

Section 3. Simple Regression

- Omitted Variable Bias

1. **Sampling Properties of OLS estimators**
2. **What's a Covariance?**
3. **More sampling properties of OLS estimators**
4. **Standard Errors of OLS estimators**
5. **Confidence intervals**
6. **CLT demonstration**
7. **Omitted Variables and Omitted Variable Bias (prelude to Section 4)**

1. Sampling Properties of OLS estimators

Review:

- Unbiased: $E(b_0) = \beta_0$, $E(b_1) = \beta_1$
- Consistent: $\text{plim}(b_0) = \beta_0$, $\text{plim}(b_1) = \beta_1$

New:

- Asymptotically Normal by the CLT

Derivation

$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$= \frac{\sum (x_i - \bar{x}) y_i - \cancel{\sum (x_i - \bar{x}) \bar{y}}}{\sum (x_i - \bar{x}) x_i - \cancel{\sum (x_i - \bar{x}) \bar{x}}}$$

$$= \frac{1}{\sum (x_i - \bar{x}) x_i} \sum (x_i - \bar{x}) [\beta_0 + \beta_1 x_i + u_i]$$

$$= \frac{1}{\sum (x_i - \bar{x}) x_i} \left[\cancel{\sum (x_i - \bar{x})} \beta_0 + \beta_1 \sum (x_i - \bar{x}) x_i + \sum (x_i - \bar{x}) u_i \right]$$

$$= 0 + \beta_1 + \frac{\sum (x_i - \bar{x}) u_i}{\sum (x_i - \bar{x}) x_i}$$

$$E(b_1 | X) = E \left(\beta_1 + \frac{\sum (x_i - \bar{x}) u_i}{\sum (x_i - \bar{x}) x_i} \mid X \right) = \beta_1 + \frac{1}{\sum (x_i - \bar{x}) x_i} \left[\sum E(x_i - \bar{x}) u_i \mid X \right]$$

$$\Rightarrow E(b_1) = \beta_1$$

$$= \beta_1 + 0 = \beta_1 \quad \text{FOR ALL POSSIBLE } X$$

LEMMA:

$$\#1) \sum_{i=1}^N (x_i - \bar{x}) = \sum x_i - \sum \frac{x_i}{N}$$

$$= \sum x_i - N \cdot \frac{\sum x_i}{N} = 0$$

$$\#2) \sum_{i=1}^N (x_i - \bar{x}) C, \quad C \text{ is a constant}$$

$$= (x_1 - \bar{x})C + (x_2 - \bar{x})C + \dots + (x_N - \bar{x})C$$

$$= 0$$

$$\#3) E(Z|X) = A \quad \text{FOR ALL } X$$

$$\Rightarrow E(Z) = A$$

Derivation

$$b_1 = \beta_1 + \frac{\sum (x_i - \bar{x}) u_i}{\sum (x_i - \bar{x})^2}$$

$$\begin{aligned} \text{plim } b_1 &= \text{plim} (\beta_1) + \text{plim} \left(\frac{\sum (x_i - \bar{x}) u_i}{\sum (x_i - \bar{x})^2 / N} \right) \\ &= \beta_1 + \frac{\text{Cov}(x, u) \times 1}{V(x) \times 1} = 0 \end{aligned}$$

$$\left. \begin{aligned} \text{Cov}(x, u) &= E(x - E(x))(u - E(u)) \\ &= E \left[E(x - E(x))(u - E(u)) | x \right] = 0 \\ \text{L.L.N} \\ E(z) &= \mu_z \Rightarrow \text{plim}(\bar{z}) = \mu_z \end{aligned} \right\}$$

$$\Rightarrow \text{plim } b_1 = \beta_1 + 0 = \beta_1$$

$$E \left[\frac{\sum (x_i - \bar{x}) u_i}{(N-2)} \frac{(N-2)}{N} \right] = \frac{N-2}{N}$$

