

# 3 CHAPTER

## THE LINEAR-REGRESSION MODEL

3.1 Classical Inference with the Linear-Regression Model: A Review	64
3.2 Pooling Two Samples	76
3.3 Bayesian Inference with the Linear-Regression Model	77
3.4 Multivariate Normal Sampling	85

This chapter reviews theories of inference with the linear-regression model. The first section summarizes the standard theorems of classical inference. The second section deals with pooling information from two samples. The third section describes Bayesian pooling of information from one sample with prior information. And the fourth section describes Bayesian inference about the mean vector and variance-covariance matrix of a multivariate normal distribution.

For more complete introductory material the reader is referred to classical treatments by Johnston (1973), Theil (1971), and Rao (1965), and to Bayesian treatments by Lindley (1965), Zellner (1971), and Box and Tiao (1973).

### 3.1 Classical Inference with the Linear-Regression Model: A Review

Theoretical descriptions often amount to nothing more than the statement that some hypothetical variable  $\eta$  might depend on some vector of hypothetical variables  $X$ . Empirical workers daringly translate this into a statement about observable variables such as

$$Y_t = \beta_1 X_{1t} + \beta_2 X_{2t} + \dots + \beta_k X_{kt} + u_t$$

where  $Y_t$  is the observable counterpart of  $\eta$ ,  $X_t = (X_{1t}, X_{2t}, \dots, X_{kt})$  is the observable counterpart of  $X_t$ , indexes a set of  $T$  observations on each variable  $t=1,$

$\dots, T$ , and where the unobservables  $\beta_1, \beta_2, \dots, \beta_k$  and  $u_t$  represent features of the precise relationship between  $\eta$  and  $X$  that are not disclosed by the theory. The functional dependence of  $Y_t$  on  $X_t$  is thereby assumed to be linear with an unknown parameter vector  $\beta' = (\beta_1, \beta_2, \dots, \beta_k)$ , and the theory is asserted to be incomplete,  $u_t \neq 0$ , or at least the linear approximation to it is incomplete.

The model may be written in matrix form as

$$Y = X\beta + u \tag{3.1}$$

where  $Y$  is a  $T \times 1$  vector of observables,  $X$  is a  $T \times k$  matrix of observables,  $\beta$  is a  $k \times 1$  vector of unobservables, and  $u$  is a  $T \times 1$  vector of unobservables. Traditionally,  $\beta$  is called a parameter vector and  $u$  a disturbance or error vector. By definition,  $u$  is all of those things that determine  $Y$ , excluding  $X\beta$ , and (3.1) is merely a tautological definition of  $u$ ,  $u \equiv Y - X\beta$ . To be more explicit, write  $Y$  as

$$Y = X\beta + Z\gamma$$

where  $Z\gamma$  is the part of  $Y$  left out of Equation (3.1), with  $Z$  a  $(T \times m)$  matrix of variables and  $\gamma$  an  $(m \times 1)$  vector of fixed effects, and where  $m$  need not be finite.

Substantive content is introduced into the analysis by assigning to  $u \equiv Z\gamma$  a frequency distribution, say, multivariate normal with mean zero and variance matrix  $\Sigma$ . It is thereby asserted that if the matrix  $X$  were held fixed and the matrix  $Z$  allowed to vary within the confines of some more-or-less well-defined experimental conditions, the vector  $Z\gamma$  would appear to have been drawn from a particular normal distribution. More importantly, it is assumed that there is no tendency for any of the  $T$  residual effects to exceed or fall short of zero on the average,  $E(Z\gamma|X) = 0$ . Or in the more familiar parlance, the left-out effects are assumed to be uncorrelated with the included effects.

The requirement  $E(Z\gamma|X) = 0$  is a crucial and quite unlikely assumption. To give a trivial example, suppose in a sample of individuals the dependent variable  $Y$  measures sunburn susceptibility and the explanatory variable  $X$  measures hair color. The correlation between these two variables may lead us erroneously to conclude that red hair causes sunburn, when in fact, genetic inheritance ( $Z$ ) determines both sunburn susceptibility and hair color.

The most foolish errors of inference derive from this kind of model misspecification; yet books about statistical inference hardly mention it. Of course, as discussed in Chapter 1, the choice of specification in this sense is not within the purview of the theory of statistical inference, which almost by definition takes the model as well defined. This enormously restricts the usefulness of statistical theory in nonexperimental research. It is sensible

sometimes during a data analysis to take the model as well specified; but it is senseless always to do so. Any inference will be thoroughly discredited by the identification of particular left-out variables. In devoting a whole chapter to post-data model construction, this book places relatively great emphasis on this topic. Even here, the heavy allocation of space to problems that presume the definition of the underlying causal model reflects what *can* be said about the various topics, not what *should* be said. More is said of this in Chapter 9, but until then the specification assumption  $E(\mathbf{u}|\mathbf{X}) = \mathbf{0}$  should be taken as given.

FINDING A COMPLETE THEORY

An initial question that may be asked of the data is whether the theory is, in fact, complete,  $\mathbf{u} = \mathbf{0}$ . In particular, is there a value of  $\beta$ , say,  $\mathbf{b}$ , that makes the theory appear to be complete,  $\mathbf{Y} - \mathbf{X}\mathbf{b} = \mathbf{0}$ ? Whenever the number of observations  $T$  exceeds the number of parameters  $k$ , such a value of  $\mathbf{b}$  is quite unlikely to exist. We can, however, find a value that makes the theory appear as complete as possible by making the residual vector  $\mathbf{e} = \mathbf{Y} - \mathbf{X}\mathbf{b}$  as small as possible. One (arbitrary) measure of smallness is the square of the length of  $\mathbf{e}$ ,  $\mathbf{e}'\mathbf{e}$ . Minimizing  $\mathbf{e}'\mathbf{e}$  by setting its derivatives to zero implies the equations

$$\mathbf{0} = \frac{\partial}{\partial \beta} (\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta) = -2\mathbf{X}'\mathbf{Y} + 2\mathbf{X}'\mathbf{X}\beta. \tag{3.2}$$

Solving this with the assumption that  $\mathbf{X}'\mathbf{X}$  is invertible yields the least squares value of  $\beta$

$$\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}.$$

THE MULTIPLE CORRELATION COEFFICIENT

Although  $\mathbf{e}'\mathbf{e}$  is one measure of the completeness of the theory, it is more common to use an  $R^2$  defined by

$$R^2 = \left[ \frac{\sum_t (Y_t - \bar{Y})^2 - \mathbf{e}'\mathbf{e}}{\sum_t (Y_t - \bar{Y})^2} \right], \tag{3.3}$$

where  $\bar{Y} = \sum Y_t / T$ . The  $R^2$  can be shown to lie between zero and one. A model with only a constant term  $Y_t = \beta_1 + u_t$  can be written in vector notation as  $\mathbf{Y} = \mathbf{1}_T\beta_1 + \mathbf{u}$ , where  $\mathbf{1}_T$  is a  $T \times 1$  vector of ones. The least-squares estimate of  $\beta_1$  is  $\hat{\beta}_1 = (\mathbf{1}_T'\mathbf{1}_T)^{-1}\mathbf{1}_T'\mathbf{Y} = \bar{Y}$ , where  $\bar{Y}$  is the mean of the observations  $Y_t$ ,  $t = 1, \dots, T$ . The corresponding residual sum of squares is  $(\mathbf{Y} - \mathbf{1}_T\hat{\beta}_1)'(\mathbf{Y} - \mathbf{1}_T\hat{\beta}_1) = \sum(Y_t - \bar{Y})^2$ . Assuming that the first column of  $\mathbf{X}$  is a vector of ones, it is always possible to make  $\mathbf{e}'\mathbf{e}$  equal to  $\sum(Y_t - \bar{Y})^2$  by setting  $\beta' = (\bar{Y}, 0, 0, \dots, 0)$ . Since we are choosing  $\mathbf{b}$  to make  $\mathbf{e}'\mathbf{e}$  as small as possible, it must be true that  $0 \leq \mathbf{e}'\mathbf{e} \leq \sum(Y_t - \bar{Y})^2$  and thus  $0 \leq R^2 \leq 1$ . An

$R^2$  of one means that the model is possibly complete,  $\mathbf{e}'\mathbf{e} = 0$ , whereas an  $R^2$  of zero means that the model is totally incapable of explaining the variability in the observed data.

