

Applied Econometrics

Derek Li

Contents

1	Introduction	2
1.1	Econometric Methodology	2
1.2	Data Types	2
1.2.1	Cross-Section	2
1.2.2	Time Series	2
1.2.3	Repeated Cross-Section	2
1.2.4	Panel Data	2
2	Review of Statistics	3
2.1	Estimator	3
2.2	Sampling Distribution	3
2.3	Confidence Interval	4
2.4	Proves of Some Theorems and Results (optional)	4
3	Simple Regression	6
3.1	Basic Property	6
3.2	Econometric Model	6
3.3	Estimator: OLS	7
3.4	Properties of OLS	7
4	Multiple Regression: Estimation	9
4.1	Econometric Model	9

1 Introduction

Econometrics are quantitative methods of analyzing and interpreting economic data, which need to combine economic theory, math and statistics, data, and statistical or econometrics software.

We focus on estimating economic relationships, testing hypothesis involving economic behavior, and forecasting the behavior of economic variables.

1.1 Econometric Methodology

- Ask a question - statement of theory or hypothesis
- Specification of economic model
- Specification of econometric model
- Collection of Data
- Estimation of the econometric model
- Hypothesis testing
- Prediction or forecasting

1.2 Data Types

There are different data structures: cross-section, time series, repeated cross-section, and panel data.

1.2.1 Cross-Section

Cross-section consists of a sample of individuals, households, firms, countries, etc, taken at a given point in time. Observations are generally independent draws from the population. It is commonly indexed by i as x_i .

1.2.2 Time Series

Time series consists of observations on a variable or several variables over time. Observations are almost never independent of each other. It is commonly indexed by t as x_t .

1.2.3 Repeated Cross-Section

Repeated cross-section consists of two or more cross-sectional data in different points in time, and is different units in different periods. It is commonly index by it as x_{it} .

Example 1.1. Suppose two cross-sectional household surveys are taking in 1985 and 1990. In 1985, a random sample of households is surveyed for variables such as income. In 1990, a new random sample of households is taken using the same survey questions. This is a repeated cross-section data set.

1.2.4 Panel Data

Panel data consists of a time series for each cross-sectional unit in the data set. Observations are independent among units and dependent over time for each unit. It is commonly index by it as x_{it} .

2 Review of Statistics

Suppose $X \sim \mathcal{N}(\mu, \sigma^2)$. Assume a random sample $\{x_1, \dots, x_n\}$, i.e., identically and independently distributed (i.i.d.).

2.1 Estimator

Definition 2.1 (Statistic). A statistic is a function of the data.

Definition 2.2 (Estimator). An estimator is a statistic that is used to estimate the parameter of interest.

Take μ as an example, the proposed estimator is sample average

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n x_i.$$

2.2 Sampling Distribution

We have

$$\mathbb{E}[\bar{X}_n] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[x_i] = \frac{1}{n} \sum_{i=1}^n \mu = \mu,$$

i.e., \bar{X}_n is an unbiased estimator of μ . Besides,

$$\text{Var}[\bar{X}_n] = \text{Var}\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{1}{n^2} \text{Var}\left[\sum_{i=1}^n x_i\right] \stackrel{\text{independent}}{=} \frac{1}{n^2} \sum_{i=1}^n \text{Var}[x_i] = \frac{\sigma^2}{n}.$$

Definition 2.3 (Consistency). Let W_n be an estimator of θ based on a sample Y_1, \dots, Y_n . Then W_n is a consistent estimator of θ if for every $\varepsilon > 0$,

$$P(|W_n - \theta| > \varepsilon) \rightarrow 0 \text{ as } n \rightarrow \infty.$$

As commonly stated, an estimator is called consistent when its sampling distribution becomes more and more concentrated around the parameters of interest as the sample size increases. Note that \bar{X}_n is a consistent estimator of μ .

Theorem 2.1. If $Y_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$, $Y_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$, then

$$Y_1 + Y_2 \sim \mathcal{N}(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2 + 2\text{Cov}(Y_1, Y_2)),$$

By theorem, we have the sampling distribution

$$\bar{X}_n \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right),$$

provided $X \sim \mathcal{N}(\mu, \sigma^2)$.

