatterplot smoother, 153  
 Logistic regression, 16, 87, 159, 388

- Logit model, 16, 87, 388  
 Log-linear model, 16, 87, 388  
 Lowess, 153  
 LSD, 123, 133  
 LSE, 28
- M**
- Mahalanobis distance, 345  
 Mallows's  $C_p$ , 424  
 Marginal distribution, 486  
 Matrix
  - design, 1
  - diagonal, 457
  - generalized inverse, 466
  - idempotent, 470
  - identity, 457
  - inverse, 458
  - model, 1
  - nonnegative definite, 462
  - orthogonal, 461
  - orthonormal, 461
  - partitioned, 458
    - inverse, 479
  - positive definite, 462
  - projection, 470
    - oblique, 471
    - perpendicular, 463
  - square, 457
  - symmetric, 457
  - zero, 473
- Maximum likelihood estimates
  - generalized least squares models, 39
  - standard models, 34
- Mean squared error, 32
  - population, 156, 394
- Mean squared groups, 115  
 M-estimates, 412  
 Milliken and Graybill test, 276  
 Minimum variance unbiased estimate, 35  
 Missing data, 265  
 Mixed model, 313  
 MLEs, 34, 425  
 Model matrix, 1  
 Models, 1
  - analysis of covariance, 255
  - analysis of variance
    - balanced incomplete block (BIB), 267
    - multifactor, 230
    - one-way, 107
    - three-way, 230
    - two-way, 197, 204, 214, 217, 219
  - balanced incomplete block (BIB) design, 267
  - cell means, 205, 223
  - cluster sampling, 314
  - completely randomized design (CRD), 242
  - estimable constraints, 104
  - experimental design, 241
  - full, 64, 68
  - general Gauss–Markov, 281
  - generalized least squares, 38, 98
  - generalized split plot, 319
  - Graeco–Latin square, 246
  - Latin square, 243
  - randomization theory, 509
  - randomized complete block (RCB)
    - design, 242, 318, 328
    - reduced, 64, 68
    - split plot design, 328
    - subsampling, 332
- Model selection, 419  
 Multicollinearity, 393  
 Multifactor structures, 230  
 Multiple comparisons, 123, 124  
 Multiple correlation coefficient, 166  
 Multiple range method, 137, 139  
 Multiple regression, 148  
 Multiple testing, 124  
 Multivariate distribution, 485  
 Multivariate normal, 7
- N**
- Newman–Keuls multiple range test, 124, 137  
 Neyman–Pearson testing, 70, 142, 491, 499  
 Noncentral
  - chi-squared distribution, 481
  - $F$  distribution, 482
  - $t$  distribution, 482
- Noncentrality parameter, 69, 98, 481  
 Nonestimable, 24  
 Nonestimable constraints, 30  
 Nonidentifiable, 23, 24  
 Nonidentifiable constraints, 75, 76, 83, 84, 114, 120, 206  
 Nonnegative definite matrix, 462  
 Nonnormality, 342  
 Nonparametric methods, 509  
 Nonparametric regression, 153, 186  
 Nonsingular case, 258  
 Nonsingular covariance matrix, 38, 98, 281  
 Nonsingular distribution, 5