ESTIMATING  $\beta$

The least squares vector  $\mathbf{b}$  has properties other than making the theory appear as complete as possible. A frequency distribution for  $\mathbf{u}$  with  $E\mathbf{u} = \mathbf{0}_T$ , and  $V\mathbf{u} = \sigma^2\mathbf{1}_T$  implies a frequency distribution for  $\mathbf{b}$ . ( $\mathbf{1}_T$  is the  $T \times T$  identity matrix.  $\mathbf{0}_T$  is a  $T \times 1$  vector of zeros.) Since  $\mathbf{b}$  is a linear function of  $\mathbf{u}$ ,

$$\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{X}\beta + \mathbf{u}) = \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u},$$

the moments of  $\mathbf{b}$  are straightforwardly calculated as

$$E\mathbf{b} = \beta + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E(\mathbf{u}) = \beta,$$

$$V\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'V(\mathbf{u})\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}.$$

It is important to understand that these moments describe the behavior of  $\mathbf{b}$  in repeated samples. This is, if the "experiment" is repeated over and over with fixed  $\mathbf{X}$  matrix, the average value of  $\mathbf{b}$  would be  $\beta$ . In that sense only, a particular value of  $\mathbf{b}$  is taken as indicating where  $\beta$  lies. These results are summarized as follows:

THEOREM 3.1 (MOMENTS OF LEAST-SQUARES ESTIMATOR). Assuming  $\mathbf{X}'\mathbf{X}$  is invertible,  $E\mathbf{u} = \mathbf{0}_T$  and  $V(\mathbf{u}) = \sigma^2\mathbf{1}_T$ , the least-squares estimator of  $\beta$  is  $\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$ , with mean  $\beta$  and variance  $\sigma^2(\mathbf{X}'\mathbf{X})^{-1}$ .

A desirable feature of the estimator  $\mathbf{b}$  from the standpoint of calculation is that it is a linear function of  $\mathbf{Y}$ . It is also an unbiased estimator of  $\beta$ ,  $E(\mathbf{b}) = \beta$ . It is of interest to determine if some other unbiased linear estimator might be better. Suppose we wished to estimate the scalar parameter  $\theta = \psi'\beta$  where  $\psi$  is a  $k$ -dimensional vector of constants. The least-squares estimator of  $\theta$  is taken to be  $\psi'\mathbf{b}$ . An alternative linear estimator is

$$\hat{\theta} = \mathbf{A}'\mathbf{Y} + a$$

where  $\mathbf{A}$  is a vector of constants, and  $a$  is a scalar constant. The estimator is unbiased only if

$$\psi'\beta = E\hat{\theta}, \quad \text{for all } \beta,$$

that is, only if

$$\psi'\beta = \mathbf{A}'E\mathbf{Y} + a = \mathbf{A}'\mathbf{X}\beta + a.$$

For this to be true for all  $\beta$  we must have

$$a = 0, \quad \psi' = A'X.$$

The variance of the estimator  $\hat{\theta}$  is

$$V\hat{\theta} = A'V(Y)A = \sigma^2 A' A.$$

Minimization of the variance of  $\hat{\theta}$  subject to the unbiasedness restriction is a simple constrained maximization problem. With  $\lambda$  as the vector of  $k$  Lagrange multipliers, this requires the derivatives of the function  $\sigma^2 A' A + 2\lambda'(\psi - X'A)$  to be set to zero; that is,

$$\begin{aligned} \mathbf{0} &= \frac{\partial f}{\partial A} = 2\sigma^2 A' - 2\lambda' X' \\ \mathbf{0} &= \frac{\partial f}{\partial \lambda} = \psi' - A' X \end{aligned} \tag{3.4}$$

Postmultiplying the first equation by  $X$  we obtain

$$\sigma^2 A' X - \lambda' X' X = \mathbf{0}$$

which implies

$$\lambda' = \sigma^2 A' X (X' X)^{-1} = \sigma^2 \psi' (X' X)^{-1}.$$

Substituting this last expression into (3.4) yields

$$\mathbf{0} = 2\sigma^2 A' - 2\sigma^2 \psi' (X' X)^{-1} X'$$

or

$$A' = \psi' (X' X)^{-1} X'.$$

Thus the minimum variance linear unbiased estimator of  $\psi/\beta$  is

$$\hat{\theta} = A' Y = \psi' (X' X)^{-1} X' Y,$$

which is, of course, just the least-squares estimator, and we have established the following result.

**THEOREM 3.2 (GAUSS-MARKOV).** *Assuming  $X'X$  is invertible,  $E(\mathbf{u}) = \mathbf{0}$  and  $V(\mathbf{u}) = \sigma^2 \mathbf{1}$  with  $\sigma^2$  finite, the least-squares estimator  $\psi'b$  of the linear combination of coefficients  $\psi/\beta$  has minimum variance among the class of unbiased linear estimators.*

If  $\mathbf{u}$  has a normal distribution, the least-squares estimator  $\mathbf{b}$  which is a linear function of  $\mathbf{u}$  is also normally distributed. Furthermore,  $\mathbf{b}$  is the maximum likelihood estimator of  $\beta$ . The distribution of  $\mathbf{Y}$  is

$$f(\mathbf{Y}|\mathbf{X}, \beta, \sigma^2) = (2\pi\sigma^2)^{-T/2} \exp \left[ -\frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{X}\beta)' (\mathbf{Y} - \mathbf{X}\beta) \right]$$

Using the fact that

$$\begin{aligned} (\mathbf{Y} - \mathbf{X}\beta)' (\mathbf{Y} - \mathbf{X}\beta) &= (\mathbf{Y} - \mathbf{X}\mathbf{b} + \mathbf{X}\mathbf{b} - \mathbf{X}\beta)' (\mathbf{Y} - \mathbf{X}\mathbf{b} + \mathbf{X}\mathbf{b} - \mathbf{X}\beta) \\ &= (\mathbf{Y} - \mathbf{X}\mathbf{b})' (\mathbf{Y} - \mathbf{X}\mathbf{b}) + (\beta - \mathbf{b})' X' X (\beta - \mathbf{b}) \end{aligned} \tag{3.5}$$

since  $X'(\mathbf{Y} - \mathbf{X}\mathbf{b}) = \mathbf{0}$ , we may write the likelihood function as

$$L(\beta, \sigma^2; \mathbf{Y}, \mathbf{X})$$

$$\propto (\sigma^2)^{-T/2} \exp \left[ -\frac{1}{2\sigma^2} (\mathbf{Y} - \mathbf{X}\mathbf{b})' (\mathbf{Y} - \mathbf{X}\mathbf{b}) \right] \exp \left[ -\frac{1}{2\sigma^2} (\beta - \mathbf{b})' X' X (\beta - \mathbf{b}) \right] \tag{3.6}$$

which attains its maximum at

$$\beta = \mathbf{b}, \quad \sigma^2 = \frac{(\mathbf{Y} - \mathbf{X}\mathbf{b})' (\mathbf{Y} - \mathbf{X}\mathbf{b})}{T} = \hat{\sigma}^2.$$

Summarizing these results in a theorem:

**THEOREM 3.3 (MAXIMUM LIKELIHOOD ESTIMATOR).** *Assuming that  $X'X$  is invertible and that  $\mathbf{u}$  is normally distributed with mean vector  $\mathbf{0}_T$  and variance matrix  $\sigma^2 \mathbf{1}_T$ , the least-squares estimator  $\mathbf{b}$  is the maximum likelihood estimator and is normally distributed with mean  $\beta$  and variance  $\sigma^2 (X'X)^{-1}$ .*

**ESTIMATING  $\sigma^2$**

The maximum likelihood estimator of  $\sigma^2$  is  $\hat{\sigma}^2 = (\mathbf{Y} - \mathbf{X}\mathbf{b})' (\mathbf{Y} - \mathbf{X}\mathbf{b}) / T$ . With  $M_x = \mathbf{I} - X(X'X)^{-1}X'$  this estimator can be written as  $\hat{\sigma}^2 = Y'M_x Y / T = uM_x u / T$ , since  $M_x X = \mathbf{0}$ . The expected value of  $\hat{\sigma}^2$  is  $E(\hat{\sigma}^2) = E(u'M_x u / T) = E(tr[u'M_x u] / T) = E(tr[M_x u u'] / T) = \sigma^2 tr[M_x] / T = \sigma^2 (T - tr[X(X'X)^{-1}X']) / T = \sigma^2 (T - k) / T$ . Therefore  $\hat{\sigma}^2$  is biased estimator,  $E(\hat{\sigma}^2) \neq \sigma^2$ . The unbiased estimator is

$$s^2 \equiv \frac{(\mathbf{Y} - \mathbf{X}\mathbf{b})' (\mathbf{Y} - \mathbf{X}\mathbf{b})}{T - k}.$$

**CONFIDENCE REGIONS AND HYPOTHESIS TESTS**

The construction of confidence intervals and hypothesis tests makes use of the following result.

**LEMMA.** *If the  $p \times 1$  vector  $\mathbf{z}$  is normally distributed with mean  $E(\mathbf{z})$  and variance  $V(\mathbf{z})$ , then the quantity  $[\mathbf{z} - E\mathbf{z}]' [V(\mathbf{z})]^{-1} [\mathbf{z} - E\mathbf{z}]$  is the sum of squares of  $p$  independent, standard, normal, random variables, which by definition has a chi-square distribution on  $p$  degrees of freedom.*

The proof of this lemma follows straightforwardly from the observation that there exists a matrix  $C$  such that  $CV(\mathbf{z})C'$  is the identity matrix, and therefore,  $Cz - CE(\mathbf{z})$  has mean vector zero and an identity variance matrix. Thus  $[Cz - CE(\mathbf{z})][Cz - CE(\mathbf{z})]' = [z - E(\mathbf{z})][Cz - E(\mathbf{z})]' = [z - E(\mathbf{z})]V^{-1}(\mathbf{z})[z - E(\mathbf{z})]$  is the sum of squares of  $p$  independent standard normal random variables.