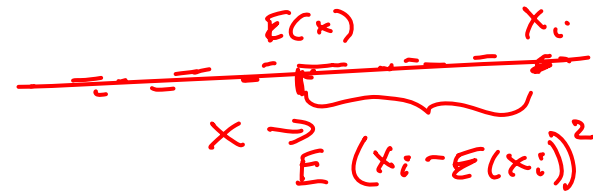
$$\begin{aligned} E \left[\frac{\sum (x_i - \bar{x})^2}{(N-1)} \frac{(N-1)}{N} \right] &= \text{Cov}(x, u) \\ &= V(x) \left(\frac{N-1}{N} \right) \end{aligned}$$

$$\text{plim} \left(\frac{N-1}{N} \right) = 1$$

$$\text{plim} \left(\frac{N-2}{N} \right) = 1$$

2. What's a Covariance?

VARIANCE: $E(x_i - E(x))^2$

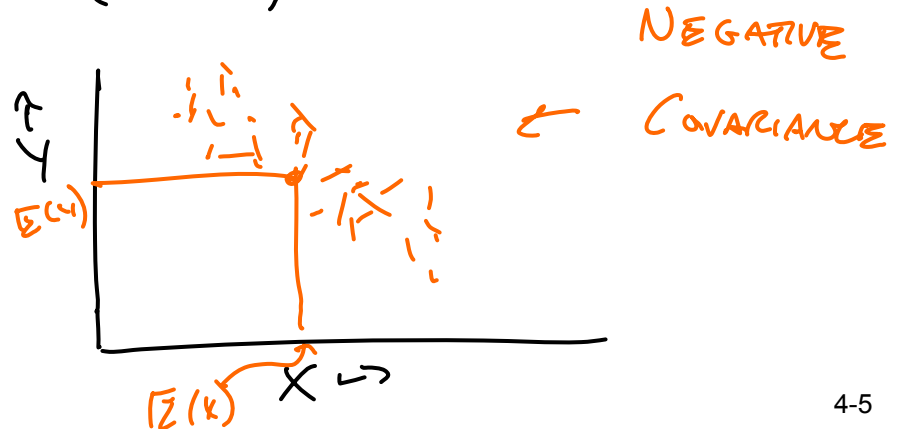
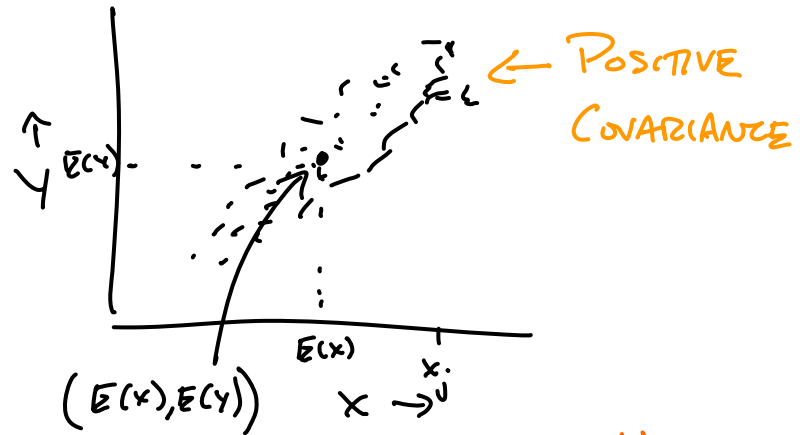


Covariance: $E([x_i - E(x)][y_i - E(y)])$

SAMPLE COVARIANCE OF X AND Y

(under random sampling) $\frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{(N-2)}$

$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$



3. Sampling Properties of OLS estimators (cont.)

Large-Sample Distributions of $\hat{\beta}_0$ and $\hat{\beta}_1$

If the least squares assumptions in Key Concept 4.3 hold, then in large samples $\hat{\beta}_0$ and $\hat{\beta}_1$ have a jointly normal sampling distribution. The large-sample normal distribution of $\hat{\beta}_1$ is $N(\beta_1, \sigma_{\hat{\beta}_1}^2)$, where the variance of this distribution, $\sigma_{\hat{\beta}_1}^2$, is

$$\hat{V}(b_1) = \frac{\sum (x_i - \bar{x})^2 e_i^2 / n}{\left[\frac{\sum (x_i - \bar{x})^2}{n-2} \right]^2} \left(\frac{1}{n} \right) \quad b_1 \sim N(\beta_1, \sigma_{\hat{\beta}_1}^2) \quad \hookrightarrow \quad \sigma_{\hat{\beta}_1}^2 = \frac{1}{n} \frac{\text{var}[(X_i - \mu_X)u_i]}{[\text{var}(X_i)]^2} = V(b_1) \quad (4.14)$$