2.3 Confidence Interval

Theorem 2.2. If Y has $\mathbb{E}[Y] = \mu$, $\text{Var}[Y] = \sigma^2$, then $Z = \frac{Y - \mu}{\sigma}$ is such that $\mathbb{E}[Z] = 0$, $\text{Var}[Z] = 1$.

By theorem, we have

$$Z = \frac{\bar{X}_n - \mu}{\sqrt{\sigma^2/n}} \sim \mathcal{N}(0, 1).$$

Therefore,

$$\begin{aligned} 1 - \alpha &= P(-z_{\frac{\alpha}{2}} \leq Z \leq z_{\frac{\alpha}{2}}) = P\left(-z_{\frac{\alpha}{2}} \leq \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \leq z_{\frac{\alpha}{2}}\right) \\ &= P\left(\bar{X}_n - z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X}_n + z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}\right), \end{aligned}$$

i.e., $(1 - \alpha)\%$ confidence interval is

$$\left[\bar{X}_n - z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}, \bar{X}_n + z_{\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \right].$$

2.4 Proves of Some Theorems and Results (optional)

Theorem 2.3 (Markov's Inequality). If X is a nonnegative random variable and $a > 0$, then

$$P(X \geq a) \leq \frac{\mathbb{E}[X]}{a}.$$

Proof. Since X is a nonnegative random variable, we have

$$\begin{aligned} \mathbb{E}[X] &= \int_0^\infty x f(x) dx = \int_0^a x f(x) dx + \int_a^\infty x f(x) dx \\ &\geq \int_a^\infty x f(x) dx \geq \int_a^\infty a f(x) dx = a \int_a^\infty f(x) dx = a P(X \geq a). \end{aligned}$$

Hence

$$P(X \geq a) \leq \frac{\mathbb{E}[X]}{a}.$$

□

Theorem 2.4 (Chebyshev Inequality). For any $b > 0$,

$$P(|X - \mathbb{E}[X]| \geq b) \leq \frac{\text{Var}[X]}{b^2}.$$

Proof. By Markov's Inequality, we have

$$P((X - \mathbb{E}[X])^2 \geq b^2) \leq \frac{\mathbb{E}[(X - \mathbb{E}[X])^2]}{b^2} = \frac{\text{Var}[X]}{b^2}.$$

Therefore,

$$P(|X - \mathbb{E}[X]| \geq b) \leq \frac{\text{Var}[X]}{b^2}.$$

□

Theorem 2.5 (Weak Law of Large Numbers). Let X_1, \dots, X_n be a sequence of independent random variables with $\mathbb{E}[X_i] = \mu$. Let $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$, then

$$\forall \varepsilon > 0, \lim_{n \rightarrow \infty} P(|\bar{X}_n - \mu| > \varepsilon) = 0.$$

Proof. By Chebyshev Inequality, for all $\varepsilon > 0$,

$$P(|\overline{X}_n - \mathbb{E}[\overline{X}_n]| > \varepsilon) \leq \frac{\text{Var}[\overline{X}]}{\varepsilon^2} \Leftrightarrow 0 \leq P(|\overline{X}_n - \mu| > \varepsilon) \leq \frac{\sigma^2}{n\varepsilon}.$$

Since

$$\lim_{n \rightarrow \infty} \frac{\sigma^2}{n\varepsilon} = 0,$$

then

$$\lim_{n \rightarrow \infty} P(|\overline{X}_n - \mu| > \varepsilon) = 0.$$

□

It follows that \overline{X}_n is a consistent estimator of μ .

3 Simple Regression

3.1 Basic Property

Property 3.1. $\sum_{i=1}^n (x_i - \bar{x}) = 0$.

Proof. $\sum_{i=1}^n (x_i - \bar{x}) = \sum_{i=1}^n x_i - n\bar{x} = n\bar{x} - n\bar{x} = 0$. □

Property 3.2. $\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n x_i(y_i - \bar{y}) = \sum_{i=1}^n (x_i - \bar{x})y_i$.

Proof. On the one hand,

$$\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n x_i(y_i - \bar{y}) - \sum_{i=1}^n \bar{x}(y_i - \bar{y}) = \sum_{i=1}^n x_i(y_i - \bar{y}).$$

On the other hand,

$$\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n (x_i - \bar{x})y_i - \sum_{i=1}^n (x_i - \bar{x})\bar{y} = \sum_{i=1}^n (x_i - \bar{x})y_i.$$

□

Corollary 3.1. $\sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i(x_i - \bar{x})$.