Partitioning the matrices conformably,

$$\begin{aligned} \mathbf{b}' &= (\mathbf{b}'_I, \mathbf{b}'_J) \\ \boldsymbol{\beta}' &= (\boldsymbol{\beta}'_I, \boldsymbol{\beta}'_J) \\ (\mathbf{X}'\mathbf{X})^{-1} &= \begin{bmatrix} [(X'X)^{-1}]_{II} & [(X'X)^{-1}]_{IJ} \\ [(X'X)^{-1}]_{JI} & [(X'X)^{-1}]_{JJ} \end{bmatrix}, \end{aligned}$$

and making use of the fact that marginal distributions of multivariate normal random variables are themselves normal, we have the following consequence of this lemma.

**THEOREM 3.4 (THE CHI-SQUARE TEST STATISTIC).** *The scalar random variable*

$$[\mathbf{b}_I - \boldsymbol{\beta}_I]' [\sigma^2(\mathbf{X}'\mathbf{X})^{-1}]^{-1}_{II} [\mathbf{b}_I - \boldsymbol{\beta}_I] = \chi^2 \quad (3.6)$$

*has a chi-square distribution with  $p$  degrees of freedom, where  $p$  is the dimension of  $\mathbf{b}_I$ .*

Thus if  $\chi^2_\alpha(p)$  is the upper  $\alpha$  percentage point of the chi-square distribution with  $p$  degrees of freedom, then the ellipsoid<sup>1</sup>

$$[\mathbf{b}_I - \boldsymbol{\beta}_I]' [\sigma^2(\mathbf{X}'\mathbf{X})^{-1}]^{-1}_{II} [\mathbf{b}_I - \boldsymbol{\beta}_I] \leq \chi^2_\alpha(p) \quad (3.7)$$

is a  $1 - \alpha\%$  confidence region for  $\boldsymbol{\beta}$ . Similarly, an  $\alpha$ -level test of the point null hypothesis  $\boldsymbol{\beta}_I = \boldsymbol{\beta}_I^*$  versus the alternative  $\boldsymbol{\beta}_I \neq \boldsymbol{\beta}_I^*$ , would reject the null hypothesis if  $\boldsymbol{\beta}_I = \boldsymbol{\beta}_I^*$  were outside this region.

If  $\sigma^2$  is unknown, these statements, although still true, lose their usefulness, since the regions described are functions of  $\sigma^2$ . It seems natural to use some estimate in place of the unknown  $\sigma^2$ . Remarkably enough, for the normal linear regression model this is approximately correct. To show that we need the following result.

<sup>1</sup>Actually, this is an ellipsoidal cylinder with  $\boldsymbol{\beta}_J$  unconstrained.

**THEOREM 3.5.** *The quantity  $(Y - X\mathbf{b})'(Y - X\mathbf{b})/\sigma^2 = (T - k)s^2/\sigma^2$  has a chi-square distribution with  $T - k$  degrees of freedom and is independent of  $\boldsymbol{\beta} - \mathbf{b}$ .*

The residual sum of squares  $(Y - X\mathbf{b})'(Y - X\mathbf{b})$  can be written as  $\mathbf{Y}'\mathbf{M}_x\mathbf{Y}$  where  $\mathbf{M}_x = \mathbf{I}_T - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ . Making use of  $\mathbf{M}_x\mathbf{X} = \mathbf{0}$ , the residual sum of squares can also be written

$$(\mathbf{Y} - \mathbf{X}\mathbf{b})'(Y - X\mathbf{b}) = \mathbf{u}'\mathbf{M}_x\mathbf{u}.$$

The idempotent matrix  $\mathbf{M}_x$  is shown in Appendix 1 to have  $k$  characteristic values equal to zero and the remaining  $T - k$  equal to one. Thus  $\mathbf{u}'\mathbf{M}_x\mathbf{u}$  can be written as  $\mathbf{u}'C\Lambda C\mathbf{u} = \mathbf{u}^*\Lambda\mathbf{u}^*$  where  $\mathbf{u}^*$  is normally distributed with mean zero and variance matrix  $\sigma^2\mathbf{I}_T$ , and where  $\Lambda$  is a diagonal matrix with  $k$  zero diagonal elements and  $T - k$  elements equal to one. We have thus written  $(Y - X\mathbf{b})'(Y - X\mathbf{b})/\sigma^2$  as the sum of squares of  $T - k$  independent, standard, normal, random variables, which by definition is distributed chi-square with  $T - k$  degrees of freedom.

The independence of the two random variables in the theorem follows from the independence of  $\mathbf{Y} - \mathbf{X}\mathbf{b} = \mathbf{M}_x\mathbf{u}$  and  $(\boldsymbol{\beta} - \mathbf{b}) = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}$ . Given the normality assumption, these quantities may be shown to be independent by computing their covariance  $E[\mathbf{M}_x\mathbf{u}][(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{u}]' = \sigma^2\mathbf{M}_x\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} = \mathbf{0}$ .

**Definition.** If the random variables  $\chi^2_1$  and  $\chi^2_2$  are independent and have chi-square distributions with  $p_1$  and  $p_2$  degrees of freedom, respectively, then the ratio

$$F = \frac{(\chi^2_1/p_1)}{(\chi^2_2/p_2)}$$

has an  $F$ -distribution with degrees of freedom parameters  $p_1$  and  $p_2$ .

**THEOREM 3.6 (THE F-TEST STATISTIC).** *Assuming that  $\mathbf{X}'\mathbf{X}$  is invertible and that  $\mathbf{u}$  is normally distributed with mean vector  $\mathbf{0}_T$  and variance matrix  $\sigma^2\mathbf{I}_T$ , the quantity*

$$\frac{[\mathbf{b}_I - \boldsymbol{\beta}_I]'[(\mathbf{X}'\mathbf{X})^{-1}]^{-1}_{II}[\mathbf{b}_I - \boldsymbol{\beta}_I]}{ps^2} = F \quad (3.8)$$

*has an  $F$  distribution with  $p$  and  $T - k$  degrees of freedom, where  $p$  is the dimension of  $\mathbf{b}_I$ .*

This is a direct consequence of Theorems 3.4 and 3.5 and should be compared with Theorem 3.4, which applies when  $\sigma^2$  is known. The statement implies  $1 - \alpha\%$  confidence ellipsoids of the form

$$[\mathbf{b}_T - \beta_T]'[(\mathbf{X}'\mathbf{X})^{-1}]_{TT}^{-1}[\mathbf{b}_T - \beta_T] \leq ps^2 F_{\alpha}(p, T - k)$$

where  $p$  is the dimension of  $\mathbf{b}_T$ ;  $p \leq k$ , and  $F_{\alpha}(p, T - k)$  is the upper  $\alpha$  percentage point of an  $F$  distribution with  $p$  and  $T - k$  degrees of freedom. In the special case when  $p = 1$ ,  $F^{1/2}$  is said to have a Student's  $t$  distribution, and this ellipsoid becomes a confidence interval for a particular coefficient

$$|b_i - \beta_i| \leq s[(\mathbf{X}'\mathbf{X})^{-1}]_{ii}^{-1} t_{\alpha}(T - k).$$

Furthermore, an  $\alpha$ -level test of the hypothesis  $\beta_T = \beta_T^*$  versus the alternative  $\beta_T \neq \beta_T^*$  would reject the hypothesis if  $\beta_T^*$  were not within the confidence region.

The quadratic form (3.3) can be written in another informative way. We may partition  $\mathbf{X}'\mathbf{X}$  as above,

$$\mathbf{X}'\mathbf{X} = \begin{bmatrix} \mathbf{X}'_T \mathbf{X}_T & \mathbf{X}'_T \mathbf{X}_J \\ \mathbf{X}'_J \mathbf{X}_T & \mathbf{X}'_J \mathbf{X}_J \end{bmatrix}$$

and, by the partitioned inverse rule, we have

$$[(\mathbf{X}'\mathbf{X})^{-1}]_{TT}^{-1} = \mathbf{X}'_T \mathbf{X}_T - \mathbf{X}'_T \mathbf{X}_J (\mathbf{X}'_J \mathbf{X}_J)^{-1} \mathbf{X}'_J \mathbf{X}_T.$$

