

9. Maximum Likelihood Estimation

I. Ordinary Least Squares Estimation:

- For a linear model

$$Y_j = \beta_0 + \beta_1 X_{1j} + \dots + \beta_r X_{rj} + \epsilon_j,$$

the OLS estimator for

$$\beta = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_r \end{bmatrix} \text{ is any } \mathbf{b} = \begin{bmatrix} b_0 \\ \vdots \\ b_r \end{bmatrix}$$

that minimizes the sum of squared residuals

$$Q(\mathbf{b}) = \sum_{j=1}^n (Y_j - b_0 - b_1 X_{1j} - \dots - b_r X_{rj})^2.$$

657

- The estimating equations (normal equations) are

$$\frac{\partial Q(\mathbf{b})}{\partial b_0} = -2 \sum_{j=1}^n (Y_j - b_0 - b_1 X_{1j} - \dots - b_r X_{rj}) = 0$$

and

$$\frac{\partial Q(\mathbf{b})}{\partial b_i} = -2 \sum_{j=1}^n X_{ij} (Y_j - b_0 - b_1 X_{1j} - \dots - b_r X_{rj}) = 0$$

for $i = 1, 2, \dots, r$

The matrix form of these equations is

$$(\mathbf{X}^T \mathbf{X}) \mathbf{b} = \mathbf{X}^T \mathbf{Y}$$

and a solution is

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

658

The OLS estimator for an estimable function $C^T \beta$ is

$$C^T \mathbf{b} = C^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

for any solution to the normal equations.

- $E(C^T \mathbf{b}) = C^T \beta$
- $Var(C^T \mathbf{b}) = C^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \Sigma \mathbf{X} [(X^T X)^{-1}]^T C$, where $\Sigma = Var(\mathbf{Y})$.
- The distribution of \mathbf{Y} is not completely specified.

659

For a Gauss-Markov model with

$$E(\mathbf{Y}) = \mathbf{X} \beta \text{ and } Var(\mathbf{Y}) = \sigma^2 \mathbf{I}$$

the OLS estimator of an estimable function $C^T \beta$ is the unique best linear unbiased estimator (b.l.u.e.) of $C^T \beta$.

- $E(C^T \mathbf{b}) = C^T \beta$
- $Var(C^T \mathbf{b}) = \sigma^2 C^T (\mathbf{X}^T \mathbf{X})^{-1} C$ is smaller than the variance of any other linear unbiased estimator for $C^T \beta$.
- The distribution of \mathbf{Y} is not completely specified.

660

II. Generalized Least Squares Estimation

Consider the Aitken model

$$E(Y) = X\beta \quad \text{and} \quad \text{Var}(Y) = \sigma^2 V$$

where V is a positive definite symmetric matrix of known constants and σ^2 is an unknown variance parameter.

- A GLS estimator for β is any b that minimizes

$$Q(b) = (Y - Xb)^T V^{-1} (Y - Xb)$$

(from Definition 3.8 with $\Sigma = \sigma^2 V$).

661

- The estimating equations are

$$(X^T V^{-1} X)b = X^T V^{-1} Y.$$

- A solution is

$$b_{GLS} = (X^T V^{-1} X)^{-1} X^T V^{-1} Y.$$

- For any estimable function $C^T \beta$ the unique b.l.u.e. is

$$C^T b_{GLS} = C^T (X^T V^{-1} X)^{-1} X^T V^{-1} Y$$

for any solution to the normal equations.

662

- $E(C^T b) = C^T \beta$ and $\text{Var}(C^T b) = \sigma^2 C^T (X^T V^{-1} X)^{-1} C$.

- The distribution of Y is not completely specified.

- An unbiased estimator for σ^2 in the Aitken model is

$$\begin{aligned} \hat{\sigma}_{GLS}^2 &= \frac{Y^T [V^{-1} - V^{-1} X (X^T V^{-1} X)^{-1} X^T V^{-1}] Y}{n - \text{rank}(X)} \\ &= \frac{(Y - X b_{GLS})^T V^{-1} (Y - X b_{GLS})}{n} \end{aligned}$$

663

- In practice, V may not be known. Then b_{GLS} and σ_{GLS}^2 can be approximated by replacing V with a consistent estimator:

- The estimator for $C^T \beta$ is not b.l.u.e.

- The estimator for σ^2 is not unbiased.

- Both estimators are consistent.

664

III. Maximum Likelihood Estimation

The model must include a specification of the joint distribution of the observations.

Example: Normal theory

Gauss-Markov model:

$$Y_j = \beta_0 + \beta_1 X_{1j} + \dots + \beta_r X_{rj} + \epsilon_j$$

where

$$\epsilon_j \sim \text{NID}(0, \sigma^2), \quad i = 1, \dots, n$$

or

$$Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix} \sim N(X\beta, \sigma^2 I)$$

665

- Find the parameter values that maximize the “likelihood” of the observed data.