Property 3.3. $\mathbb{E}[Y] = \mathbb{E}[\mathbb{E}[Y|X]]$.

3.2 Econometric Model

Let (Y, X, U) be random variables with joint distribution s.t. $Y = g(X, U)$, where Y is dependent variable, X is explanatory (regressor/covariate) variable, and U is unobservable variable.

Assumption 3.1 (Linear in Parameters). $Y = \beta_0 + \beta_1 X + U$.

Example 3.1. $y = \beta_0 + \beta_1 x^2 + u$, then $\frac{\partial y}{\partial x} = 2\beta_1 x$.

Example 3.2. $\ln y = \beta_0 + \beta_1 \ln x + u$, then $\frac{\partial \ln y}{\partial \ln x} = \beta_1 \approx \frac{\Delta y\%}{\Delta x\%}$.

Assumption 3.2 (Zero Conditional Mean). $\mathbb{E}[U|X] = \mathbb{E}[U]$, and $\mathbb{E}[U] = 0$.

Example 3.3. Let Y be wage, X be training program, and U be ability. If training is assigned randomly, then X and U are fully independent. Nevertheless, if let X be education, then the education may influence the ability so that $\mathbb{E}[U|X=0] \neq \mathbb{E}[U|X=1]$, then A2 is violated.

From A1 and A2, we have

$$\begin{aligned} \mathbb{E}[Y|X] &\stackrel{A1}{=} \mathbb{E}[\beta_0 + \beta_1 X + U|X] = \beta_0 + \beta_1 \mathbb{E}[X|X] + \mathbb{E}[U|X] \\ &\stackrel{A2}{=} \beta_0 + \beta_1 X + \mathbb{E}[U] = \beta_0 + \beta_1 X. \end{aligned}$$

Assumption 3.3 (Random Sample). $\{(x_i, y_i), i = 1, \dots, n\}$ is i.i.d.

Assumption 3.4 (Sample Variation). $\{x_1, \dots, x_n\}$ are not all the same.

Assumption 3.5 (Homoscedasticity). $\text{Var}[U|X] = \sigma_U^2$.

From A1 and A5, we have

$$\text{Var}[Y|X] \stackrel{A1}{=} \text{Var}[\beta_0 + \beta_1 X + U|X] = \text{Var}[U|X] \stackrel{A5}{=} \sigma_U^2.$$

3.3 Estimator: OLS

We want to solve

$$\min_{\hat{\beta}_0, \hat{\beta}_1} \sum_{i=1}^n Q(\hat{\beta}_0, \hat{\beta}_1) = \sum_{i=1}^n \left(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i \right)^2.$$

Let

$$\begin{aligned} \frac{\partial Q}{\partial \hat{\beta}_0} &= - \sum_{i=1}^n 2 \left(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i \right) = 0, \\ \frac{\partial Q}{\partial \hat{\beta}_1} &= - \sum_{i=1}^n 2 \left(y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i \right) x_i = 0. \end{aligned}$$

Therefore,

$$\hat{\beta}_0 = \bar{y} - \beta_1 \bar{x}.$$

Besides,

$$\begin{aligned} \sum_{i=1}^n \left(x_i y_i - \hat{\beta}_0 x_i - \hat{\beta}_1 x_i^2 \right) &= \sum_{i=1}^n (x_i y_i - x_i \bar{y} + \hat{\beta}_1 x_i \bar{x} - \hat{\beta}_1 x_i^2) \\ &= \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \bar{y} + \sum_{i=1}^n \hat{\beta}_1 x_i \bar{x} - \sum_{i=1}^n \hat{\beta}_1 x_i^2 \\ &= \sum_{i=1}^n x_i y_i - n \bar{x} \bar{y} + n \hat{\beta}_1 \bar{x}^2 - \hat{\beta}_1 \sum_{i=1}^n x_i^2 = 0, \end{aligned}$$

and thus

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2} = \frac{\sum_{i=1}^n (x_i y_i - x_i \bar{y})}{\sum_{i=1}^n (x_i^2 - x_i \bar{x})} = \frac{\sum_{i=1}^n x_i (y_i - \bar{y})}{\sum_{i=1}^n x_i (x_i - \bar{x})} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}.$$