The large-sample normal distribution of $\hat{\beta}_0$ is $N(\beta_0, \sigma_{\hat{\beta}_0}^2)$, where

We like precision
So we like
BIG n

$$\sigma_{\hat{\beta}_0}^2 = \frac{1}{n} \frac{\text{var}(H_i u_i)}{[E(H_i^2)]^2}, \quad \text{where } H_i = 1 - \left(\frac{\mu_X}{E(X_i^2)} \right) X_i. \quad (4.15)$$

Why asymptotically $N(\cdot)$?

$$b_i = \frac{\sum (x_i - \bar{x}) y_i}{\sum (x_i - \bar{x})^2} = \frac{1}{\sum (x_i - \bar{x})^2} \sum (x_i - \bar{x}) [\beta_0 + \beta_1 x_i + u_i] = 0 + \beta_1 + \frac{\sum (x_i - \bar{x}) u_i}{\sum (x_i - \bar{x})^2}$$

WHAT'S THE SAMPLING DISTRIBUTION OF b_i ?

$$v_i = (x_i - \bar{x}) u_i \quad E(v_i) = 0 \quad E(v_i | X) = E((x_i - \bar{x}) u_i | X) = (x_i - \bar{x}) E(u_i | X) = 0$$

v_i iid? ✓

$V(v_i) = V[(x_i - \bar{x}) u_i]$ is finite. ✓

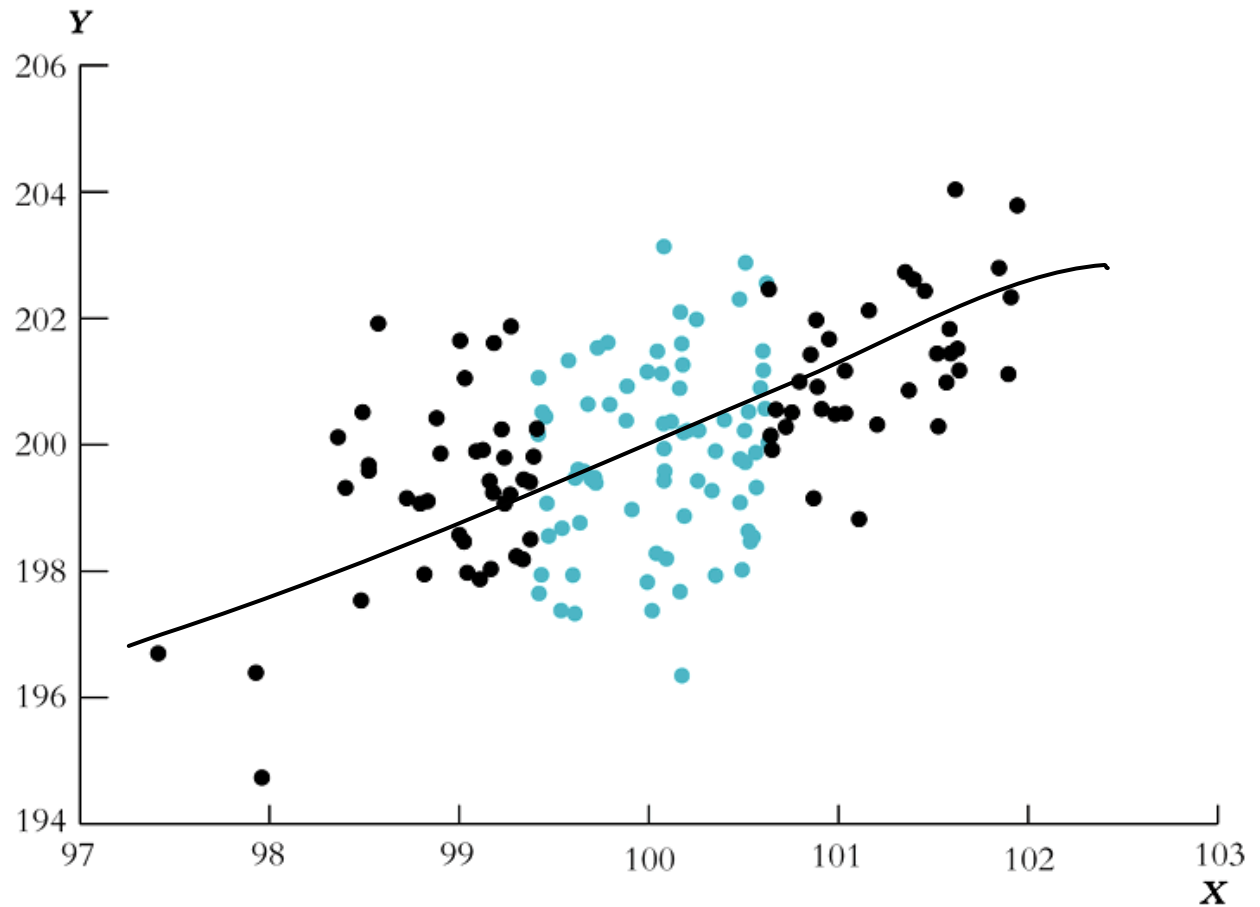
$$\sum v_i / n = \bar{v} \stackrel{A}{\sim} N(0, V(\bar{v}))$$

