

A Primer in Econometric Theory

Lecture 11: Ordinary Least Squares

John Stachurski

Lectures by Akshay Shanker

March 26, 2017

Ordinary Least Squares

In chapter 11 we imposed only mild regularity conditions on the data

To provide additional interpretation of the estimated coefficients, we now make more assumptions — the classical OLS assumptions

The first assumption is a repetition of assumption 11.1.1 from chapter 11

Assumption.(12.1.1) \mathbf{X} has full column rank with probability one

Assumption.(12.1.2) The observations satisfy $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$ for some unknown K -vector of parameters $\boldsymbol{\beta}$ and unobservable vector of shocks \mathbf{u}

Assumption.(12.1.3) $\mathbb{E}[\mathbf{u} | \mathbf{X}] = \mathbf{0}$

Assumption.(12.1.4) $\mathbb{E}[\mathbf{u}\mathbf{u}^\top | \mathbf{X}] = \sigma^2\mathbf{I}$ for some unknown $\sigma > 0$

Assumption 12.1.2 can be decomposed into the separate equations

$$y_n = \mathbf{x}_n^\top \boldsymbol{\beta} + u_n, \quad n = 1, \dots, N \quad (1)$$

Examples of models that produce relationships in the form of (1):

Example. The **Cobb–Douglas production function** relates capital and labor inputs with output via $y = Ak^a\ell^\delta$, where A is a random, firm-specific productivity term and a and δ are parameters

Taking logs yields the linear regression model

$$\ln y = \gamma + a \ln k + \delta \ln \ell + u$$

where the random term $\ln A$ is represented by $\gamma + u$

Example. The **gravity model** relates international trade flows between country ℓ and country n via the equation

$$T_{\ell n} = \lambda \xi_{\ell n} G_{\ell}^{\alpha} G_n^{\beta} / D_{\ell n}^{\gamma}$$

- $T_{\ell n}$ is exports from country ℓ to country n
- λ is a constant term
- $\xi_{\ell n}$ is a shock
- G_{ℓ} and G_n are GDP in country ℓ and n respectively and $D_{\ell n}$ is distance between them

Taking logs gives

$$\ln T_{\ell n} = \ln \lambda + \alpha \ln G_{\ell} + \beta \ln G_n - \gamma \ln D_{\ell n} + \ln \xi_{\ell n} \quad (2)$$

Fact. (12.1.1) If the linearity assumption 12.1.2 holds, then

1. $\mathbf{My} = \mathbf{Mu}$
2. $\mathbf{Py} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Pu}$
3. $\text{RSS} = \mathbf{u}^\top \mathbf{Mu}$

For a proof, see exercise 12.4.2

Fact. (12.1.2) If assumption 12.1.3 holds, then

1. $\mathbb{E}[\mathbf{u}] = \mathbf{0}$
2. $\mathbb{E}[u_m | x_{nk}] = 0$ for any m, n, k
3. $\mathbb{E}[u_m x_{nk}] = 0$ for any m, n, k (orthogonality)
4. $\text{cov}[u_m, x_{nk}] = 0$ for any m, n, k

For a proof, see exercise 12.4.3

Fact. (12.1.3) If assumption 12.1.4 holds, then

1. $\text{var}[\mathbf{u}] = \mathbb{E}[\mathbf{u}\mathbf{u}^T] = \sigma^2\mathbf{I}$,
2. $\mathbb{E}[u_i^2 | \mathbf{X}] = \mathbb{E}[u_j^2 | \mathbf{X}] = \sigma^2$ for any i, j in $1, \dots, N$, and
3. $\mathbb{E}[u_i u_j | \mathbf{X}] = 0$ whenever $i \neq j$.