- Nonsingular matrix, 458  
 Normal distribution, 7  
 Normal equations, 43  
 Normality, test for, 356  
 Normal plot, 354  
 Normal score, 354  
 Null hypothesis, 491  
 Null model, 70, 491, 500  
 Null space, 458, 479
- O**
- Oblique projection, 478  
 Oblique projection operator, 471  
 Odds ratio, 240  
 Offset, 57, 72, 73, 84, 105, 424  
 OLS, 39  
 One sample, 38, 73  
 One-way analysis of variance (ANOVA), 107  
 Optimal allocation of  $x$  values, 193  
 Ordinary least squares, 39, 42, 300, 328  
 Ordinary residual, 6, 263, 342  
 Orthogonal, 165, 450, 471  
   basis, 450  
   complement, 452, 476  
   constraints, 90, 91  
   contrasts, 94, 119  
   distance regression, 413  
   matrix, 461  
   polynomials, 187, 215  
   projection, 165, 463, 471  
 Orthonormal  
   basis, 450, 451, 453, 460, 461, 465  
   matrix, 359, 461, 479  
 Outliers  
   in dependent variable, 380, 384  
   in the design space, 347, 384  
   in the estimation space, 347, 384  
 Overfitting, 437
- P**
- Parameter, 1, 22, 491  
 Parameterization, 22  
 Partial correlation coefficient, 170, 193, 429  
 Partial determination, coefficient of, 171  
 Partially identifiable, 23  
 Partitioned matrices, 458  
   inverse, 479  
 Partitioned model, 256  
 PC, 404  
 PCR, 404, 408  
 Penalized estimation, 405
- Percentile, 481  
 Perfect estimation, 304  
 Perpendicular, 165, 450, 471  
   projection, 157, 163  
   projection operator (ppo), 165, 344, 463, 471  
 Poisson distribution, 16, 388  
 Polynomial contrasts, 188, 215  
 Polynomial regression, 147, 186, 214  
 Positive definite matrix, 462  
 Power, 70, 497  
 Power transformations, 388  
 Predicted  $R^2$ , 423  
 Predicted residual, 350, 380  
 Predicted RESidual Sum of Squares (PRESS), 381  
 Predicted values, 31, 369  
 Prediction, 155, 172  
   best linear predictor (BLP), 160, 173  
   best linear unbiased predictor (BLUP), 163, 172, 174, 343  
   best predictor (BP), 156, 173  
 Prediction interval, 57  
 Predictor, 174  
 PRESS, 381  
 Principal component, 404  
 Principal component regression (PCR), 404, 408  
 Probability distribution, 485  
 Projection  
   oblique, 471  
   perpendicular, 157, 163, 165, 463  
 Projection operator, 470  
 Proportional numbers, 95, 217, 230  
 Proportions, 388  
 Pure error, 177, 203
- Q**
- q-q plot, 354  
 Quadratic forms, 11, 474  
   distribution, 12, 13  
   expectation, 11  
   independence, 14, 15  
 Quadratic formula, 415  
 Quantile, 354  
 Quantitative factors, 188, 214
- R**
- Random effects, 176  
 Randomization, 241  
 Randomization theory, 509

- Randomized complete block design, 241, 242, 318, 328
- Random matrix, 4
- Random vector, 5
- Range, 136
- Range space, 448
- Rank, 449
- Rankit plot, 354
- Rao's simple covariance structure, 43, 101
- RCB, 242
- Recovery of interblock information, 268
- Reduced model, 44, 64, 68, 500
- Reduction in sums of squares, 95
- Reference distribution, 491
- Reflexive generalized inverse, 467
- Regression analysis
  - definition, 2
  - in canonical form, 401
  - multiple regression, 148
  - nonparametric, 153
  - polynomial, 186, 214
  - simple linear regression, 146
- Regression model, 145, 420
- Regularization, 405
- Rejection region, 70, 482, 493
- Reparameterization, 36, 62, 63, 72, 73, 75, 83, 146, 151, 264
- Residual mean square, 394, 422
- Residual plots, 342
  - heteroscedasticity, 369
  - lack of fit, 377
  - normality, 354
  - serial correlation, 364, 380
- Residuals, 6, 31
  - definition, 2
  - deleted, 350, 380
  - predicted, 350, 380
  - standardized, 342
  - standardized deleted, 384
  - standardized predicted residuals, 384
  - Studentized, 342
- Residual sum of squares, 394
- Response surface, 192
- Ridge regression, 405, 410
- Ridge trace, 406
- Robust regression, 411
- Rotation, 479
- Row structure, 178
  
- S**
- Sample partial correlation coefficient, 171
- Sandwich estimator, 376
- Scaling the model matrix, 397
- Scheffé's method, 123, 128, 148
- Sequential fitting, 90, 95, 149
- Sequential sums of squares, 90, 95, 149
- Serial correlation, 363
  - test, 367
- Side conditions, 30, 75, 76, 83, 84, 114, 120, 206
  - estimation under, 82, 84, 104
- Significance testing, 70, 142, 491, 493, 499, 500
- Simple covariance structure, 43, 101
- Simple least squares, 39
- Simple linear regression, 146
- Simultaneous confidence intervals, 132
- Simultaneous inference, 123, 148
- Singular covariance matrix, 281
- Singular distribution, 5
- Singular value, 459
- Singular value decomposition, 401, 462
- Skew symmetric additive effects, 251
- Spanning set, 448
- Spanning space, 448
- Spatial data, 176
- Split plot designs, 313, 328
  - generalized, 319
- Split plot model, 103
- Square predictive correlation, 168
- Square matrix, 457
- Standard error, 491
- Standardized deleted residual, 384
- Standardized predicted residual, 384
- Standardized residuals, 342
- Standard linear model, 1, 6, 10
- Stepwise regression, 428
- Stochastically larger, 483
- Studentized range, 136
- Studentized residuals, 342, 384
- Student's  $t$ , 37, 50, 57, 482, 492
- Subplot analysis, 313, 328
- Subplot error, 323
- Subsampling, 332
- Subspace, 447
- Summation convention, 474
- Sum of squares
  - contrast, 118
  - error, 32
  - for regressing, 90
  - for testing  $H_0 : \lambda' \beta = 0$ , 81
  - groups, 115
  - reduction in, 95
  - regression, 150
  - total, 114