$$\begin{aligned} b_i &= \beta_1 + \frac{\sum (x_i - \bar{x}) u_i}{\sum (x_i - \bar{x})^2} \\ &= \beta_1 + \frac{\sum v_i}{\sum (x_i - \bar{x})^2} = \beta_1 + \frac{\sum v_i / n}{\sum (x_i - \bar{x})^2 / n} \end{aligned}$$

$$\Rightarrow b_i \stackrel{A}{\sim} N(\beta_1, V(b_i))$$

FIGURE 4.5 The Variance of $\hat{\beta}_1$ and the Variance of X

The colored dots represent a set of X_i 's with a small variance. The black dots represent a set of X_i 's with a large variance. The regression line can be estimated more accurately with the black dots than with the colored dots.



$$Q \sim N(\mu_Q, V(Q))$$

4. Standard Errors

$$Z = \frac{Q - \mu_Q}{\sqrt{V(Q)}} \sim N(0,1)$$

- CLT: b_1 approx $\sim N(\beta_1, V(b_1))$

$$b_0 \text{ approx } \sim N(\beta_0, V(b_0))$$

$$t = \frac{Q - \mu_Q}{\sqrt{\hat{V}(Q)}} \overset{\Delta}{\sim} N(0,1)$$

- **“Standard Errors” are consistent estimators of standard deviations of b_0 and b_1 . (S&W p. 133, 151, 180).**

- **So “t-statistics” have a standard normal distribution.**
- $$t = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\hat{V}(b_1)}} = \frac{\hat{\beta}_1 - \beta_1}{\text{Std. error}(b_1)} \overset{\Delta}{\sim} N(0,1)$$

5. Confidence Intervals

Confidence Intervals for β_1

A 95% two-sided confidence interval for β_1 is an interval that contains the true value of β_1 with a 95% probability, that is, it contains the true value of β_1 in 95% of all possible randomly drawn samples. Equivalently, it is also the set of values of β_1 that cannot be rejected by a 5% two-sided hypothesis test. When the sample size is large, it is constructed as

$$95\% \text{ confidence interval for } \beta_1 = (\hat{\beta}_1 - 1.96SE(\hat{\beta}_1), \hat{\beta}_1 + 1.96SE(\hat{\beta}_1)). \quad (4.27)$$

..and the same for β_0

Confidence Intervals

$$b_i \overset{A}{\sim} N(\beta_i, V(b_i)) \text{ by CLT}$$

Confidence intervals for β_0, β_1

$$P \left(-1.96 \leq \left(\frac{b_i - \beta_i}{\text{se.}(b_i)} \right) \leq 1.96 \right) \approx 0.95$$