Parts 1 and 2 are called **homoskedasticity** and zero correlation respectively

Combining assumptions 12.1.3 and 12.1.4 gives

$$\text{var}[\mathbf{u} \mid \mathbf{X}] := \mathbb{E}[\mathbf{u}\mathbf{u}^T \mid \mathbf{X}] - \mathbb{E}[\mathbf{u} \mid \mathbf{X}]\mathbb{E}[\mathbf{u}^T \mid \mathbf{X}] = \mathbb{E}[\mathbf{u}\mathbf{u}^T \mid \mathbf{X}]$$

Thus the conditional variance–covariance matrix is the diagonal matrix $\sigma^2\mathbf{I}$

The standard estimator of β in (1) is the least squares estimator $\hat{\beta}$ defined by Equation (11.9) in ET:

$$\hat{\beta} := (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (3)$$

Unlike chapter 11, $\hat{\beta}$ now understood as an estimator of the unknown parameter vector β

Substituting $\mathbf{y} = \mathbf{X}\beta + \mathbf{u}$ into (3) and cancelling terms gives the useful comparison

$$\hat{\beta} = \beta + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{u} \quad (4)$$

$\hat{\beta}$ is also called the **OLS estimator** of β

The usual OLS estimator of the parameter σ^2 introduced in assumption 12.1.4:

$$\hat{\sigma}^2 := \frac{\text{RSS}}{N - K} \quad (5)$$

Example that runs the gravity model regression (2) on world trade data using Python — full set of code and data can be found at johnstachurski.net/emet

First import some Python libraries often used for statistics:

```
import pandas as pd
import statsmodels.formula.api as smf
```

Rad in the data to a pandas DataFrame from local CSV file:

```
data = pd.read_csv("./data/gravity_dataset_2013.csv")
```

Build a model using a formula to indicate the regression we want to run symbolically (similar to R):

```
formula = "log(value) ~ log(egdp) \
          + log(igdp) + log(dist)"
model = smf.ols(formula, data)
```

Perform the estimation and print a table summarizing results:

```
result = model.fit(cov_type='HC1')
print(result.summary())
```

OLS Regression Results

```

=====
Dep. Variable:          log(value)  R-squared:          0.652
Model:                  OLS        Adj. R-squared: 0.652
Method:                 Least Squares  F-statistic:      1.203e+04
No. Observations:      19655       Prob (F-statistic): 0.00
Df Residuals:          19651       Log-Likelihood: -47185
Df Model:               3          AIC:               9.438e+04
Covariance Type:       HC1         BIC:               9.441e+04\docume
=====

```

	coef	std err	z	P> z
Intercept	-30.2350	0.394	-76.773	0.000
log(egdp)	1.2783	0.008	153.772	0.000
log(igdp)	1.0287	0.009	118.885	0.000
log(dist)	-1.3483	0.023	-58.113	0.000

Let's check the values for $\hat{\beta}$ reported under `coef` coincide with the expression for $\hat{\beta}$ given in (3)

First build the design matrix \mathbf{X} (see the URL above for details) and then compute as follows:

```
betahat = inv(X.T @ X) @ X.T @ y
```

The output agrees with the output in the table

```
[-30.23498073    1.27825004    1.02865139   -1.34830012]
```

A **partial regression plot** for the linear model just described:

```
import matplotlib.pyplot as plt
import statsmodels.api as sm
fig = plt.figure()
fig = sm.graphics.plot_partregress_grid(result, fig=fig)
```

- horizontal axis — residuals from regressing $\log(\text{igdp})$ on all other columns in \mathbf{X}
- vertical axis — residuals from regressing the dependent variable $\log(\text{value})$ on all other columns in \mathbf{X}

(Refer to discussion surrounding equations (11.26) and (11.27) in ET)