For the normal-theory Gauss-Markov model, the likelihood function is

$$\begin{aligned} L(\beta, \sigma^2; Y_1, \dots, Y_n) \\ = \frac{1}{(2\pi)^{n/2} \sigma^n} e^{-\frac{1}{2\sigma^2} (Y - X\beta)^T (Y - X\beta)} \end{aligned}$$

Find values of β and σ^2 that maximize this likelihood function.

666

- This is equivalent to finding values of β and σ^2 that maximize the log-likelihood.

$$\begin{aligned} \ell(\beta, \sigma^2; Y_1, \dots, Y_n) \\ = \log(L(\beta, \sigma^2; Y_1, \dots, Y_n)) \\ = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) \\ \quad - \frac{1}{2\sigma^2} (Y - X\beta)^T (Y - X\beta) \\ = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) \\ \quad - \frac{1}{2\sigma^2} \sum_{j=1}^n (Y_j - \beta_0 - \dots - \beta_r X_{rj})^2 \end{aligned}$$

↗

this is minimized by an OLS estimator for β regardless of the value of σ^2

667

Solve the likelihood equations:

$$\begin{aligned} 0 &= \frac{\partial \ell(\beta, \sigma^2; Y)}{\partial \beta_0} \\ &= \frac{1}{\sigma^2} \sum_{j=1}^n (Y_j - \beta_0 - \dots - \beta_r X_{rj}) \\ 0 &= \frac{\partial \ell(\beta, \sigma^2; Y)}{\partial \beta_i} \\ &= \frac{1}{\sigma^2} \sum_{j=1}^n X_{ij} (Y_j - \beta_0 - \dots - \beta_r X_{rj}) \quad \text{for } i = 1, 2, \dots, r \\ 0 &= \frac{\partial \ell(\beta, \sigma^2; Y)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} \\ &\quad + \frac{1}{2(\sigma^2)^2} \sum_{j=1}^n (Y_j - \beta_0 - \dots - \beta_r X_{rj})^2 \end{aligned}$$

668

Solution:

$$\hat{\beta} = b_{OLS} = (X^T X)^{-1} X^T Y$$

and

$$\begin{aligned} \hat{\sigma}^2 &= \frac{1}{n} \sum_{j=1}^n (Y_j - \hat{\beta}_0 - \dots - \hat{\beta}_r X_{rj})^2 \\ &= \frac{1}{n} Y^T (I - P_X) Y \\ &= \frac{1}{n} SSE \end{aligned}$$

- This is a biased estimator for σ^2 .
- $\frac{1}{n - \text{rank}(X)} SSE$ is an unbiased estimator for σ^2 .
- $\frac{1}{n} SSE$ and $\frac{1}{n - \text{rank}(X)} SSE$ are asymptotically equivalent.

669

Normal-theory Aitken model

$$Y = X\beta + \epsilon$$

where $\epsilon \sim N(0, \sigma^2 V)$ and V is a known positive definite matrix.

The multivariate normal likelihood function is

$$L(\beta; Y) =$$

$$\frac{1}{(2\pi\sigma^2)^{n/2} |V|^{1/2}} e^{-\frac{1}{2\sigma^2} (Y - X\beta)^T V^{-1} (Y - X\beta)}$$

670

The log-likelihood function is

$$\begin{aligned} \ell(\beta; Y) &= -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(|V|) \\ &\quad - \frac{n}{2} \log(\sigma^2) \\ &\quad - \frac{1}{2\sigma^2} (Y - X\beta)^T V^{-1} (Y - X\beta) \end{aligned}$$

For any value of σ , the log-likelihood is maximized by finding a β that minimizes

$$(Y - X\beta)^T V^{-1} (Y - X\beta)$$

The estimating equations are

$$(X^T V^{-1} X)\beta = X^T V^{-1} Y$$

Solutions are of the form

$$\hat{\beta} = b_{GLS} = (X^T V^{-1} X)^{-1} X^T V^{-1} Y$$

When V is known the mle for β is also the generalized least squares estimator.

671

The additional estimating equation corresponding to σ^2 is

$$0 = \frac{\partial \ell(\beta, \sigma^2; Y)}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} (Y - X\beta)^T V^{-1} (Y - X\beta)$$

Substituting the solution to the other estimating equations for β , the solution is

$$\hat{\sigma}^2 = \frac{1}{n} (Y - Xb_{\text{GLS}})^T V^{-1} (Y - Xb_{\text{GLS}})$$

↗
This is a biased estimator for σ^2 .

