

Linear Models

One of the most useful statistical models is

$$X_i = \beta^\tau Z_i + \varepsilon_i, \quad i = 1, \dots, n,$$

where X_i is the i th observation and is often called the i th response; β is a p -vector of unknown parameters (main parameters of interest), $p < n$;

Z_i is the i th value of a p -vector of explanatory variables (or covariates); $\varepsilon_1, \dots, \varepsilon_n$ are random errors (not observed).

Data: $(X_1, Z_1), \dots, (X_n, Z_n)$.

Z_i 's are nonrandom or given values of a random p -vector, in which case our analysis is conditioned on Z_1, \dots, Z_n .

A matrix form of the model is

$$X = Z\beta + \varepsilon, \tag{1}$$

where $X = (X_1, \dots, X_n)$, $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$, and $Z =$ the $n \times p$ matrix whose i th row is the vector Z_i , $i = 1, \dots, n$.

Definition 3.4 (LSE).

Suppose that the range of β in model (6) is $B \subset \mathcal{R}^p$.

A *least squares estimator* (LSE) of β is defined to be any $\hat{\beta} \in B$ such that

$$\|X - Z\hat{\beta}\|^2 = \min_{b \in B} \|X - Zb\|^2.$$

For any $l \in \mathcal{R}^p$, $l^\tau \hat{\beta}$ is called an LSE of $l^\tau \beta$.

Throughout this book, we consider $B = \mathcal{R}^p$ unless otherwise stated.

Differentiating $\|X - Zb\|^2$ w.r.t. b , we obtain that any solution of

$$Z^\tau Zb = Z^\tau X$$

is an LSE of β .

Full rank Z

If the rank of the matrix Z is p , in which case $(Z^\tau Z)^{-1}$ exists and Z is said to be of full rank, then there is a unique LSE, which is

$$\hat{\beta} = (Z^\tau Z)^{-1} Z^\tau X.$$

Not full rank Z

If Z is not of full rank, then there are infinitely many LSE's of β . Any LSE of β is of the form

$$\hat{\beta} = (Z^T Z)^- Z^T X,$$

where $(Z^T Z)^-$ is called a *generalized inverse* of $Z^T Z$ and satisfies

$$Z^T Z (Z^T Z)^- Z^T Z = Z^T Z.$$

Generalized inverse matrices are not unique unless Z is of full rank, in which case $(Z^T Z)^- = (Z^T Z)^{-1}$

Assumptions

To study properties of LSE's of β , we need some assumptions on the distribution of X or ε (conditional on Z if Z is random and ε and Z are independent).

A1: ε is distributed as $N_n(0, \sigma^2 I_n)$ with an unknown $\sigma^2 > 0$.

A2: $E(\varepsilon) = 0$ and $\text{Var}(\varepsilon) = \sigma^2 I_n$ with an unknown $\sigma^2 > 0$.

A3: $E(\varepsilon) = 0$ and $\text{Var}(\varepsilon)$ is an unknown matrix.

Remarks

- Assumption A1 is the strongest and implies a parametric model.
- We may assume a slightly more general assumption that ε has the $N_n(0, \sigma^2 D)$ distribution with unknown σ^2 but a known positive definite matrix D .

Let $D^{-1/2}$ be the inverse of the square root matrix of D .

Then model (6) with assumption A1 holds if we replace X , Z , and ε by the transformed variables $\tilde{X} = D^{-1/2}X$, $\tilde{Z} = D^{-1/2}Z$, and $\tilde{\varepsilon} = D^{-1/2}\varepsilon$, respectively.

- A similar conclusion can be made for assumption A2.
- Under assumption A1, the distribution of X is $N_n(Z\beta, \sigma^2 I_n)$, which is in an exponential family \mathcal{P} with parameter $\theta = (\beta, \sigma^2) \in \mathcal{R}^p \times (0, \infty)$.
- However, if the matrix Z is not of full rank, then \mathcal{P} is not identifiable (see §2.1.2), since $Z\beta_1 = Z\beta_2$ does not imply $\beta_1 = \beta_2$.