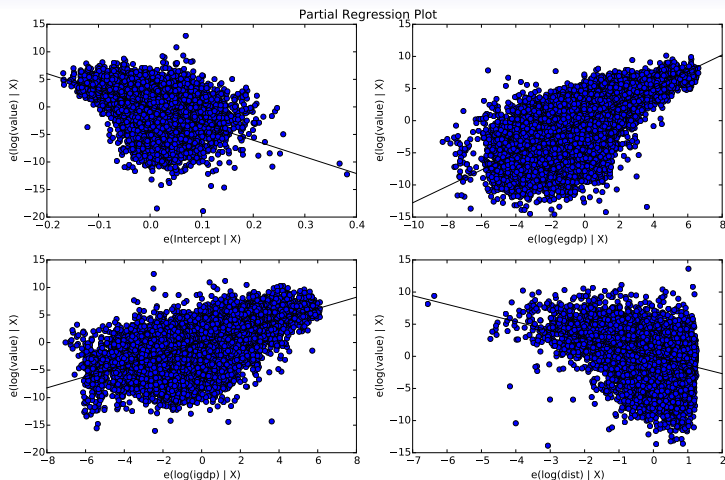


Figure: Partial regression plot for the gravity model

Finite Sample Properties

Theorem. (12.1.1) If assumptions 12.1.1–12.1.3 hold, then

$$\mathbb{E} \hat{\boldsymbol{\beta}} = \mathbb{E} [\hat{\boldsymbol{\beta}} | \mathbf{X}] = \boldsymbol{\beta} \quad (6)$$

If assumption 12.1.4 also holds, then

$$\mathbb{E} \hat{\sigma}^2 = \mathbb{E} [\hat{\sigma}^2 | \mathbf{X}] = \sigma^2 \quad (7)$$

Proof. By (12.4) and assumption 12.1.3, we have

$$\mathbb{E}[\hat{\boldsymbol{\beta}} | \mathbf{X}] = \boldsymbol{\beta} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}[\mathbf{u} | \mathbf{X}] = \boldsymbol{\beta}$$

Taking the unconditional expectation gives (6)

Regarding (7), we use the expression for RSS in fact 12.1.1:

$$\mathbb{E}[\text{RSS} | \mathbf{X}] = \mathbb{E}[\mathbf{u}^T \mathbf{M} \mathbf{u} | \mathbf{X}] = \sigma^2 \mathbb{E}[\boldsymbol{\zeta}^T \mathbf{M} \boldsymbol{\zeta} | \mathbf{X}] \quad \text{where} \quad \boldsymbol{\zeta} := \sigma^{-1} \mathbf{u}$$

Proof.(cont.) From fact 5.1.3, $\mathbb{E}[\text{RSS} \mid \mathbf{X}] = \sigma^2 \text{trace } \mathbf{M}$

Hence

$$\mathbb{E}[\hat{\sigma}^2 \mid \mathbf{X}] = \frac{\sigma^2 \text{trace}(\mathbf{M})}{N - K}$$

By fact 3.3.4, we have $\text{trace } \mathbf{M} = N - K$

Inserting this into the preceding equation and taking unconditional expectations gives (7) \square

Theorem. (12.1.2) If assumptions 12.1.1–12.1.4 hold, then

$$\text{var}[\hat{\beta} | \mathbf{X}] = \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1} \quad (8)$$

For a proof, see exercise 12.4.5

Under the stated conditions, the OLS estimator $\hat{\beta}$ is best linear unbiased:

Theorem. (12.1.3) (Gauss-Markov) Let assumptions 12.1.1–12.1.4 hold and let \mathbf{b} be an estimator of β . If \mathbf{b} is linear and unbiased, then

$$\text{var}[\hat{\beta} | \mathbf{X}] \leq \text{var}[\mathbf{b} | \mathbf{X}]$$

Here $\text{var}[\hat{\boldsymbol{\beta}} | \mathbf{X}] \leq \text{var}[\mathbf{b} | \mathbf{X}]$ means $\text{var}[\mathbf{b} | \mathbf{X}] - \text{var}[\hat{\boldsymbol{\beta}} | \mathbf{X}]$ is nonnegative definite — one way to assert $\text{var}[\mathbf{b} | \mathbf{X}]$ is “larger”