Let the minimum error sum of squares be

$$ESS = (\mathbf{Y} - \mathbf{X}\mathbf{b})'(\mathbf{Y} - \mathbf{X}\mathbf{b}),$$

and the constrained minimum error sum of squares with constraint  $\beta_T = \beta_T^*$  be

$$ESS_0 = (\mathbf{Y} - \mathbf{X}_J \beta_T^* - \mathbf{X}_J \mathbf{b}_J^*)'(\mathbf{Y} - \mathbf{X}_J \beta_T^* - \mathbf{X}_J \mathbf{b}_J^*)$$

with  $\mathbf{b}_J^* = (\mathbf{X}'_J \mathbf{X}_J)^{-1} \mathbf{X}'_J (\mathbf{Y} - \mathbf{X}_J \beta_T^*)$ . By some tedious manipulations we may show that

$$(\beta_T^* - \mathbf{b}_T)'[(\mathbf{X}'\mathbf{X})^{-1}]_{TT}^{-1}(\beta_T^* - \mathbf{b}_T) = ESS_0 - ESS \quad (3.9)$$

and therefore the numerator of the  $F$  statistic is just the increase in the error sum of squares implied by the restriction  $\beta_T = \beta_T^*$

$$F = \frac{ESS_0 - ESS}{ps^2} \quad (3.10)$$

Furthermore, using the definition of  $R^2$  given by Equation (3.3) and letting  $R^2$  correspond to the unconstrained regression and  $R_0^2$  correspond to the

constrained regression, the  $F$  statistic may also be written as

$$F = \left( \frac{R^2 - R_0^2}{1 - R^2} \right) \left( \frac{T - k}{p} \right) \quad (3.11)$$

CONSTRAINED ESTIMATION

The consequences of constraints other than  $\beta_T = \beta_T^*$  are easy to compute. Minimization of the sum of squares  $(\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta)$  subject to restriction  $\mathbf{R}\beta = \mathbf{r}$  is a Lagrangian problem requiring the derivatives of the following expression to be set to zero.

$$(\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta) - 2\lambda(\mathbf{R}\beta - \mathbf{r}).$$

That is

$$\mathbf{R}\beta - \mathbf{r} = 0 \quad (3.12)$$

$$\mathbf{X}'(\mathbf{Y} - \mathbf{X}\beta) - \mathbf{R}'\lambda = 0. \quad (3.13)$$

Premultiplying (3.13) by  $\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}$  yields

$$\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y} - \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} (\mathbf{X}'\mathbf{X})\beta = \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}'\lambda$$

and using (3.12)

$$\lambda = (\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}')^{-1} (\mathbf{R}\mathbf{b} - \mathbf{r}),$$

which can be reinserted into (3.13) to yield, finally, the constrained least-squares value

$$\begin{aligned} \hat{\beta}(\mathbf{R}, \mathbf{r}) &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y} - (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}'(\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}')^{-1} (\mathbf{R}\mathbf{b} - \mathbf{r}) \\ &= \mathbf{b} - (\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}'(\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}')^{-1} (\mathbf{R}\mathbf{b} - \mathbf{r}). \end{aligned} \quad (3.14)$$

Using this value of  $\hat{\beta}$  in the sum of squares expression we have

$$(\mathbf{Y} - \mathbf{X}\hat{\beta})'(\mathbf{Y} - \mathbf{X}\hat{\beta}) = (\mathbf{Y} - \mathbf{X}\mathbf{b})'(\mathbf{Y} - \mathbf{X}\mathbf{b}) + (\mathbf{R}\mathbf{b} - \mathbf{r})'(\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{R}')^{-1} (\mathbf{R}\mathbf{b} - \mathbf{r}) \quad (3.15)$$

where the last term in the last line is, therefore, the increase in the error sum of squares implied by the restriction  $\mathbf{R}\beta = \mathbf{r}$ . By inspection, this establishes the validity of Equation (3.9) with  $\mathbf{r} = \beta_T^*$  and  $\mathbf{R} = (\mathbf{I}_p \mathbf{0})$ . Computation of the mean and variance of  $\hat{\beta}$  is left as an exercise.

TREATMENT OF THE CONSTANT TERM

A regression function ordinarily has a constant as one of the "variables." Equivalently, one of the columns of the  $\mathbf{X}$  matrix is a vector of ones. Although everything that has been said still applies, it is informative to treat the constant somewhat differently.

Let us remove from the  $\mathbf{X}$  matrix the vector of ones and write the regression as

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{1}_T\alpha + \mathbf{u}$$

where  $\mathbf{1}_T$  is a vector of ones and  $\alpha$  is the scalar constant. The error sum of squares may be minimized with respect to  $\alpha$  as

$$\begin{aligned} 0 &= \frac{\partial}{\partial \alpha} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{1}_T\alpha)'(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{1}_T\alpha) \\ &= -2\mathbf{1}'_T(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) - 2\mathbf{1}'_T\mathbf{1}_T\alpha \end{aligned}$$

thus as a function of  $\boldsymbol{\beta}$ , the least-squares estimate of  $\alpha$  is

$$a(\boldsymbol{\beta}) = \mathbf{1}'_T(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})(\mathbf{1}'_T\mathbf{1}_T)^{-1}$$

inserting this value into the error sum of squares yields

$$\begin{aligned} &(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{1}_T\mathbf{1}'_T[\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}]\mathbf{T}^{-1})(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{1}_T\mathbf{1}'_T[\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}]\mathbf{T}^{-1}) \\ &= (\mathbf{M}\mathbf{Y} - \mathbf{M}\mathbf{X}\boldsymbol{\beta})'(\mathbf{M}\mathbf{Y} - \mathbf{M}\mathbf{X}\boldsymbol{\beta}) \end{aligned}$$

where  $\mathbf{M} = \mathbf{I}_T - \mathbf{1}_T(\mathbf{1}'_T\mathbf{1}_T)^{-1}\mathbf{1}'_T$ . The least-squares estimate of  $\boldsymbol{\beta}$  can, therefore, be computed by transforming the observations  $(\mathbf{Y}, \mathbf{X})$  by the matrix  $\mathbf{M}$  and minimizing the transformed error sum of squares. It is straightforward to show that this yields

$$\mathbf{b} = (\mathbf{X}'\mathbf{M}'\mathbf{M}\mathbf{X})^{-1}(\mathbf{X}'\mathbf{M}'\mathbf{M}\mathbf{Y})$$

But now notice the effect of the matrix  $\mathbf{M}$  on a vector  $\mathbf{Y}$ :

$$\mathbf{M}\mathbf{Y} = \mathbf{Y} - \mathbf{1}_T(\mathbf{1}'_T\mathbf{Y})\mathbf{T}^{-1} = \mathbf{Y} - \mathbf{1}_T\bar{Y}$$

where  $\bar{Y}$  is the average value of the observed variables. In words, the matrix  $\mathbf{M}$  subtracts out the mean of a variable, and the least-squares estimate of the slope coefficients can be computed by first subtracting the means from all the variables and then computing the least-squares value in the usual way.

Parenthetically, if the regression process is written  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma} + \mathbf{u}$  where  $\mathbf{X}$  and  $\mathbf{Z}$  are observable matrices, the least-squares estimate of  $\boldsymbol{\beta}$  is  $\mathbf{b} = (\mathbf{X}'\mathbf{M}'_2\mathbf{M}_2\mathbf{X})^{-1}\mathbf{X}'\mathbf{M}'_2\mathbf{M}_2\mathbf{Y} = (\mathbf{X}'\mathbf{M}'_2\mathbf{X})^{-1}\mathbf{X}'\mathbf{M}_2\mathbf{Y}$  where  $\mathbf{M}_2 = \mathbf{I}_T - \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'$ . The residual vector formed from the regression of  $\mathbf{Y}$  on  $\mathbf{Z}$  is  $\mathbf{Y} - \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{Y} = \mathbf{M}_2\mathbf{Y}$  and the residual matrix formed from the regression of each of the columns of  $\mathbf{X}$  on  $\mathbf{Z}$  is  $\mathbf{X} - \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{X} = \mathbf{M}_2\mathbf{X}$ . The least-squares estimate  $\mathbf{b}$  is then just the regression of the residual vector  $\mathbf{M}_2\mathbf{Y}$  on the residual matrix  $\mathbf{M}_2\mathbf{X}$ .

#### THE ADJUSTED $R^2$ AND REGRESSION SELECTION STRATEGIES

The  $R^2$  given by (3.3) has been suggested as a measure of the "completeness" of a theory. For that function it has the defect that if there are as many independent variables as there are observations, the  $R^2$  takes on the

value one. Moreover, as variables are added to the equation, the  $R^2$  necessarily increases. This creates an unwarranted preference for models with many explanatory variables, and some adjustment for the number of explanatory variables seems desirable. An analogous problem occurs when estimating  $\sigma^2$ . The maximum likelihood estimator of  $\sigma^2$  is biased upward, and the unbiased estimator,  $s^2 = ESS/(T - k)$ , is usually preferred.<sup>2</sup> If we think of  $1 - R^2$  as estimating  $\sigma^2/\sigma_y^2$ , it is by the same logic natural to define the adjusted  $R^2$  by the equation

$$1 - \bar{R}^2 \equiv \frac{s^2}{s_y^2} = \frac{ESS/(T - k)}{\Sigma(y_i - \bar{y})^2/T - 1} \quad (3.16)$$

The relationship between  $R^2$  and  $\bar{R}^2$  is then

$$1 - \bar{R}^2 = \frac{(T - 1)}{(T - k)}(1 - R^2) \quad (3.17)$$

No book on specification searches would be complete without the following two results. The first, due to Theil (1971, p. 543), has been used to justify search strategies that maximize the  $\bar{R}^2$ . The second defines a simple algorithm for increasing the  $\bar{R}^2$ .

