

# Chapter 3

## Least Squares

Notation:  $y_i$  is a scalar, and  $x_i = (x_{i1}, \dots, x_{iK})'$  is a  $p \times 1$  vector.  $Y = (y_1, \dots, y_n)'$  is an  $n \times 1$  vector, and

$$X = \begin{bmatrix} x'_1 \\ x'_2 \\ \vdots \\ x'_n \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{22} & \cdots & x_{np} \end{bmatrix}$$

is an  $n \times p$  matrix.  $I_n$  is an  $n \times n$  identity matrix.

Ordinary least squares (OLS) is the most basic estimation technique in econometrics. It is simple and transparent. Understanding it thoroughly paves the way to study more sophisticated linear estimators. Moreover, many nonlinear estimators resemble the behavior of linear estimators in a neighborhood of the true value. In this lecture, we learn a series of facts from the linear algebra operation.

### 3.1 Estimator

As we have learned from the linear projection model, the projection coefficient  $\beta$  in the regression

$$y = x'\beta + e$$

can be written as

$$\beta = (E [xx'])^{-1} E [xy]. \quad (3.1)$$

We draw a pair of  $(y, x)$  from the joint distribution, and we mark it as  $(y_i, x_i)$  for  $i = 1, \dots, n$  repeated experiments. We possess a sample  $(y_i, x_i)_{i=1}^n$ .

*Remark 3.1.* Is  $(y_i, x_i)$  random or deterministic? Before we make the observation, they are treated as random variables whose realized values are uncertain.  $(y_i, x_i)$  is treated as random when we talk about statistical properties — statistical properties of a fixed number is meaningless. After we make the observation, they become deterministic values which cannot vary anymore.

*Remark 3.2.* In reality, we have at hand fixed numbers (more recently, words, photos, audio clips, video clips, etc., which can all be represented in digital formats with 0 and 1)

to feed into a computational operation, and the operation will return one or some numbers. All statistical interpretation about these numbers are drawn from the probabilistic thought experiments. A *thought experiment* is an academic jargon for a *story* in plain language. Under the axiomatic approach of probability theory, such stories are mathematical consistent and coherent. But mathematics is a tautological system, not science. The scientific value of a probability model depends on how close it is to the *truth* or implications of the truth. In this course, we suppose that the data are generated from some mechanism, which is taken as the truth. In the linear regression model for example, the joint distribution of  $(y, x)$  is the truth, while we are interested in the linear projection coefficient  $\beta$ , which is an implication of the truth as in (3.1). Probabilists suppose there is a dragon and try to tell the dragon's behaviors. Statisticians observe many snakes on earth, and try to tell what a dragon looks like.

The sample mean is a natural estimator of the population mean. Replace the population mean  $E[\cdot]$  in (3.1) by the sample mean  $\frac{1}{n} \sum_{i=1}^n \cdot$ , and the resulting estimator is

$$\hat{\beta} = \left( \frac{1}{n} \sum_{i=1}^n x_i x_i' \right)^{-1} \frac{1}{n} \sum_{i=1}^n x_i y_i = \left( \frac{X' X}{n} \right)^{-1} \frac{X' y}{n} = (X' X)^{-1} X' y$$

if  $X' X$  is invertible. This is one way to motivate the OLS estimator.

## 3.2 Geometry of OLS

There is natural geometry interpretation of the OLS estimator in a  $n$ -dimensional Euclidean space. Notice  $\mathcal{X} = \{Xb : b \in \mathbb{R}^p\}$  is the linear space spanned by the  $p$  columns of  $X = [X_1, \dots, X_p]$ , which is of  $p$ -dimension if the columns are linearly independent. The OLS estimator is the minimizer of  $\min_{b \in \mathbb{R}^p} \|Y - Xb\|$  (Square the Euclidean norm or not does not change the minimizer because  $a^2$  is a monotonic transformation for  $a \geq 0$ ). In other words,  $X\hat{\beta}$  is the point in  $\mathcal{X}$  such that it is the closest to the vector  $Y$  in terms of the Euclidean norm.

Define

$$\hat{Y} = X\hat{\beta} = X(X' X)^{-1} X' Y = P_X Y,$$

where  $P_X = X(X' X)^{-1} X'$  is the projector to the columns space of  $X$ . On the other hand, define

$$\hat{e} = Y - \hat{Y} = (I_n - P_X)Y = P_X^\perp Y,$$

where  $P_X^\perp$  is the projector to the null space of  $X$ . Since the null space and the column space are orthogonal, it is easy to verify that

$$\langle X\hat{\beta}, \hat{e} \rangle = \hat{\beta}' X' \hat{e} = 0'_p.$$

The Pythagorean theorem implies

$$\|Y\|^2 = \|X\hat{\beta}\|^2 + \|\hat{e}\|^2.$$

If  $X$  has full column rank and we use a SVD

$$X = \begin{matrix} U & S & V' \\ (n \times n) & (p \times p) & (p \times p) \end{matrix}$$

where  $U$  and  $V$  have orthonormal columns and  $S$  is diagonal with the  $p$  singular values of  $X$ , then the fitted values can be written as

$$\hat{Y} = X\hat{\beta} = X(X'X)^{-1}X'Y = UU'Y.$$

In other words, in terms of in-sample fitting, OLS using regressors  $X$  is equivalent to regressing  $Y$  on the orthonormal basis  $U$  for the column space of  $X$  (the two procedures have the same projection  $\hat{Y}$ ). The corresponding coefficients satisfy  $\hat{\beta} = VS^{-1}U'Y$ .