- an implication: $\text{var}[b_k | \mathbf{X}] \geq \text{var}[\hat{\beta}_k | \mathbf{X}]$ for all k

Linearity of \mathbf{b} means \mathbf{b} is linear as a function of \mathbf{y} , taking \mathbf{X} as fixed

- equivalent to requiring $\mathbf{b} = \mathbf{C}\mathbf{y}$ for some matrix \mathbf{C} ; \mathbf{C} can depend on \mathbf{X} but not \mathbf{y}

Saying \mathbf{b} is unbiased means $\mathbb{E}[\mathbf{b} | \mathbf{X}] = \mathbb{E}[\mathbf{C}\mathbf{y} | \mathbf{X}] = \boldsymbol{\beta}$ for all $\boldsymbol{\beta} \in \mathbb{R}^K$

Proof.[proof for theorem (12.1.3)]

Let $\mathbf{b} = \mathbf{C}\mathbf{y}$, as described above, and let $\mathbf{D} := \mathbf{C} - \mathbf{A}$, where $\mathbf{A} := (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$

Then

$$\mathbf{b} = \mathbf{C}\mathbf{y} = \mathbf{D}\mathbf{y} + \mathbf{A}\mathbf{y} = \mathbf{D}(\mathbf{X}\boldsymbol{\beta} + \mathbf{u}) + \hat{\boldsymbol{\beta}} = \mathbf{D}\mathbf{X}\boldsymbol{\beta} + \mathbf{D}\mathbf{u} + \hat{\boldsymbol{\beta}} \quad (9)$$

Take conditional expectations and use the fact that \mathbf{D} is a function of \mathbf{X} :

$$\begin{aligned}\mathbb{E}[\mathbf{b} \mid \mathbf{X}] &= \mathbb{E}[\mathbf{D}\mathbf{X}\boldsymbol{\beta} \mid \mathbf{X}] + \mathbb{E}[\mathbf{D}\mathbf{u} \mid \mathbf{X}] + \mathbb{E}[\hat{\boldsymbol{\beta}} \mid \mathbf{X}] \\ &= \mathbf{D}\mathbf{X}\mathbb{E}[\boldsymbol{\beta} \mid \mathbf{X}] + \mathbf{D}\mathbb{E}[\mathbf{u} \mid \mathbf{X}] + \mathbb{E}[\hat{\boldsymbol{\beta}} \mid \mathbf{X}] \\ &= \mathbf{D}\mathbf{X}\boldsymbol{\beta} + \mathbf{0} + \boldsymbol{\beta}\end{aligned}$$

Since \mathbf{b} is unbiased and, in particular, $\mathbb{E}[\mathbf{b} \mid \mathbf{X}] = \boldsymbol{\beta}$ for any given $\boldsymbol{\beta}$, we have

$$\boldsymbol{\beta} = \mathbf{DX}\boldsymbol{\beta} + \boldsymbol{\beta} \quad \text{for all } \boldsymbol{\beta} \in \mathbb{R}^K$$

$$\therefore \mathbf{0} = \mathbf{DX}\boldsymbol{\beta} \quad \text{for all } \boldsymbol{\beta} \in \mathbb{R}^K$$

In light of exercise 3.5.13, we conclude $\mathbf{DX} = \mathbf{0}$

Combining this result with (9), we obtain $\mathbf{b} = \mathbf{Du} + \hat{\boldsymbol{\beta}}$

Hence \mathbf{b} is equal to the OLS estimator plus zero-mean noise

To complete the proof, observe

$$\begin{aligned}\text{var}[\mathbf{b} \mid \mathbf{X}] &= \text{var}[\mathbf{D}\mathbf{u} + \hat{\boldsymbol{\beta}} \mid \mathbf{X}] \\ &= \text{var}[(\mathbf{D} + \mathbf{A})\mathbf{u} \mid \mathbf{X}] = (\mathbf{D} + \mathbf{A}) \text{var}[\mathbf{u} \mid \mathbf{X}] (\mathbf{D} + \mathbf{A})^\top\end{aligned}$$