**THEOREM 3.7 (EXPECTED  $\bar{R}^2$ ).** *Given two normal regression models, one of which is assumed to be true, the expected value of  $s^2$  for the true model is less than or equal to the expected value of  $s^2$  for the other model.*

To prove this, suppose that  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$  is the true model and  $\mathbf{Y} = \mathbf{Z}\boldsymbol{\gamma} + \mathbf{u}$  is the alternative, where  $\mathbf{X}$  and  $\mathbf{Z}$  are  $T \times k_x$  and  $T \times k_z$  matrices. Then

$$s_x^2 = (T - k_x)^{-1}\mathbf{Y}'\mathbf{M}_x\mathbf{Y}, \quad s_z^2 = (T - k_z)^{-1}\mathbf{Y}'\mathbf{M}_z\mathbf{Y}$$

where

$$\mathbf{M}_x = \mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \quad \text{and} \quad \mathbf{M}_z = \mathbf{I} - \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'$$

The expected error sum of squares for the false model is

$$\begin{aligned} E(\mathbf{Y}'\mathbf{M}_z\mathbf{Y}|\boldsymbol{\beta}) &= E((\mathbf{X}\boldsymbol{\beta} + \mathbf{u})'\mathbf{M}_z(\mathbf{X}\boldsymbol{\beta} + \mathbf{u})|\boldsymbol{\beta}) \\ &= \boldsymbol{\beta}'\mathbf{X}'\mathbf{M}_z\mathbf{X}\boldsymbol{\beta} + E(\mathbf{u}'\mathbf{M}_z\mathbf{u}|\boldsymbol{\beta}) = \boldsymbol{\beta}'\mathbf{X}'\mathbf{M}_z\mathbf{X}\boldsymbol{\beta} + (T - k_z)\sigma^2 \\ &\geq (T - k_z)\sigma^2 \end{aligned}$$

Thus  $E(\mathbf{Y}'\mathbf{M}_z\mathbf{Y}|\boldsymbol{\beta})/(T - k_z) \geq \sigma^2 = E(\mathbf{Y}'\mathbf{M}_x\mathbf{Y}|\boldsymbol{\beta})/(T - k_x)$ .

Proof of the next result is left to the reader.

**THEOREM 3.8 (INCREASING THE  $\bar{R}^2$ ).** *Omitting a variable from an equation will increase the  $\bar{R}^2$  if and only if the square of the  $t$ -statistic (the  $F$ ) for that coefficient is less than one.*

<sup>2</sup>It is also preferred by this Bayesian, but not of course because of the bias. See section 7.1.

3.2 Pooling Two Samples

Bayesian inference is different from classical inference in that it makes use of information that is not contained in the sample under study. Bayesian theory is concerned with the optimal pooling of sample information with nonsample information. Classical theory can also be used to pool information from more than one source, provided the information is generated by experiments that admit a frequency interpretation.

Suppose, in particular, that a given data set was arbitrarily split into two parts. Then the regression equation could be written as

$$\begin{bmatrix} Y^* \\ Y \end{bmatrix} = \begin{bmatrix} X^* \\ X \end{bmatrix} \beta + \begin{bmatrix} u^* \\ u \end{bmatrix}$$

where the \* indicates the first subset of the data. The usual least-squares estimate of  $\beta$  would be

$$\begin{aligned} \mathbf{b}^{**} &= \left( \begin{bmatrix} X^* \\ X \end{bmatrix} \begin{bmatrix} X^* \\ X \end{bmatrix} \right)^{-1} \begin{bmatrix} X^* \\ X \end{bmatrix} \begin{bmatrix} Y^* \\ Y \end{bmatrix} = (X^*X^* + XX)^{-1}(X^*Y^* + XY) \\ &= (X^*X^* + XX)^{-1}(X^*X^*b^* + XXb) \end{aligned} \tag{3.18}$$

where  $\mathbf{b}^* = (X^*X^*)^{-1}X^*Y^*$  and  $\mathbf{b} = (XX)^{-1}XY$ . In words, the least-squares estimate of  $\beta$  is a matrix-weighted average of the pair of estimates computed from each of two subsets of the data,  $\mathbf{b}^*$  and  $\mathbf{b}$ .

Another possibility is that the variance of the residuals in the first part of the data set is different from the variance in the second. Indicating these variances by  $\sigma^{*2}$  and  $\sigma^2$ , respectively, the model may be transformed to make the residual covariance matrix equal the identity matrix by dividing  $Y^*$  and  $X$  by  $\sigma$  and  $Y^*$  and  $X^*$  by  $\sigma^*$ . The resulting estimate of  $\beta$  is

$$\mathbf{b}^{**} = (\sigma^{*2}X^*X^* + \sigma^2XX)^{-1}(\sigma^{*2}X^*X^*b^* + \sigma^2XY) \tag{3.19}$$

The notation just used may seem confusing, but it is chosen to anticipate the Bayesian analysis of the next section, in which personal prior opinions end up being equivalent to the preliminary observations  $Y^*, X^*$ . The prior information is pooled with the current information exactly in accordance with formula (3.19).<sup>3</sup>

The variances  $\sigma^2$  and  $\sigma^{*2}$  may also be estimated. The likelihood function implied by the model is

$$L(\beta, \sigma^2, \sigma^{*2}; Y, Y^*, X, X^*) \propto (\sigma^2)^{-T/2} \exp \left[ -\frac{1}{2\sigma^2} (Y - X\beta)(Y - X\beta) \right] (\sigma^{*2})^{-T^*/2} \exp \left[ -\frac{1}{2\sigma^{*2}} (Y^* - X^*\beta)(Y^* - X^*\beta) \right]$$

<sup>3</sup>A schizophrenic attitude toward probabilities allows Theil and Goldberger (1961) to introduce prior information into classical inference in this way. An earlier paper of the same genre is Durbin (1953).

The reader may verify that maximizing this likelihood function implies choosing  $\beta$  in accord with Equation (3.19) and choosing the variances as

$$\hat{\sigma}^2 = \frac{(Y - Xb^{**})(Y - Xb^{**})}{T} \tag{3.20}$$

$$\hat{\sigma}^{*2} = \frac{(Y^* - X^*b^{**})(Y^* - X^*b^{**})}{T^*} \tag{3.21}$$

3.3 Bayesian Inference with the Linear-Regression Model

Bayesian inference with the linear-regression model is now introduced and shown to be formally equivalent to the pooling problem just discussed. Bayes' rule applied to the uncertain parameters  $(\beta, \sigma^2)$  is

$$f(\beta, \sigma^2 | Y, X) = \frac{f(Y, X | \beta, \sigma^2) f(\beta, \sigma^2)}{f(Y, X)}$$

where  $Y$  and  $X$  are observable data. The linear-regression model  $Y = X\beta + u$ , where  $u$  is normally distributed with mean  $0$  and variance  $\sigma^2I$ , implies a conditional distribution for  $Y$ :  $f_N(Y | X\beta, \sigma^2I)$ , indicating that  $Y$  is normally distributed with mean  $X\beta$  and variance  $\sigma^2I$ . If, furthermore,  $X$  is distributed independently of  $\beta$  and  $\sigma^2$ , the posterior distribution can be written as

$$\begin{aligned} f(\beta, \sigma^2 | Y, X) &= \frac{f_N(Y | X\beta, \sigma^2I) f(X) f(\beta, \sigma^2)}{f(Y | X) f(X)} \\ &= \frac{f_N(Y | X\beta, \sigma^2I) f(\beta, \sigma^2)}{f(Y | X)} \end{aligned}$$

Before discussing the choice of prior,  $f(\beta, \sigma^2)$ , a few things may be said about the immediately preceding formula. It implies that randomness in the matrix  $X$  is irrelevant for inference if  $X$  is distributed independently of  $\beta$  and  $\sigma^2$ . Classical inference, in contrast, cannot easily ignore the randomness in  $X$ , since sampling properties are necessarily affected. Another important point is that both the vectors of unobservables,  $\beta$  and  $u$ , are assigned personal probability distributions. Classically, the residual vector  $u$  is thought to have a frequency distribution and the parameter vector  $\beta$  is thought to be a fixed vector of constants. But from the Bayesian point of view, your personal opinion about both  $\beta$  and  $u$  is described in probabilistic terms, and the extreme distinction between  $\beta$  and  $u$  is regarded as unwarranted.

Having assumed that the residual vector is normally distributed, only the choice of a prior distribution for  $(\beta, \sigma^2)$  remains. Two alternatives are suggested here. The resulting posterior distributions are described in the following pair of theorems, and further discussion follows. Both results

ake use of the assumption that  $Y$  given  $X$ ,  $\beta$  and  $\sigma^2$  is normal with mean  $\beta$  and variance  $\sigma^2 I$ .