### 3.3 FWL Theorem

The Frisch-Waugh-Lovell (FWL) theorem is an algebraic fact about the formula of a sub-vector of the OLS estimator. Decompose  $X = (X_1, X_2)$  where  $X_1$  is  $n \times (p - 1)$  matrix and  $X_2$  is an  $n \times 1$  vector, so that

$$Y = X_1\hat{\beta}_1 + X_2\hat{\beta}_2 + \hat{e}. \quad (3.2)$$

To find  $\hat{\beta}_2$ , we first decompose  $X_2$  into its projection onto  $\text{span}(X_1)$  and its projection residual:

$$X_2 = P_1X_2 + P_1^\perp X_2,$$

where  $P_1 = X_1(X_1'X_1)^{-1}X_1'$  is the projector onto  $\text{span}(X_1)$  and  $P_1^\perp = I_n - P_1$ . Pre-multiply  $P_1^\perp$  on both sides of (3.2), we have

$$P_1^\perp Y = P_1^\perp X_1\hat{\beta}_1 + P_1^\perp X_2\hat{\beta}_2 + P_1^\perp \hat{e} = P_1^\perp X_2\hat{\beta}_2 + \hat{e}$$

since  $P_1^\perp X_1 = 0$  and  $P_1^\perp \hat{e} = \hat{e}$ .

Because  $\tilde{X}_2 = P_1^\perp X_2$  is orthogonal to  $\text{span}(X_1)$ , the coefficient on  $X_2$  in the full regression is pinned down entirely by  $\tilde{X}_2$ :

$$\hat{\beta}_2 = (\tilde{X}_2'\tilde{X}_2)^{-1}\tilde{X}_2'Y = (X_2'P_1^\perp X_2)^{-1}X_2'P_1^\perp Y.$$

### 3.4 Omitted Variable Bias

Suppose that researcher A runs a full regression  $Y = X_1\hat{\beta}_1 + X_2\hat{\beta}_2 + \hat{e}$ . When he sends the data to researcher B, he omits  $X_2$  by accident. Research B can only estimate the coefficient for  $X_1$ . What would be its different from researcher A's  $\hat{\beta}_1$ ?

Notice in researcher A's world:

$$\begin{aligned} Y &= X_1\hat{\beta}_1 + X_2\hat{\beta}_2 + \hat{e} \\ &= X_1\hat{\beta}_1 + (P_1X_2 + P_1^\perp X_2)\hat{\beta}_2 + \hat{e} \\ &= (X_1\hat{\beta}_1 + P_1X_2\hat{\beta}_2) + (P_1^\perp X_2\hat{\beta}_2 + \hat{e}) \\ &= X_1(\hat{\beta}_1 + \hat{\pi}\hat{\beta}_2) + (P_1^\perp X_2\hat{\beta}_2 + \hat{e}), \end{aligned}$$

where  $\hat{\pi} = (X_1' X_1)^{-1} X_1' X_2$ . In other words, researcher B will obtain the coefficient  $(\hat{\beta}_1 + \hat{\pi}\hat{\beta}_2)$  associated with  $X_1$ , and his residual is  $(P_1^\perp X_2 \hat{\beta}_2 + \hat{e})$ .

Intuitively, regressing on  $X_1$  only forces the component  $P_1 X_2 = X_1 \hat{\pi}$  to be absorbed by the coefficient on  $X_1$ , while the orthogonal component  $P_1^\perp X_2 \hat{\beta}_2$  becomes part of the residual.

## 3.5 Summary

The derivations above are finite-sample linear algebra in  $\mathbb{R}^n$ :  $X_1$  and  $X_2$  span subspaces of  $\mathbb{R}^n$ , and  $P_1$  and  $P_1^\perp$  are ordinary projection matrices. There are population counterparts of both the FWL theorem and the omitted variable bias formula. In that setting,  $Y$  and the components of  $X$  are random variables that live in a Hilbert space (typically an  $L^2$  space), and “projection” is defined by the linear projection (orthogonality) operator rather than an  $n \times n$  matrix. The same logic goes through with population inner products (expectations) replacing sample inner products, but the geometry is more abstract, so we omit the full display.

**Historical notes:** Carl Friedrich Gauss (1777–1855) claimed he had come up with the operation of OLS in 1795. With only three data points at hand, Gauss successfully applied his method to predict the location of the dwarf planet Ceres in 1801. While Gauss did not publish the work on OLS until 1809, Adrien-Marie Legendre (1752–1833) presented this method in 1805. Today people tend to attribute OLS to Gauss, assuming that a giant like Gauss had no need to tell a lie to steal Legendre’s discovery.

Zhentao Shi. January 14, 2026