Using assumption 12.1.4 and fact 3.2.4, the right-hand side of this expression becomes

$$\sigma^2(\mathbf{D} + \mathbf{A})(\mathbf{D}^\top + \mathbf{A}^\top) = \sigma^2(\mathbf{D}\mathbf{D}^\top + \mathbf{D}\mathbf{A}^\top + \mathbf{A}\mathbf{D}^\top + \mathbf{A}\mathbf{A}^\top)$$

Since

$$\mathbf{D}\mathbf{A}^\top = \mathbf{D}\mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1} = \mathbf{0}(\mathbf{X}^\top\mathbf{X})^{-1} = \mathbf{0}$$

and

$$\mathbf{A}\mathbf{A}^\top = (\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1} = (\mathbf{X}^\top\mathbf{X})^{-1}$$

we conclude

$$\text{var}[\mathbf{b} \mid \mathbf{X}] = \sigma^2[\mathbf{D}\mathbf{D}^\top + (\mathbf{X}^\top\mathbf{X})^{-1}] = \sigma^2\mathbf{D}\mathbf{D}^\top + \text{var}[\hat{\boldsymbol{\beta}} \mid \mathbf{X}]$$

The matrix $\sigma^2\mathbf{D}\mathbf{D}^\top$ is nonnegative definite, so the proof is now done \square

Precision of Estimators

By theorem 12.1.2, the variance — covariance matrix of $\hat{\beta}$ given \mathbf{X} is $\sigma^2(\mathbf{X}^T\mathbf{X})^{-1}$

- scalar variances of the individual OLS coefficient estimate $\hat{\beta}_1, \dots, \hat{\beta}_K$ given by the principal diagonal of this matrix

Since each $\hat{\beta}_k$ is unbiased (theorem 12.1.1), small variance means probability mass concentrated around the true parameter β_k — the estimator has high **precision**

Gauss–Markov Theorem: the OLS estimates will have highest precision

Hold the estimation technique fixed, as well as the sample size, and vary the application

- which problems will have high precision estimates?
- which will have low precision estimates?

Consider variance of a fixed coefficient β_k

Write $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$ as

$$\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \text{col}_k(\mathbf{X})\beta_k + \mathbf{u} \quad (10)$$

- $\text{col}_k(\mathbf{X})$ is the vector of observations of the k th regressor
- \mathbf{X}_1 contains as its columns the observations of the other regressors
- $\hat{\boldsymbol{\beta}}_1$ is the OLS estimates of the corresponding coefficients

By the FWL theorem:

$$\hat{\beta}_k = (\text{col}_k(\mathbf{X})^\top \mathbf{M}_1 \text{col}_k(\mathbf{X}))^{-1} \text{col}_k(\mathbf{X})^\top \mathbf{M}_1 \mathbf{y} \quad (11)$$

- \mathbf{M}_1 is the residual projection $\mathbf{M}_1 := \mathbf{I} - \mathbf{P}_1$
- $\mathbf{P}_1 := \mathbf{X}_1(\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top$ projects onto colspace \mathbf{X}_1

Apply \mathbf{M}_1 to both sides of (10)

$$\mathbf{M}_1 \mathbf{y} = \mathbf{M}_1 \text{col}_k(\mathbf{X}) \beta_k + \mathbf{M}_1 \mathbf{u}$$

Substituting into (11) gives

$$\hat{\beta}_k = \beta_k + (\text{col}_k(\mathbf{X})^\top \mathbf{M}_1 \text{col}_k(\mathbf{X}))^{-1} \text{col}_k(\mathbf{X})^\top \mathbf{M}_1 \mathbf{u} \quad (12)$$