CONJUGATE PRIOR

THEOREM 3.9. Given the normal linear-regression model and a normal-gamma prior distribution for  $(\beta, \sigma^{-2})$ :

$$f(\beta, \sigma^{-2}) = f_N(\beta | b^*, \sigma^2(N^*)^{-1}) f_\gamma(\sigma^{-2} | s^{*2}, \nu^*),$$

which indicates that  $\beta$  given  $\sigma^2$  is distributed normally with mean  $b^*$  and variance  $\sigma^2(N^*)^{-1}$  and that  $\sigma^{-2}$  is distributed gamma with parameters  $s^{*2}$  and  $\nu^*$ , then the posterior distribution of  $(\beta, \sigma^{-2})$  is in the normal-gamma family with parameters

$$b^{**} = (N^* + XX')^{-1}(N^*b^* + XX'b) \tag{3.22}$$

$$N^{**} = N^* + XX' \tag{3.23}$$

$$\nu^{**} = \nu^* + T \tag{3.24}$$

$$(s^{**})^2 = [\nu^{**}]^{-1} [\nu^* s^{*2} + ESS + (b - b^*)(N^*(XX + N^*))^{-1} XX(b - b^*)]. \tag{3.25}$$

The marginal posterior distribution of  $\beta$  is multivariate Student with parameters  $b^{**}$ ,  $(s^{**})^2(N^{**})^{-1}$ , and  $\nu^{**}$ .

roof: Let  $h = \sigma^{-2}$  and use the likelihood function (3.6) together with the prior to obtain

$$f(\beta, h | Y, X) \propto h^{T/2} \exp[-\frac{1}{2}hESS] \exp[-\frac{1}{2}h(\beta - b)(\beta - b)] h^{k/2} \exp[-\frac{1}{2}h(\beta - b^*)N^*(\beta - b^*)] h^{\frac{\nu^* - 1}{2}} \exp[-\frac{1}{2}\nu^* s^{*2}h].$$

sing result (T10) in Appendix 1, this can be rewritten as

$$f(\beta, h | Y, X) \propto h^{k/2} \exp[-\frac{1}{2}h(\beta - b^{**})N^{**}(\beta - b^{**})] h^{\frac{\nu^{**} - 1}{2}} \exp[-\frac{1}{2}\nu^{**} s^{**2}h],$$

hich by inspection is the normal-gamma distribution described in the eorem above.

The normal-gamma prior is Raiffa and Schlaifer's (1961) "natural conjugate" prior. The member of this class of distributions that is sometimes eed to represent prior ignorance is

$$f(\beta, \sigma^{-2}) \propto (\sigma^{-2})^{k/2}$$

ith  $N^* = 0$  and  $\nu^* = 0$ . Given these prior parameters, the posterior param-

ters become  $b$ ,  $XX'$ ,  $T$ , and  $ESS/T$ . The marginal posterior distribution of  $\beta$ , therefore, reproduces the classical result, except that the degrees of freedom parameter is  $T$  instead of  $T - k$ , a point further discussed below.

Observe also that the posterior mean (3.22) is formally the same as Equation (3.18), which describes the pooling of two samples of the same regression process. Thus the normal gamma prior is equivalent to a previous sample of the same process with estimate equal to  $b^*$ ,  $XX'$  matrix equal to  $N^*$ , degrees of freedom equal to  $\nu^*$ , and error sum of squares equal to  $\nu^* s^{*2}$ .

The posterior distribution for  $\beta$  associated with this normal-gamma prior has the dual features that it is unimodal and located at a fixed, weighted average of the sample location  $b$  and the prior location  $b^*$ . In that sense it never distinguishes sample information from prior information, no matter how strong their apparent conflict. This is so because a conjugate prior treats prior information as if it were a previous sample of the same process. It may be argued that most prior information is distinctly different from sample information, and when they are apparently in conflict, the posterior distribution ought to be multimodal with modes at both the sample location and the prior location. A distribution that has this feature is described in the following theorem, which uses a prior suggested by Dickey (1975).

STUDENT PRIOR

THEOREM 3.10. If  $\beta$  has a multivariate Student prior distribution independent of  $\sigma^2$ , and if  $\sigma^{-2}$  has a gamma prior distribution

$$f(\beta, \sigma^{-2}) = f_S^k(\beta | b^*, H^{*-1}, \nu_\beta^*) f_\gamma(\sigma^{-2} | s^{*2}, \nu_\sigma^*),$$

then the marginal posterior distribution for  $\beta$  is proportional to the product of two Student functions<sup>4</sup>

$$f(\beta | Y) \propto f_S^k(\beta | b, H^{-1}, \nu_\sigma^{**}) f_S^k(\beta | b, H^{*-1}, \nu_\beta)$$

where  $b$  is the least-squares estimate, and

$$H = (s^{**})^{-2} XX'$$

$$\nu_\sigma^{**} = \nu_\sigma^* + T - k$$

$$(s^{**})^2 = (\nu_\sigma^{**})^{-1} (\nu_\sigma^* s^{*2} + ESS).$$

<sup>4</sup>The product of two Student distributions can be written as a one-dimensional mixture of Student distributions (Dickey, 1975).



THE LINEAR-REGRESSION MODEL

of: Integrate  $\sigma^2$  from the product of the likelihood function and prior:

$$\begin{aligned} \beta|\mathbf{Y} &\propto \int_{-\infty}^{\infty} f(\mathbf{Y}|\beta, \sigma^2) f(\beta, \sigma^{-2}) d\sigma^{-2} \\ &= \int_{\sigma^{-2}} f_N(\mathbf{Y}|\mathbf{X}\beta, \sigma^2 I) f_T(\sigma^{-2} | s^{*2}, \nu_\sigma^*) f_S(\beta|\mathbf{b}^*, \mathbf{H}^{*-1}, \nu_\beta^*) d\sigma^{-2} \\ &= f_S^T(\mathbf{Y}|\mathbf{X}\beta, s^{*2} \mathbf{I}, \nu_\sigma^*) f_S(\beta|\mathbf{b}^*, \mathbf{H}^{*-1}, \nu_\beta^*). \end{aligned}$$

e Student distribution in  $\mathbf{Y}$  can be written as a Student function in  $\beta$  as

$$\begin{aligned} f_S^T(\mathbf{Y}|\mathbf{X}\beta, s^{*2} \mathbf{I}, \nu_\sigma^*) &\propto [\nu_\sigma^* s^{*2} + (\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta)]^{-(\nu_\sigma^* + T)/2} \\ &\propto [\nu_\sigma^* s^{*2} + ESS + (\beta - \mathbf{b})'\mathbf{X}'\mathbf{X}(\beta - \mathbf{b})]^{-(\nu_\sigma^* + k)/2} \\ &\propto f_S^k(\beta|\mathbf{b}, s^{*2}(\mathbf{X}'\mathbf{X})^{-1}, \nu_\sigma^*) \end{aligned}$$

h  $\nu_\sigma^*$  and  $s^{*2}$  defined as before.

The difference between the prior of Theorem (3.9) and the prior of eorem (3.10) is in the latter case  $\beta$  and  $\sigma^2$  are independent. The former ijugate prior implies that if information is obtained about  $\sigma^2$ , opinions out  $\beta$  change. Either  $\beta$  becomes more uncertain or more certain. Of irse, your priors are your own business, but I can say that I tend to fer to have  $\beta$  and  $\sigma^2$  independent. A counterargument is that if you cover that the process is noisy ( $\sigma^2$  large), you may come to doubt the idity of your prior information.<sup>5</sup>

As pointed out earlier, the conjugate prior cannot reproduce classical st squares because it does not allow loss of degrees of freedom. The ident-gamma prior with  $\mathbf{H}^* = 0$  and  $\nu_\sigma^* = 0$  does exactly reproduce classi- least squares with

$$t = \frac{\beta_i - b_i}{s[(\mathbf{X}'\mathbf{X})^{-1}]_{ii}}$$

en  $\mathbf{Y}$ , having a  $t$  distribution with  $T - k$  degrees of freedom. The natural conjugate prior was shown to be equivalent to a previous nple of the same process. The Student prior for  $\beta$ , together with the fuse prior for  $\sigma^2$  ( $\nu_\sigma^* = 0$ ), is equivalent to a previous sample of a regres- n process with a different variance. Bayesian pooling of two such nples without any other information would lead to the "double-Student" sterior of Theorem 3.10. Also, the modes of the posterior distribution

<sup>5</sup>To anticipate Chapter 9, it may make sense to write the regression process as  $\mathbf{Y} = \mathbf{X}(\beta + \mathbf{u})$  where  $\beta^c$  reflects the specification error. Although your prior for the true parameter  $\beta$  y be independent of  $\sigma^2$ , your prior for  $\beta^c$  may not be.

may be found by setting its logarithmic derivatives to zero

$$\mathbf{0} = \partial \ln f(\beta|\mathbf{Y}) / \partial \beta = \lambda \mathbf{X}'\mathbf{X}(\beta - \mathbf{b}) + \lambda^* \mathbf{H}^*(\beta - \mathbf{b}^*)$$

where

$$\lambda = \frac{\nu_\sigma^* + T}{\nu_\sigma^* s^{*2} + ESS + (\beta - \mathbf{b})'\mathbf{X}'\mathbf{X}(\beta - \mathbf{b})} \tag{3.26}$$

$$\lambda^* = \frac{\nu_\beta^* + k}{\nu_\beta^* + (\beta - \mathbf{b}^*)'\mathbf{H}^*(\beta - \mathbf{b}^*)} \tag{3.27}$$

With a suitable choice of  $\nu_\sigma^*$ ,  $\nu_\beta^*$ , and  $\mathbf{H}^*$ , these equations are equivalent to (3.19), (3.20), and (3.21) which describe pooling information from two different processes.