Remarks

- Suppose that the rank of Z is $r \leq p$.
Then there is an $n \times r$ submatrix Z_* of Z such that

$$Z = Z_* Q \quad (2)$$

and Z_* is of rank r , where Q is a fixed $r \times p$ matrix, and

$$Z\beta = Z_* Q\beta.$$

- \mathcal{P} is identifiable if we consider the reparameterization $\tilde{\beta} = Q\beta$.
- The new parameter $\tilde{\beta}$ is in a subspace of \mathcal{R}^p with dimension r .
- In many applications, we are interested in estimating some linear functions of β , i.e., $\vartheta = l^\tau \beta$ for some $l \in \mathcal{R}^p$.
- From the previous discussion, however, estimation of $l^\tau \beta$ is meaningless unless $l = Q^\tau c$ for some $c \in \mathcal{R}^r$ so that

$$l^\tau \beta = c^\tau Q\beta = c^\tau \tilde{\beta}.$$

The following result shows that $l^T \beta$ is estimable if $l = Q^T c$, which is also necessary for $l^T \beta$ to be estimable under assumption A1.

Theorem 3.6

Assume model (6) with assumption A3.

- (i) A necessary and sufficient condition for $l \in \mathcal{R}^p$ being $Q^T c$ for some $c \in \mathcal{R}^r$ is $l \in \mathcal{R}(Z) = \mathcal{R}(Z^T Z)$, where Q is given by (2) and $\mathcal{R}(A)$ is the smallest linear subspace containing all rows of A .
- (ii) If $l \in \mathcal{R}(Z)$, then the LSE $l^T \hat{\beta}$ is unique and unbiased for $l^T \beta$.
- (iii) If $l \notin \mathcal{R}(Z)$ and assumption A1 holds, then $l^T \beta$ is not estimable.

Proof

(i) Note that $a \in \mathcal{R}(A)$ iff $a = A^T b$ for some vector b .
If $l = Q^T c$, then

$$l = Q^T c = Q^T Z_*^T Z_* (Z_*^T Z_*)^{-1} c = Z^T [Z_* (Z_*^T Z_*)^{-1} c].$$

Hence $l \in \mathcal{R}(Z)$.

Proof (continued)

If $I \in \mathcal{R}(Z)$, then $I = Z^\tau \zeta$ for some ζ and

$$I = (Z_* Q)^\tau \zeta = Q^\tau c, \quad c = Z_*^\tau \zeta.$$

(ii) If $I \in \mathcal{R}(Z) = \mathcal{R}(Z^\tau Z)$, then $I = Z^\tau Z \zeta$ for some ζ and by

$$\widehat{\beta} = (Z^\tau Z)^{-1} Z^\tau X,$$

$$E(I^\tau \widehat{\beta}) = E[I^\tau (Z^\tau Z)^{-1} Z^\tau X] = \zeta^\tau Z^\tau Z (Z^\tau Z)^{-1} Z^\tau Z \beta = \zeta^\tau Z^\tau Z \beta = I^\tau \beta.$$

If $\bar{\beta}$ is any other LSE of β , then, by $Z^\tau Z \bar{\beta} = Z^\tau X$,

$$I^\tau \widehat{\beta} - I^\tau \bar{\beta} = \zeta^\tau (Z^\tau Z) (\widehat{\beta} - \bar{\beta}) = \zeta^\tau (Z^\tau X - Z^\tau X) = 0.$$

(iii) Under A1, if there is an estimator $h(X, Z)$ unbiased for $I^\tau \beta$, then

$$I^\tau \beta = \int_{\mathcal{R}^n} h(x, Z) (2\pi)^{-n/2} \sigma^{-n} \exp \left\{ -\frac{1}{2\sigma^2} \|x - Z\beta\|^2 \right\} dx.$$

Differentiating w.r.t. β and applying Theorem 2.1 lead to

$$I^\tau = Z^\tau \int_{\mathcal{R}^n} h(x, Z) (2\pi)^{-n/2} \sigma^{-n-2} (x - Z\beta) \exp \left\{ -\frac{1}{2\sigma^2} \|x - Z\beta\|^2 \right\} dx,$$

which implies $I \in \mathcal{R}(Z)$.

Example 3.12 (Simple linear regression)

Let $\beta = (\beta_0, \beta_1) \in \mathcal{R}^2$ and $Z_i = (1, t_i)$, $t_i \in \mathcal{R}$, $i = 1, \dots, n$.