Some calculations show (ex. 12.4.11)

$$\text{var}[\hat{\beta}_k | \mathbf{X}] = \sigma^2 (\text{col}_k(\mathbf{X})^\top \mathbf{M}_1 \text{col}_k(\mathbf{X}))^{-1} = \sigma^2 \|\mathbf{M}_1 \text{col}_k(\mathbf{X})\|^{-2}$$

Variance of $\hat{\beta}_k$ depends on:

1. variance σ^2 of the shock u — unavoidable
2. norm of the vector $\mathbf{M}_1 \text{col}_k \mathbf{X}$

$\mathbf{M}_1 \text{col}_k \mathbf{X}$ is the residuals from regressing $\text{col}_k \mathbf{X}$ on \mathbf{X}_1

If $\|\mathbf{M}_1 \text{col}_k \mathbf{X}\|$ is small, then variance of $\hat{\beta}_k$ large

- small when $\text{col}_k \mathbf{X}$ is “almost” a linear combination of the other regressors and hence close to $\text{colspace } \mathbf{X}_1$

Then

$$\|\mathbf{M}_1 \text{col}_k \mathbf{X}\| = \|\text{col}_k \mathbf{X} - \mathbf{P}_1 \text{col}_k \mathbf{X}\| \approx 0$$

Situation referred to as **multicollinearity**

Figure to show effect of multicollinearity on the variance of OLS estimate $\hat{\beta}_2$

\mathbf{X} has two columns related by $\text{col}_2 \mathbf{X} = \delta \text{col}_1 \mathbf{X} + (1 - \delta)\mathbf{z}$

- \mathbf{z} is a vector of N independent draws from the standard normal distribution

Larger δ means more dependence between $\text{col}_1 \mathbf{X}$ and $\text{col}_2 \mathbf{X}$

Figure shows the distribution of $\hat{\beta}_2$ for different values of δ

- true parameter is $\beta_2 = 1$
- $\hat{\beta}_2$ is unbiased but its variance increases with δ

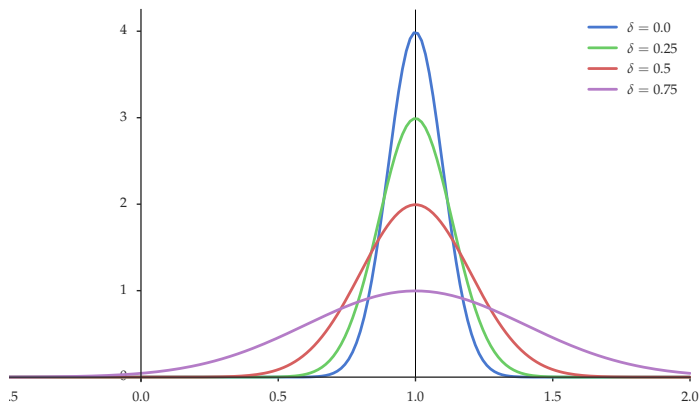


Figure: Distribution of $\hat{\beta}_2$ when $\text{col}_2 \mathbf{X} = \delta \text{col}_1 \mathbf{X} + (1 - \delta)\mathbf{z}$

Inference with Normal Errors

We want to compute confidence intervals or test hypotheses about the coefficients in finite samples (i.e., without appealing to asymptotics)

Assumption.(12.1.5) \mathbf{X} and \mathbf{u} are independent and $\mathcal{L}(\mathbf{u}) = \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$

Assumption 12.1.5 implies both assumption 12.1.3 and assumption 12.1.4

Theorem. (12.1.4) If assumptions 12.1.1–12.1.2 and 12.1.5 hold, then

$$\mathcal{L} \left[\hat{\boldsymbol{\beta}} \mid \mathbf{X} \right] = N(\boldsymbol{\beta}, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1})$$

The theorem implies the distribution of individual coefficient $\hat{\beta}_k$ given \mathbf{X} is also normal, with

$$\mathcal{L} [\hat{\beta}_k \mid \mathbf{X}] = \mathcal{L} \left[\mathbf{e}_k^\top \hat{\boldsymbol{\beta}} \mid \mathbf{X} \right] = N(\beta_k, \sigma^2 v_k(\mathbf{X})) \quad (13)$$