A GRAPHICAL PRESENTATION

A graphical analysis of inference with the linear regression model is instructive. The graphs now to be discussed apply to the results just reported, but they also apply to a wider class of distributions. We need only assume that  $\mathbf{A}\beta$  and  $\mathbf{u}$  have independent, spherically symmetric distributions.<sup>6</sup>

A random variable  $\mathbf{z}$  is said to have a spherically symmetric distribution if its density function depends only on the length of  $\mathbf{z}$ :  $f(\mathbf{z}) = cg(\mathbf{z}'\mathbf{z})$ . Thus the assumption of spherical symmetry implies that the densities of  $\mathbf{Y}$  and  $\beta$  may be written as

$$\begin{aligned} f(\mathbf{Y}|\beta) &\propto g_u[(\mathbf{Y} - \mathbf{X}\beta)'(\mathbf{Y} - \mathbf{X}\beta)] \\ f(\beta) &\propto g_\beta[(\beta - \mathbf{b}^*)'\mathbf{A}'\mathbf{A}(\beta - \mathbf{b}^*)]. \end{aligned}$$

Using (3.5) and setting  $\mathbf{A}'\mathbf{A} = \mathbf{N}^*$ , the posterior distribution can be written as

$$\begin{aligned} f(\beta|\mathbf{Y}) &\propto f(\mathbf{Y}|\beta) f(\beta) \\ &\propto g_u[ESS + (\beta - \mathbf{b})'\mathbf{X}'\mathbf{X}(\beta - \mathbf{b})] g_\beta[(\beta - \mathbf{b}^*)'\mathbf{N}^*(\beta - \mathbf{b}^*)]. \end{aligned}$$

The full posterior distribution, of course, depends on the functions  $g_u$  and  $g_\beta$ . The mode, however, necessarily lies on a curve. To find the equation for this curve, we may differentiate the logarithm of the posterior

$$\frac{\partial \ln f(\beta|\mathbf{Y})}{\partial \beta} = g_u^{-1} g_u' 2\mathbf{N}(\beta - \mathbf{b}) + g_\beta^{-1} g_\beta'(2\mathbf{N}^*(\beta - \mathbf{b}^*)). \quad (\mathbf{N} = \mathbf{X}'\mathbf{X})$$

<sup>6</sup>Hill (1969).

Setting this equal to zero and solving for  $\beta$  yields

$$\beta = (\mathbf{N} + \lambda \mathbf{N}^*)^{-1} (\mathbf{N} \mathbf{b} + \lambda \mathbf{N}^* \mathbf{b}^*) \tag{3.28}$$

where

$$\lambda = \frac{g_\beta^{-1} g_\beta'}{g_u^{-1} g_u'} \tag{3.29}$$

As  $\lambda$  varies from zero to infinity, Equation (3.28) sweeps out a curve in  $k$ -space affectionately called by Dickey (1975) the *curve d'occoltage*. It is anchored by the least-squares point  $\mathbf{b}$  at one end and by the "prior point"  $\mathbf{b}^*$  at the other. Choice of a point along the curve depends on the functions  $g_u$  and  $g_\beta$ , but spherical symmetry is sufficient to imply the curve. Geometrically, the curve d'occoltage is the locus of tangencies between the sample family of ellipsoids,  $(\beta - \mathbf{b})' \mathbf{N} (\beta - \mathbf{b}) = c$ , and the prior family  $(\beta - \mathbf{b}^*)' \mathbf{N}^* (\beta - \mathbf{b}^*) = c^*$ . To find this locus of tangencies, set the derivatives of  $(\beta - \mathbf{b})' \mathbf{N} (\beta - \mathbf{b}) + \lambda (\beta - \mathbf{b}^*)' \mathbf{N}^* (\beta - \mathbf{b}^*)$  to zero:  $2\mathbf{N}(\beta - \mathbf{b}) + 2\lambda \mathbf{N}^*(\beta - \mathbf{b}^*) = \mathbf{0}$ , which is just Equation (3.28).

A curve d'occoltage is illustrated in Figure 3.1. Ellipses around the least-squares point  $\mathbf{b}$  are isolikelihood ellipses. The data prefer the point  $\mathbf{b}$ , and the data are indifferent between any points on a given likelihood ellipse. The most preferred point from the standpoint of the prior is  $\mathbf{b}^*$ , and an ellipse around the prior point is a prior isodensity ellipse. The curve d'occoltage contains all points jointly preferred by data and prior. Given any point not on the curve, there is a better point on the curve, in the sense that neither the likelihood value nor the prior density is less, and at least one is greater.<sup>7</sup>

The location of modes on the curve d'occoltage depends on the functions  $g_u$  and  $g_\beta$ . One possibility is the exponential family indexed by the precision parameter  $h$ :

$$g(z^2|h) = e^{-\frac{1}{2}hz^2}$$

with

$$g'g^{-1} = -\frac{h}{2}$$

Thus letting the data and prior densities have different precision parameters,  $g_u(z^2) = g(z^2|h)$  and  $g_\beta(z^2) = g(z^2|h^*)$ , we would have, from (3.29),  $\lambda = h^*/h$ ; the posterior mode, since  $\lambda$  is independent of  $\beta$ , is just (3.28), using this value of  $\lambda$ .

<sup>7</sup>Economists are reminded of the analogous Edgeworth-Bowley diagram. This analogy is pursued in Chapter 5, and the curve d'occoltage is there referred to as the information-contrast curve.

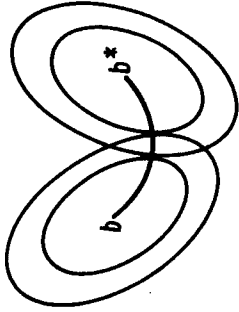


Fig. 3.1 Curve d'occoltage.

Describing this more traditionally, we have just assumed that the data and prior densities are normal  $\mathbf{u} \sim N(0, h^{-1}(\mathbf{B}\mathbf{B})^{-1})$ ,  $\beta \sim N(\mathbf{b}^*, h^{*-1}\mathbf{N}^{*-1})$ , with  $h/h^*$  known. This same result has been described in Theorem 3.9.

A somewhat less restrictive class of functions for labeling the information indifference curves is the Student family

$$g_S(z^2|a, \nu) = (a + z^2)^{-\nu/2} \tag{3.29}$$

with

$$g_S'g_S^{-1} = -\frac{\nu}{2(a + z^2)}$$

Letting  $g_u(z^2) = g_S(z^2|a, \nu)$  and  $g_\beta(z^2) = g_S(z^2|a^*, \nu^*)$  we have from (3.29)

$$\lambda = \frac{\nu^* / [a^* + (\beta - \mathbf{b}^*)' \mathbf{N} (\beta - \mathbf{b}^*)]}{\nu / [a + ESS + (\beta - \mathbf{b})' \mathbf{N} (\beta - \mathbf{b})]}$$

This family of labeling distributions has the property that the logarithmic derivative  $g'g^{-1}$  decreases with  $z^2$ . It is thus relatively steep around the origin and relatively flat elsewhere as indicated in Figure 3.2. This means

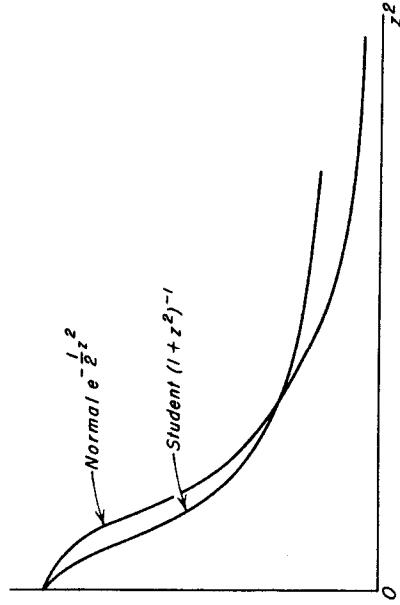


Fig. 3.2 Normal and Student labeling functions.

that an information source will be relatively resistant to moving the estimate of  $\beta$  away from its most preferred point in the neighborhood of that point and relatively indifferent to adjustments that occur farther from that point. Note also that the logarithmic derivative of the likelihood function depends not only on  $a$  and  $\nu$  but on  $ESS$  and  $\beta$  as well. The location of modes on the curve décolletage is, therefore, data dependent. These assumptions are described more traditionally in Theorem 3.10.

**SWEEPING OUT THE MEANS**

Ordinarily, there is a constant in a regression function about which we know very little relative to the prior information on the slopes. We show here that the effect of this structure of information is merely to sweep out the means of the observed variables.