Then model (6) is called a *simple linear regression* model.

It turns out that

$$\begin{pmatrix} n & \sum_{i=1}^n t_i \\ \sum_{i=1}^n t_i & \sum_{i=1}^n t_i^2 \end{pmatrix}.$$

This matrix is invertible iff some t_i 's are different.

Thus, if some t_i 's are different, then the unique unbiased LSE of $l^T \beta$ for any $l \in \mathcal{R}^2$ is $l^T (Z^T Z)^{-1} Z^T X$, which has the normal distribution if assumption A1 holds.

The result can be easily extended to the case of *polynomial regression* of order p in which $\beta = (\beta_0, \beta_1, \dots, \beta_{p-1})$ and $Z_i = (1, t_i, \dots, t_i^{p-1})$.

Example 3.13 (One-way ANOVA)

Suppose that $n = \sum_{j=1}^m n_j$ with m positive integers n_1, \dots, n_m and that

$$X_i = \mu_j + \varepsilon_i, \quad i = k_{j-1} + 1, \dots, k_j, \quad j = 1, \dots, m,$$

where $k_0 = 0$, $k_j = \sum_{l=1}^j n_l$, $j = 1, \dots, m$, and $(\mu_1, \dots, \mu_m) = \beta$.

Let J_m be the m -vector of ones.

Then the matrix Z in this case is a block diagonal matrix with J_{n_j} as the j th diagonal column.

Consequently, $Z^T Z$ is an $m \times m$ diagonal matrix whose j th diagonal element is n_j .

Thus, $Z^T Z$ is invertible and the unique LSE of β is the m -vector whose j th component is

$$\frac{1}{n_j} \sum_{i=k_{j-1}+1}^{k_j} X_i, \quad j = 1, \dots, m.$$

Sometimes it is more convenient to use the following notation:

$$X_{ij} = X_{k_{i-1}+j}, \quad \varepsilon_{ij} = \varepsilon_{k_{i-1}+j}, \quad j = 1, \dots, n_i, \quad i = 1, \dots, m,$$

and

$$\mu_i = \mu + \alpha_i, \quad i = 1, \dots, m.$$

Then our model becomes

$$X_{ij} = \mu + \alpha_i + \varepsilon_{ij}, \quad j = 1, \dots, n_i, \quad i = 1, \dots, m, \quad (3)$$

which is called a *one-way analysis of variance* (ANOVA) model.

Under model (3), $\beta = (\mu, \alpha_1, \dots, \alpha_m) \in \mathcal{R}^{m+1}$.

The matrix Z under model (3) is not of full rank.

An LSE of β under model (3) is

$$\hat{\beta} = (\bar{X}, \bar{X}_{1.} - \bar{X}, \dots, \bar{X}_{m.} - \bar{X}),$$

where \bar{X} is still the sample mean of X_{ij} 's and $\bar{X}_{i.}$ is the sample mean of the i th group $\{X_{ij}, j = 1, \dots, n_i\}$.

The notation used in model (3) allows us to generalize the one-way ANOVA model to any s -way ANOVA model with a positive integer s under the so-called factorial experiments.

Example 3.14 (Two-way balanced ANOVA)

Suppose that

$$X_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + \varepsilon_{ijk}, \quad i = 1, \dots, a, j = 1, \dots, b, k = 1, \dots, c, \quad (4)$$

where a , b , and c are some positive integers.

Model (4) is called a two-way balanced ANOVA model.

If we view model (4) as a special case of model (6), then the parameter vector β is

$$\beta = (\mu, \alpha_1, \dots, \alpha_a, \beta_1, \dots, \beta_b, \gamma_{11}, \dots, \gamma_{1b}, \dots, \gamma_{a1}, \dots, \gamma_{ab}). \quad (5)$$

One can obtain the matrix Z and show that it is $n \times p$, where $n = abc$ and $p = 1 + a + b + ab$, and is of rank $ab < p$.