Here $v_k(\mathbf{X}) :=$ the (k, k) th element of $(\mathbf{X}^\top \mathbf{X})^{-1}$

(Recall (5.12) in ET: if $\mathbf{x} = (x_1, \dots, x_N)$ is multivariate normal, then

Regarding the distribution of $\hat{\sigma}^2$, convenient to work with a transformation:

$$Q := \frac{\text{RSS}}{\sigma^2} = (N - K) \frac{\hat{\sigma}^2}{\sigma^2}$$

Theorem. (12.1.5) If assumptions 12.1.1–12.1.2 and 12.1.5 hold, then

$$\mathcal{L}[Q | \mathbf{X}] = \chi^2(N - K)$$

See page 333 in ET for a proof

Fact. (12.1.4) If assumptions 12.1.1–12.1.2 and 12.1.5 hold, then the random elements $\hat{\sigma}^2$ and $\hat{\beta}$ are independent given \mathbf{X}

See ex. 12.4.6

The t -Test

Consider problem of testing hypothesis about an individual coefficient β_k

The null hypothesis

$$H_0: \beta_k = \beta_k^0$$

where β_k^0 is any number

If σ^2 is known, construct a test of H_0 based on:

$$z_k := \frac{\hat{\beta}_k - \beta_k}{\sigma \sqrt{v_k(\mathbf{X})}} \implies \mathcal{L}[z_k | \mathbf{X}] = N(0, 1) \quad (14)$$

We don't know σ^2 — replace σ^2 with its estimator $\hat{\sigma}^2$

Some notation:

$$\text{se}(\hat{\beta}_k) := \sqrt{\hat{\sigma}^2 v_k(\mathbf{X})}$$

The term $\text{se}(\hat{\beta}_k)$ is called the **standard error** of $\hat{\beta}_k$

Replace standard deviation with its sample estimate $\text{se}(\hat{\beta}_k)$ and β_k with β_k^0 in (14) and obtain the **t-statistic**

$$t_k := \frac{\hat{\beta}_k - \beta_k^0}{\text{se}(\hat{\beta}_k)}$$

associated with the hypothesis H_0

Theorem. (12.1.6) Let assumptions 12.1.1–12.1.2 and 12.1.5 hold. If H_0 is true, then

$$\mathcal{L}[t_k | \mathbf{X}] = \text{Student's } t \text{ with } N - K \text{ degrees of freedom}$$

See page 334 for a proof.

Let $T := |t_k|$ and let a desired size α for our test of H_0 be given

To generate a test of size α , we choose $c = c_\alpha$ to solve $\alpha = \mathbb{P}_\theta\{T > c\}$, or

$$1 - \alpha = \mathbb{P}_\theta\{|t_k| \leq c\}$$

The solution is $c_\alpha = F^{-1}(1 - \alpha/2)$

- F is the Student's t CDF with $N - K$ degrees of freedom (Equation (4.36) in ET)

The corresponding p -value is $2F(-|t_k|)$

Example. A common implementation of the t -test is the test that a given coefficient is equal to zero

For the k th coefficient β_k , this leads to the statistic

$$t_k := \frac{\hat{\beta}_k}{\text{se}(\hat{\beta}_k)}$$

This statistic is sometimes called the **Z-score**

The F -Test

Return to the setting of when we discussed the FWL theorem (§11.2.2):

Let \mathbf{y} and \mathbf{X} be given and let $\hat{\boldsymbol{\beta}}$ be the least squares estimator. Let K_1 be an integer with $1 \leq K_1 < K$, and let