3.4 Properties of OLS

Property 3.4. $\mathbb{E}[\hat{\beta}_1] = \beta_1.$

Proof. We have

$$\begin{aligned} \hat{\beta}_1 &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n (x_i - \bar{x}) y_i}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(\beta_0 + \beta_1 x_i + u_i)}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ &= \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \left[\sum_{i=1}^n (x_i - \bar{x}) \beta_0 + \sum_{i=1}^n (x_i - \bar{x}) \beta_1 x_i + \sum_{i=1}^n (x_i - \bar{x}) u_i \right] \\ &= \beta_1 + \frac{\sum_{i=1}^n (x_i - \bar{x}) u_i}{\sum_{i=1}^n (x_i - \bar{x})^2}. \end{aligned}$$

By A4, $\sum_{i=1}^n (x_i - \bar{x})^2 \neq 0$, and thus,

$$\begin{aligned}
\mathbb{E} \left[\hat{\beta}_1 | x_1, \dots, x_n \right] &= \mathbb{E} \left[\beta_1 + \frac{\sum_{i=1}^n (x_i - \bar{x}) u_i}{\sum_{i=1}^n (x_i - \bar{x})^2} \middle| x_1, \dots, x_n \right] \\
&= \beta_1 + \mathbb{E} \left[\frac{\sum_{i=1}^n (x_i - \bar{x}) u_i}{\sum_{i=1}^n (x_i - \bar{x})^2} \middle| x_1, \dots, x_n \right] \\
&= \beta_1 + \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) \mathbb{E}[u_i | x_1, \dots, x_n] \\
&\stackrel{A3}{=} \beta_1 + \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) \mathbb{E}[u_i | x_i] \\
&\stackrel{A2}{=} \beta_1 + \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_i - \bar{x}) \cdot 0 = \beta_1.
\end{aligned}$$

Therefore, $\mathbb{E} [\hat{\beta}_1] = \mathbb{E} [\mathbb{E} [\hat{\beta}_1 | x_1, \dots, x_n]] = \mathbb{E} [\beta_1] = \beta_1$. As a consequence, OLS is unbiased. \square

Property 3.5. $\text{Var} [\hat{\beta}_1 | x_1, \dots, x_n] = \frac{\sigma_u^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$

Proof. We have

$$\begin{aligned}
\text{Var} [\hat{\beta}_1] &= \text{Var} \left[\beta_1 + \frac{\sum_{i=1}^n (x_i - \bar{x}) u_i}{\sum_{i=1}^n (x_i - \bar{x})^2} \middle| x_1, \dots, x_n \right] \\
&= \frac{1}{\left(\sum_{i=1}^n (x_i - \bar{x})^2 \right)^2} \text{Var} \left[\sum_{i=1}^n (x_i - \bar{x}) u_i \middle| x_1, \dots, x_n \right] \\
&\stackrel{A3}{=} \frac{1}{\left(\sum_{i=1}^n (x_i - \bar{x})^2 \right)^2} \sum_{i=1}^n (x_i - \bar{x})^2 \text{Var} [u_i | x_1, \dots, x_n] \\
&\stackrel{A3}{=} \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \text{Var} [u_i | x_i] \stackrel{A5}{=} \frac{\sigma_u^2}{\sum_{i=1}^n (x_i - \bar{x})^2}.
\end{aligned}$$

\square

Notice that larger $\sum_{i=1}^n (x_i - \bar{x})^2$ implies smaller $\text{Var} [\hat{\beta}_1]$.

Theorem 3.1 (Gauss-Markov Theorem). Under A1 to A5, OLS is the best linear unbiased estimator.

4 Multiple Regression: Estimation

4.1 Econometric Model

Assumption 4.1 (Linear in Parameters). $Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_k X_k + U$.

Assumption 4.2 (Zero Conditional Mean). $\mathbb{E}[U|X_1, \cdots, X_k] = 0$.

Assumption 4.3 (Random Sample). $\{(y_i, x_{1i}, \cdots, x_{ki}), i = 1, \cdots, n\}$ is i.i.d.

Assumption 4.4 (No Perfect Collinearity). In the sample (and therefore in the population), none of the independent variables is constant, and there are no exact linear relationships among the independent variables.