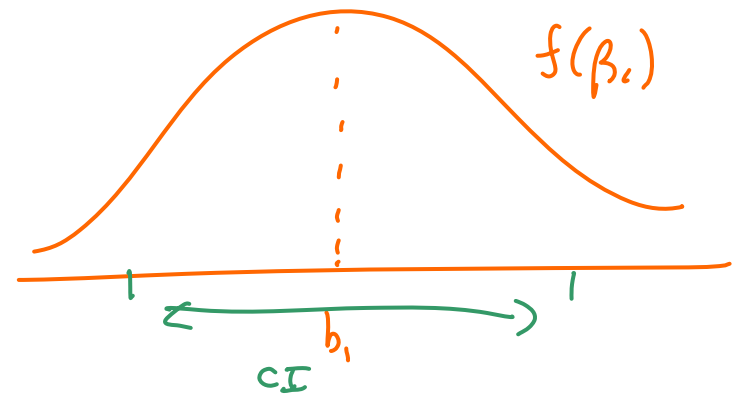
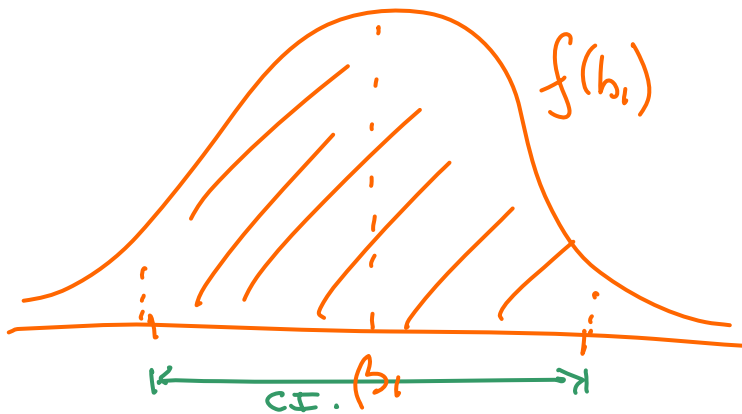
$$P \left[-1.96 \text{ se.}(b_i) \leq b_i - \beta_i \leq 1.96 \text{ s.e.}(b_i) \right] \approx 0.95$$

⇓
PROBABILITY.

$$P \left[\beta_i - 1.96 \text{ se.}(\cdot) \leq b_i \leq \beta_i + 1.96 \text{ se.}(\cdot) \right] \approx 0.95$$

⇒ STATISTICS

$$P \left[b_i - 1.96 \text{ se.}(\cdot) \leq \beta_i \leq b_i + 1.96 \text{ se.}(\cdot) \right] \approx 0.95$$

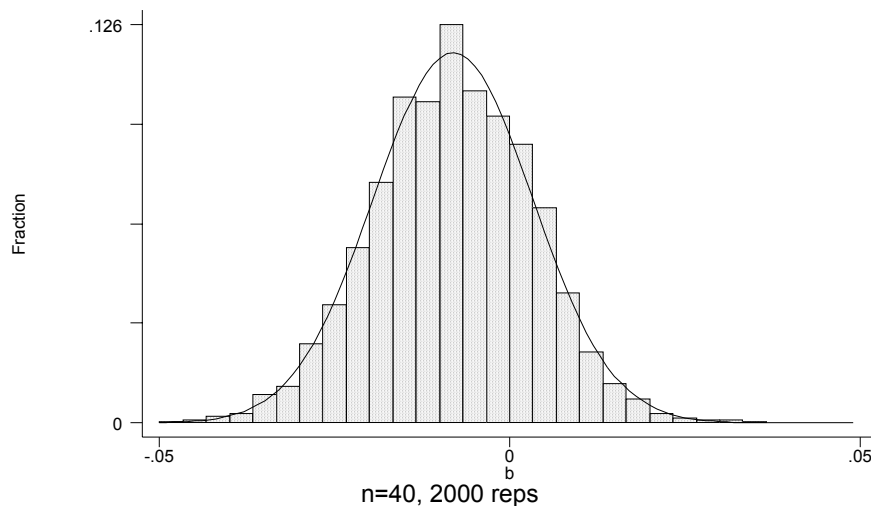
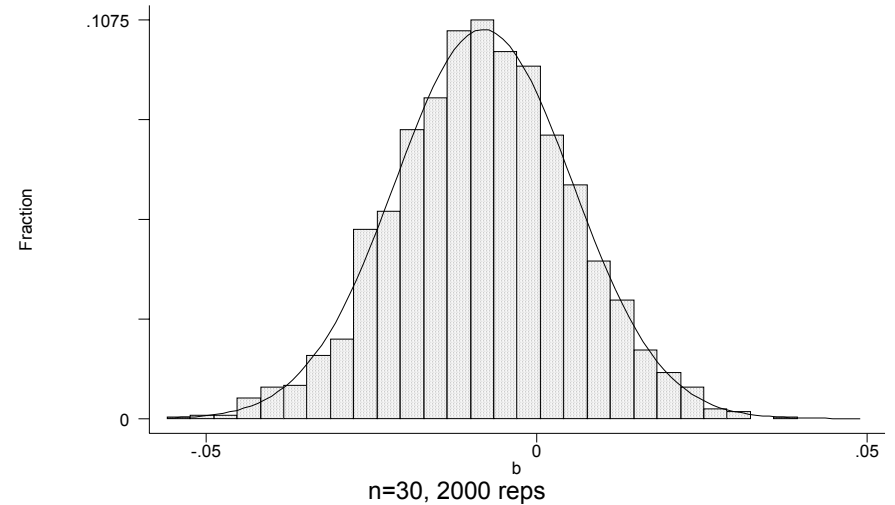
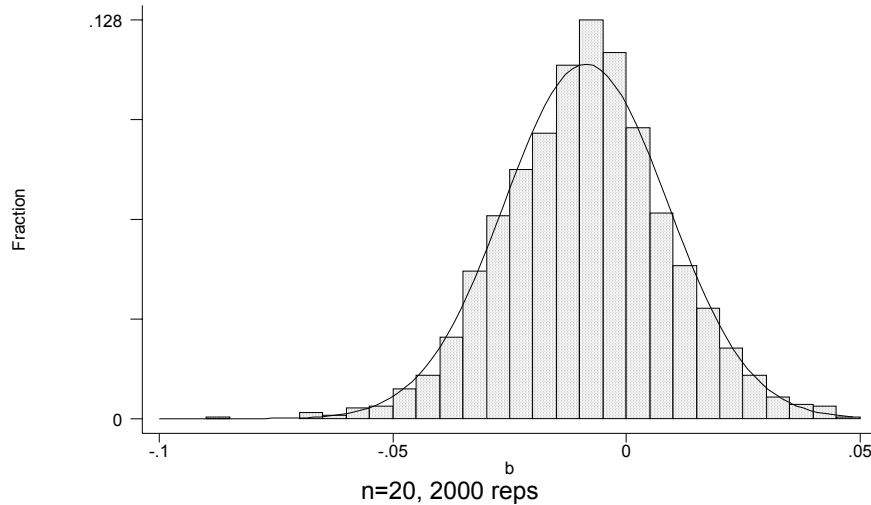