Writing the constant explicitly, the regression process becomes

$$Y = X\beta + \mathbf{1}_T\alpha + u$$

with definitions as before, but with  $\mathbf{1}_T$  a vector of ones and  $\alpha$  the scalar "process level." Let  $u$  be normal with mean zero and covariance matrix  $\sigma^2\mathbf{1}_T$ , and let  $\alpha$  be normal with mean zero and variance  $v$ . We may then write

$$Y = X\beta + \varepsilon$$

where  $\varepsilon$  is normal with mean zero and covariance matrix  $\sigma^2\mathbf{1}_T + \mathbf{1}_T v \mathbf{1}'_T$ . The precision matrix for  $\varepsilon$  is

$$(\sigma^2\mathbf{1}_T + \mathbf{1}_T v \mathbf{1}'_T)^{-1} = \sigma^{-2} \left( \mathbf{1}_T - \mathbf{1}_T \left[ \mathbf{1}'_T \mathbf{1}_T + \frac{\sigma^2}{v} \right]^{-1} \mathbf{1}'_T \right).$$

In the limit as  $v$  increases the conditional distribution of  $Y$  given  $\beta$  becomes improper with singular precision matrix

$$\sigma^{-2}\mathbf{M} = \sigma^{-2} (\mathbf{1}_T - \mathbf{1}_T (\mathbf{1}'_T \mathbf{1}_T)^{-1} \mathbf{1}'_T).$$

The determinant of the variance matrix of  $Y$ , using (T18) in Appendix 1, is  $(\sigma^2)^T (1 + T v \sigma^{-2})$ , which for large  $v$  behaves like  $(\sigma^2)^{T-b}$ . Thus the limiting likelihood function is

$$f(Y|\beta) \propto (\sigma^2)^{-(T-1)/2} \exp \left[ -\frac{1}{2\sigma^2} (Y - X\beta)' \mathbf{M} (Y - X\beta) \right].$$

But  $\mathbf{M}$  is the idempotent matrix that was shown in Section 3.1 to sweep out the means of the variables. That is, given that the constant  $\alpha$  is diffuse and independent of the other random variables, inferences about the slope  $\beta$  may proceed by first subtracting out the means of the observables  $Y$  and  $X$  and then proceeding as above.

**3.4 Multivariate Normal Sampling**

In this section, Bayesian inference about the parameters of a multivariate normal distribution is discussed. (Several results here are referred to occasionally, but it is not necessary for the reader at this point to master this material.) The  $(k \times 1)$  vector  $z_i$  is to be normally distributed with mean vector  $\mu$  and variance matrix  $\Sigma$ . A sample of  $T$  independently distributed vectors  $z_i$  implies the likelihood function

$$L(\mu, \Sigma; z_1, z_2, \dots, z_T) \propto |\Sigma|^{-T/2} \exp \left[ -\frac{1}{2} \sum_i (z_i - \mu)' \Sigma^{-1} (z_i - \mu) \right].$$

Letting  $\bar{z} = \sum_i z_i / T$

$$S = \sum_i (z_i - \bar{z})(z_i - \bar{z})',$$

the exponent in the likelihood function can be written as

$$\begin{aligned} \sum_i (z_i - \mu)' \Sigma^{-1} (z_i - \mu) &= \sum_i [(z_i - \bar{z})' \Sigma^{-1} (z_i - \bar{z}) + (\bar{z} - \mu)' \Sigma^{-1} (\bar{z} - \mu)] \\ &= \sum_i \text{tr} [\Sigma^{-1} (z_i - \bar{z})(z_i - \bar{z})'] + T(\bar{z} - \mu)' \Sigma^{-1} (\bar{z} - \mu) \\ &= \text{tr}(\Sigma^{-1} S) + T(\bar{z} - \mu)' \Sigma^{-1} (\bar{z} - \mu). \end{aligned}$$

A normal-Wishart distribution for  $(\mu, \Sigma^{-1})$  is a convenient prior:

$$f_{NW}(\mu, \Sigma^{-1} | \bar{z}^*, S^*, T^*, \nu^*) \propto |\Sigma^{-1}|^{1/2} \exp \left[ -\frac{T^*}{2} (\mu - \bar{z}^*)' \Sigma^{-1} (\mu - \bar{z}^*) \right] \cdot |\Sigma^{-1}|^{(\nu^* - k - 1)/2} \exp \left[ -\frac{1}{2} \text{tr} \Sigma^{-1} S^* \right].$$

Combining this prior with the likelihood function and using (T10) in Appendix 1 the posterior distribution is seen to be in the normal-Wishart family with parameters

$$\begin{aligned} \bar{z}^{**} &= (T + T^*)^{-1} (T\bar{z} + T^*\bar{z}^*) \\ S^{**} &= S + S^* + (\bar{z} - \bar{z}^*)' \Sigma^{-1} (\bar{z} - \bar{z}^*) T T^* (T + T^*)^{-1} \\ T^{**} &= T + T^* \\ \nu^{**} &= T + \nu^*. \end{aligned}$$

The predictive distribution of the next sample vector,  $z_F$ , is used in Chapter 6. Conditional on  $\mu$  and  $\Sigma$ ,  $z_F$  is normal with mean vector  $\mu$  and covariance matrix  $\Sigma$ . Conditional on  $\Sigma$  and the data  $(\bar{z}, S, T)$ ,  $\mu$  is normal with mean  $\bar{z}^{**}$  and covariance matrix  $\Sigma / T^{**}$ . Thus integrating with

respect to  $\mu$ , the vector  $\mathbf{z}_F$  is normal with mean  $\bar{\mathbf{z}}^{**}$  and covariance matrix  $\Sigma(1 + (T^{**})^{-1})$ . Integrating with respect to  $\Sigma^{-1}$  as in Appendix 2, the distribution for  $\mathbf{z}_F$  is Student

$$f(\mathbf{z}_F | \bar{\mathbf{z}}, \mathbf{S}, T) = f_S^k(\mathbf{z}_F | \bar{\mathbf{z}}^{**}, \mathbf{S}^{**}(1 + (T^{**})^{-1}) / v^{**}, v^{**} - k + 1).$$

The results needed for Chapter 6 are the conditional moments of  $\mathbf{z}_F'$  given  $\mathbf{z}_F'$  where the partition of  $\mathbf{z}_F$  is  $\mathbf{z}_F' = (\mathbf{z}_F', \mathbf{z}_F')$ . Using the diffuse prior assumption that  $\mathbf{S}^* = 0$ ,  $T^* = 0$ ,  $v^* = 0$ , the conditional moments are

$$E(\mathbf{z}_F' | \mathbf{z}_F') = \bar{\mathbf{z}}' + \mathbf{S}_{IJ} \mathbf{S}_{JJ}^{-1} (\mathbf{z}_F' - \bar{\mathbf{z}}')$$

$$V(\mathbf{z}_F' | \mathbf{z}_F') = (\mathbf{S}_{II} - \mathbf{S}_{IJ} \mathbf{S}_{JJ}^{-1} \mathbf{S}_{JI})(1 + T^{-1})(T - k + 1 + k_j)^{-1}$$

where  $k_j$  is the number of elements in  $\mathbf{z}_F'$ .

# 4

CHAPTER

## HYPOTHESES-TESTING SEARCHES

4.1	Hypothesis Testing: A Judicial Analogy	93
4.2	Testing a Point Null Hypothesis Against a Point Alternative	99
4.3	Testing a Point Null Hypothesis Against a Composite Alternative	100
4.4	Weighted Likelihoods: Conjugate Priors	108
4.5	Weighted Likelihoods: Diffuse Priors	110
4.6	Conclusion	114

The first variety of specification search that we discuss corresponds to the familiar hypothesis-testing problem. We assume the existence of a set of  $M$  "models" or hypotheses of the form

$$H_i: \mathbf{Y} = \mathbf{X}_i \boldsymbol{\beta}_i + \mathbf{u}_i, \quad i = 1, \dots, M, \quad (4.1)$$

where  $\mathbf{Y}$  is a  $(T \times 1)$  vector of observable variables,  $\mathbf{X}_i$  is a  $(T \times k_i)$  matrix of observable explanatory variables,  $\boldsymbol{\beta}_i$  is a  $(k_i \times 1)$  vector of parameters, and  $\mathbf{u}_i$  is a  $(T \times 1)$  vector of unobservable disturbances assumed to be normally distributed with mean zero and variance-covariance matrix  $\sigma_i^2 \mathbf{I}_T$ . The statistical problem is to determine which of these  $M$  models did, in fact, generate the data and at the same time to make inferences about the coefficient vectors  $\boldsymbol{\beta}_i$ .

The formal classical theory of hypothesis testing describes the decision problem of selecting an action from among a set of feasible actions. An action is either right or wrong, depending on the "true state of nature," and the statistician is interested in being wrong as infrequently as possible. The problem being considered in this chapter involves a set of  $M$  actions of the form "act as if hypothesis  $H_i$  were true" and a set of  $M$  states of nature of the form "hypothesis  $H_i$  is true." An error occurs when action  $i$  is chosen but hypothesis  $j$  ( $i \neq j$ ) is true.