It can also be shown that an LSE of β is given by the right-hand side of (5) with μ , α_i , β_j , and γ_{ij} replaced by $\hat{\mu}$, $\hat{\alpha}_i$, $\hat{\beta}_j$, and $\hat{\gamma}_{ij}$, respectively, where

$$\hat{\mu} = \bar{X}_{...},$$

$$\hat{\alpha}_i = \bar{X}_{i..} - \bar{X}_{...},$$

$$\hat{\beta}_j = \bar{X}_{.j.} - \bar{X}_{...},$$

$$\hat{\gamma}_{ij} = \bar{X}_{ij.} - \bar{X}_{i..} - \bar{X}_{.j.} + \bar{X}_{...},$$

and a dot is used to denote averaging over the indicated subscript, e.g., with a fixed j ,

$$\bar{X}_{.j.} = \frac{1}{ac} \sum_{i=1}^a \sum_{k=1}^c X_{ijk}$$

Theorem 3.7 (UMVUE).

Consider model

$$X = Z\beta + \varepsilon \quad (6)$$

with assumption A1 (ε is distributed as $N_n(0, \sigma^2 I_n)$ with an unknown $\sigma^2 > 0$).

Then

- (i) The LSE $I^r \hat{\beta}$ is the UMVUE of $I^r \beta$ for any estimable $I^r \beta$.
- (ii) The UMVUE of σ^2 is $\hat{\sigma}^2 = (n-r)^{-1} \|X - Z\hat{\beta}\|^2$, where r is the rank of Z .

Proof of (i)

Let $\hat{\beta}$ be an LSE of β .

By $Z^r Zb = Z^r X$,

$$(X - Z\hat{\beta})^r Z(\hat{\beta} - \beta) = (X^r Z - X^r Z)(\hat{\beta} - \beta) = 0$$

and, hence,

$$\begin{aligned}
\|X - Z\beta\|^2 &= \|X - Z\hat{\beta} + Z\hat{\beta} - Z\beta\|^2 \\
&= \|X - Z\hat{\beta}\|^2 + \|Z\hat{\beta} - Z\beta\|^2 \\
&= \|X - Z\hat{\beta}\|^2 - 2\beta^\tau Z^\tau X + \|Z\beta\|^2 + \|Z\hat{\beta}\|^2.
\end{aligned}$$

Using this result and assumption A1, we obtain the following joint Lebesgue p.d.f. of X :

$$(2\pi\sigma^2)^{-n/2} \exp \left\{ \frac{\beta^\tau Z^\tau X}{\sigma^2} - \frac{\|X - Z\hat{\beta}\|^2 + \|Z\hat{\beta}\|^2}{2\sigma^2} - \frac{\|Z\beta\|^2}{2\sigma^2} \right\}.$$

By Proposition 2.1 and the fact that $Z\hat{\beta} = Z(Z^\tau Z)^{-1}Z^\tau X$ is a function of $Z^\tau X$, the statistic $(Z^\tau X, \|X - Z\hat{\beta}\|^2)$ is complete and sufficient for $\theta = (\beta, \sigma^2)$.

Note that $\hat{\beta}$ is a function of $Z^\tau X$ and, hence, a function of the complete sufficient statistic.

If $I^\tau \beta$ is estimable, then $I^\tau \hat{\beta}$ is unbiased for $I^\tau \beta$ (Theorem 3.6) and, hence, $I^\tau \hat{\beta}$ is the UMVUE of $I^\tau \beta$.

Proof of (ii)

From $\|X - Z\beta\|^2 = \|X - Z\hat{\beta}\|^2 + \|Z\hat{\beta} - Z\beta\|^2$ and $E(Z\hat{\beta}) = Z\beta$ (Theorem 3.6),

$$\begin{aligned} E\|X - Z\hat{\beta}\|^2 &= E(X - Z\beta)^\tau(X - Z\beta) - E(\beta - \hat{\beta})^\tau Z^\tau Z(\beta - \hat{\beta}) \\ &= \text{tr}\left(\text{Var}(X) - \text{Var}(Z\hat{\beta})\right) \\ &= \sigma^2[n - \text{tr}(Z(Z^\tau Z)^- Z^\tau Z(Z^\tau Z)^- Z^\tau)] \\ &= \sigma^2[n - \text{tr}((Z^\tau Z)^- Z^\tau Z)]. \end{aligned}$$

Since each row of $Z \in \mathcal{R}(Z)$, $Z\hat{\beta}$ does not depend on the choice of $(Z^\tau Z)^-$ in $\hat{\beta} = (Z^\tau Z)^- Z^\tau X$ (Theorem 3.6).