6. CLT demonstrations using Stata

- Stata doesn't mind running a regression a few thousand times,
 - which allows us to observe a sampling distribution for \mathbf{b}_1e.g., bootreg00.do

..and the same for \mathbf{b}_0

CLT in action: sampling distributions for b_1



Note: In developing countries
this slope is about -0.2 children
per year of education. V_g

7. Omitted Variables and Omitted Variable Bias

- **What if you left out an important variable?**
- **Many interesting relationships have more than 2 dimensions**
- **Multivariate regression:**



7a. OLS Multivariate regression

The OLS Estimators, Predicted Values, and Residuals in the Multiple Regression Model

The OLS estimators $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$ are the values of b_0, b_1, \dots, b_k that minimize the sum of squared prediction mistakes $\sum_{i=1}^n (Y_i - b_0 - b_1 X_{1i} - \dots - b_k X_{ki})^2$. The OLS predicted values \hat{Y}_i and residuals \hat{u}_i are:

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \dots + \hat{\beta}_k X_{ki}, \quad i = 1, \dots, n, \text{ and} \quad (5.11)$$

$$\hat{u}_i = Y_i - \hat{Y}_i, \quad i = 1, \dots, n. \quad (5.12)$$

The OLS estimators $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$ and residual \hat{u}_i are computed from a sample of n observations of $(X_{1i}, \dots, X_{ki}, Y_i)$, $i = 1, \dots, n$. These are estimators of the unknown true population coefficients $\beta_0, \beta_1, \dots, \beta_k$ and error term, u_i .

Look familiar? Same criterion with more variables.

7b. Properties of OLS estimators in Multivariate Regression

- Consistent
- Unbiased
- Approximately $N(\cdot)$ in large samples
- Claim:

7c. Omitted Variable “Bias”

- Short regression

$$y = b_0^s + b_1^s x_1 + e^s \quad (\text{SR})$$

- Long regression

$$y = b_0^L + b_1^L x_1 + b_2^L x_2 + e^L \quad (\text{LR})$$

- Claim:

$$b_1^s = b_1^L + b_2^L b_{21},$$

b_{21} is slope of a regression of x_2 on x_1