672

When V contains unknown parameters:

- You could maximize the log-likelihood

$$\begin{aligned} \ell(\beta, \Sigma; Y) &= -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(|V|) \\ &\quad - \frac{n}{2} \log(\sigma^2) \\ &\quad - \frac{1}{2\sigma^2} (Y - X\beta)^T V^{-1} (Y - X\beta) \end{aligned}$$

with respect to β , σ^2 and the parameters in V .

673

- There may be no algebraic formulas for the solutions to the joint likelihood equations.
- The MLE's for σ^2 and the parameters in V are usually biased (too small).
- REML estimates are often used.

General Properties of MLE's

Regularity Conditions:

- (i) The parameter space has finite dimension, is closed and compact, and the true parameter vector is in the interior of the parameter space.
- (ii) Probability distributions defined by any two different values of the parameter vector are distinct (an identifiability condition).

674

(iii) First three partial derivatives of the log-likelihood function, with respect to the parameters

- (a) exist
- (b) are bounded by a function with a finite expectation.

675

(iv) The expectation of the negative of the matrix of second partial derivatives of the log-likelihood is

- (a) finite
- (b) positive definite

in a neighborhood of the true value of the parameter vector. This matrix is called the *Fisher information matrix*.

676

Suppose Y_1, \dots, Y_n are independent vectors of observations, with

$$Y_j = \begin{bmatrix} Y_{1j} \\ \vdots \\ Y_{pj} \end{bmatrix},$$

and the density function (or probability function) is

$$f(Y_j; \theta)$$

Then, the joint likelihood function is

$$L(\theta; Y_1, \dots, Y_n) = \prod_{j=1}^n f(Y_j; \theta)$$

677

The log-likelihood function is

$$\begin{aligned} \ell(\theta; Y_1, \dots, Y_n) &= \log(L(\theta; Y_1, \dots, Y_n)) \\ &= \sum_{j=1}^n \log(f(Y_j; \theta)). \end{aligned}$$

678

The score function

$$u(\theta) = \begin{bmatrix} u_1(\theta) \\ \vdots \\ u_r(\theta) \end{bmatrix} = \begin{bmatrix} \frac{\partial \ell(\theta; Y_1, \dots, Y_n)}{\partial \theta_1} \\ \vdots \\ \frac{\partial \ell(\theta; Y_1, \dots, Y_n)}{\partial \theta_r} \end{bmatrix}$$

is the vector of first partial derivatives of the log-likelihood function with respect to the elements of

$$\theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_r \end{bmatrix}.$$

The likelihood equations are

$$u(\theta; Y_1, \dots, Y_n) = 0$$

679

The maximum likelihood estimator (MLE)

$$\hat{\theta} = \begin{bmatrix} \hat{\theta}_1 \\ \vdots \\ \hat{\theta}_r \end{bmatrix}$$

is a solution to the likelihood equations, that maximizes the log-likelihood function.

680

Fisher information matrix:

$$\begin{aligned} i(\theta) &= \text{Var}(u(\theta; Y_1, \dots, Y_n)) \\ &= E(u(\theta; Y_1, \dots, Y_n) [u(\theta; Y_1, \dots, Y_n)]^T) \\ &= -E\left(\left[\frac{\partial^2 \ell(\theta; Y_1, \dots, Y_n)}{\partial \theta_r \partial \theta_k}\right]\right) \end{aligned}$$

681

Let

θ denote the parameter vector

$i(\theta)$ denote the Fisher information matrix

$\hat{\theta}$ denote the MLE for θ .

Then, if the Regularity Conditions are satisfied, we have the following results:

682

Result 9.1: $\hat{\theta}$ is a consistent estimator.

$$Pr \{(\hat{\theta} - \theta)^T(\hat{\theta} - \theta) > \epsilon\} \rightarrow 0$$

as $n \rightarrow \infty$, for any $\epsilon > 0$.

683

Result 9.2: Asymptotic normality

$$\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{\text{dist}} N(0, \lim_{n \rightarrow \infty} n[i(\theta)]^{-1})$$

as $n \rightarrow \infty$.

With a slight abuse of notation we may express this as

$$\hat{\theta} \underset{\circ}{\sim} N(\theta, [i(\theta)]^{-1})$$

for “large” sample sizes.

684

Result 9.3: If $\hat{\theta}$ is the mle for θ , then the mle for $g(\theta)$ is $g(\hat{\theta})$ for any function $g(\cdot)$.

685

References:

Anderson, T.W. (1984). **An Introduction to Multivariate Statistical Analysis**, (2nd ed.), Wiley, New York.

Cox, C. (1984). **American Statistician**, 38, pp. 283–287.

Cox, D.R. and Hinkley, D.V. (1974). **Theoretical Statistics**, Chapman & Hall, London (Chapters 8 and 9).

Rao, C.R. (1973). **Linear Statistical Inference**, Wiley, New York (Chapter 5).

686