Supplemental observations, 261  
 Support vector machine, 413  
 Sweep operator, 264, 440  
 Symmetric additive effects, 251  
 Symmetric matrix, 457

## T

Tensor, 474  
 Tests, 61, 493  
 $\alpha$  level, 70, 493, 500  
 Durbin–Watson, 367  
 generalized likelihood ratio, 70, 73, 104  
 independence, 362  
 lack of fit, 177  
 Milliken and Graybill, 276  
 models  
   cluster sampling, 317  
   general Gauss–Markov, 291  
   generalized least squares, 98  
   generalized split plot, 326  
   standard, 64  
 multiple comparisons, 123  
 normality, 356  
 one-parameter, 491  
 parametric functions  
   cluster sampling models, 317  
   generalized least squares models, 100  
   generalized split plot models, 322, 325  
   standard models, 74  
 serial correlation, 367  
 single degree of freedom, 37, 89  
 Tukey’s one degree of freedom, 276  
 variances, 104  
 Wilk–Shapiro, 356  
 Test space, 91  
 Test statistic, 67, 70, 72, 77, 81, 84, 89, 91, 99, 100, 493, 499, 500  
 Three-way ANOVA, 230  
 Toeplitz matrix, 363  
 Tolerance, 399, 429  
 Tolerance point, 57  
 Topics in Design, 242, 268  
 Trace, 462  
 Transformations, 388  
   Box–Cox, 388  
   Grizzle, Starmer, Koch, 388  
   power, 388  
   variance stabilizing, 388  
 Transpose, 447, 457  
 Tukey’s biweight, 413

Tukey’s HSD, 124, 135  
 Tukey’s one degree of freedom, 276  
 Tukey’s one degree of freedom for nonadditivity, 388  
 Two independent samples, 37, 73  
 Two-phase linear regression, 192  
 Two-stage sampling, 314  
 Two-way ANOVA, 197, 204, 214, 217, 219

## U

UMVU, 39  
 Unbalanced ANOVA, 219, 232  
 Unbiased estimate, 27, 31, 33, 35, 282, 293  
 Unbiased predictor, 174  
 Unequal numbers, 219, 232  
 Uniformly minimum variance unbiased (UMVU), 35  
 Uniformly minimum variance unbiased (UMVU) estimate, 35  
 Uniformly most powerful invariant (UMPI) test, 70  
 Unreplicated experiments, 360  
 Updating formulae, 380  
 Usual constraints, 114  
 Usual multiple regression model, 149  
 Utts’s rainbow test, 184, 266, 386

## V

Variable selection, 419, 421, 428, 433  
 Variance component, 313  
 Variance-covariance matrix, 5  
 Variance estimation  
   Bayesian, 50, 55  
   general Gauss–Markov models, 292, 293  
   generalized least squares models, 41  
   standard models, 31, 35, 36  
 Variance inflation factor, 399  
 Variance stabilizing transformations, 388  
 Vec operator, 458, 473  
 Vector, 447  
   angle between, 471  
   space, 447  
 VIF, 399

## W

Weakest link pruning, 427  
 Weak experimentwise error rate, 123, 124  
 Weighted least squares, 42  
 Well defined parameterization, 489

Whole plot, [313](#)  
Whole-plot analysis, [318](#), [319](#), [325](#), [326](#), [328](#)  
Whole-plot error, [324](#)  
Wilk–Shapiro test, [356](#)  
WLS, [42](#)  
Working–Hotelling confidence bands, [148](#)

**Y**  
Youden squares, [275](#)

**Z**  
Zero matrix, [447](#), [473](#)