Hence, we can evaluate $\text{tr}((Z^\tau Z)^- Z^\tau Z)$ using a particular $(Z^\tau Z)^-$.

From the theory of linear algebra, there exists a $p \times p$ matrix C such that $CC^\tau = I_p$ and

$$C^\tau(Z^\tau Z)C = \begin{pmatrix} \Lambda & 0 \\ 0 & 0 \end{pmatrix},$$

Then, a particular choice of $(Z^T Z)^-$ is

$$(Z^T Z)^- = C \begin{pmatrix} \Lambda^{-1} & 0 \\ 0 & 0 \end{pmatrix} C^T \quad (7)$$

and

$$(Z^T Z)^- Z^T Z = C \begin{pmatrix} I_r & 0 \\ 0 & 0 \end{pmatrix} C^T$$

whose trace is r .

Hence $\hat{\sigma}^2$ is the UMVUE of σ^2 , since it is a function of the complete sufficient statistic and

$$E\hat{\sigma}^2 = (n-r)^{-1} E\|X - Z\hat{\beta}\|^2 = \sigma^2.$$

Residual vector

- The vector $X - Z\hat{\beta}$ is called the *residual vector* and $\|X - Z\hat{\beta}\|^2$ is called the *sum of squared residuals* and is denoted by *SSR*.
- The estimator $\hat{\sigma}^2$ is then equal to $SSR/(n-r)$.

- The Fisher information matrix is

$$\frac{1}{\sigma^2} \begin{pmatrix} Z^T Z & 0 \\ 0 & \frac{n}{2\sigma^2} \end{pmatrix}$$

- The UMVUE $I^T \hat{\beta}$ attains the information lower bound, but not $\hat{\sigma}^2$.

- Since

$$X - Z\hat{\beta} = [I_n - Z(Z^T Z)^- Z^T]X$$

and

$$I^T \hat{\beta} = I^T (Z^T Z)^- Z^T X$$

are linear in X , they are normally distributed under assumption A1.

- Also, using the generalized inverse matrix in (7), we obtain that

$$[I_n - Z(Z^T Z)^- Z^T]Z(Z^T Z)^- = Z(Z^T Z)^- - Z(Z^T Z)^- Z^T Z(Z^T Z)^- = 0,$$

which implies that $\hat{\sigma}^2$ and $I^T \hat{\beta}$ are independent (Exercise 58 in §1.6) for any estimable $I^T \beta$.

- $Z(Z^T Z)^- Z^T$ is a projection matrix, $[Z(Z^T Z)^- Z^T]^2 = Z(Z^T Z)^- Z^T$, hence

$$SSR = X^T [I_n - Z(Z^T Z)^- Z^T]X.$$

- The rank of $Z(Z^T Z)^{-1} Z^T$ is $\text{tr}(Z(Z^T Z)^{-1} Z^T) = r$.
- Similarly, the rank of the projection matrix $I_n - Z(Z^T Z)^{-1} Z^T$ is $n - r$.
- From

$$X^T X = X^T [Z(Z^T Z)^{-1} Z^T] X + X^T [I_n - Z(Z^T Z)^{-1} Z^T] X$$

and Theorem 1.5 (Cochran's theorem), SSR/σ^2 has the chi-square distribution $\chi_{n-r}^2(\delta)$ with

$$\delta = \sigma^{-2} \beta^T Z^T [I_n - Z(Z^T Z)^{-1} Z^T] Z \beta = 0.$$

Thus, we have proved the following result.

Theorem 3.8.

Consider model (6) with assumption A1. For any estimable parameter $l^T \beta$, the UMVUE's $l^T \hat{\beta}$ and $\hat{\sigma}^2$ are independent; the distribution of $l^T \hat{\beta}$ is $N(l^T \beta, \sigma^2 l^T (Z^T Z)^{-1} l)$; and $(n - r) \hat{\sigma}^2 / \sigma^2$ has the chi-square distribution χ_{n-r}^2 .