- \mathbf{X}_1 be a matrix consisting of the first K_1 columns of \mathbf{X} ,
- \mathbf{X}_2 be a matrix consisting of the remaining $K_2 := K - K_1$ columns,
- $\hat{\boldsymbol{\beta}}_1$ be the $K_1 \times 1$ vector consisting of the first K_1 elements of $\hat{\boldsymbol{\beta}}$.
- $\hat{\boldsymbol{\beta}}_2$ be the $K_2 \times 1$ vector consisting of the remaining K_2 elements of $\hat{\boldsymbol{\beta}}$,
- $\mathbf{P}_1 := \text{proj}(\text{colspace } \mathbf{X}_1)$, and
- $\mathbf{M}_1 := \mathbf{I} - \mathbf{P}_1$ = the corresponding residual projection

Hypotheses concerning multiple regressors

Null hypothesis:

$$H_0: \beta_2 = \mathbf{0}$$

Since

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{u} \quad (15)$$

Under the null hypothesis,

$$\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{u} \quad (16)$$

Let

$$\text{USSR} := \|\mathbf{M}\mathbf{y}\|^2 \quad \text{and} \quad \text{RSSR} := \|\mathbf{M}_1\mathbf{y}\|^2$$

be the residual sums of squares for the unrestricted regression (15) and restricted regression (16)

Test statistic for our null hypothesis

$$F := \frac{(\text{RSSR} - \text{USSR})/K_2}{\text{USSR}/(N - K)} \quad (17)$$

Large residuals in the restricted regression (16) relative to those in (15) result in large values for F — evidence against the null hypothesis

Theorem. (12.1.7) Let assumptions 12.1.1–12.1.2 and 12.1.5 hold

If H_0 is true, then, given \mathbf{X} , the statistic F defined in (17) has the F distribution with parameters $(K_2, N - K)$

Proof.

Let $Q_1 := (\text{RSSR} - \text{USSR})/\sigma^2$ and let $Q_2 := \text{USSR}/\sigma^2$, so that

$$F = \frac{Q_1/K_2}{Q_2/(N - K)}$$

Recall that if x_1 and x_2 are independent with $\mathcal{L}(x_i) = \chi^2(k_i)$ for $i = 1, 2$, then

$\frac{x_1/k_1}{x_2/k_2}$ is F distributed with parameters k_1, k_2

Proof.(cont.)

Now suffices to show that, under the null hypothesis

- (a) Q_1 is chi-squared with K_2 degrees of freedom,
- (b) Q_2 is chi-squared with $N - K$ degrees of freedom, and
- (c) Q_1 and Q_2 are independent.

Part (b) was established in theorem 12.1.5

Regarding part (a), under the null hypothesis,

- $USSR = \|\mathbf{M}\mathbf{y}\|^2 = \|\mathbf{M}(\mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{u})\|^2 = \|\mathbf{M}\mathbf{u}\|^2 = \mathbf{u}^\top \mathbf{M}\mathbf{u}$ and
- $RSSR = \|\mathbf{M}_1\mathbf{y}\|^2 = \|\mathbf{M}_1(\mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{u})\|^2 = \|\mathbf{M}_1\mathbf{u}\|^2 = \mathbf{u}^\top \mathbf{M}_1\mathbf{u}$.

It follows that

$$RSSR - USSR = \mathbf{u}^\top \mathbf{M}_1\mathbf{u} - \mathbf{u}^\top \mathbf{M}\mathbf{u} = \mathbf{u}^\top (\mathbf{M}_1 - \mathbf{M})\mathbf{u}$$

Proof.(cont.)

Using the definitions of \mathbf{M} and \mathbf{M}_1 , we obtain

$$\begin{aligned} Q_1 &= \frac{\text{RSSR} - \text{USSR}}{\sigma^2} \\ &= \frac{\mathbf{u}^\top (\mathbf{I} - \mathbf{P}_1 - \mathbf{I} + \mathbf{P}) \mathbf{u}}{\sigma^2} = (\sigma^{-1} \mathbf{u})^\top (\mathbf{P} - \mathbf{P}_1) (\sigma^{-1} \mathbf{u}) \end{aligned}$$

Exercise: show $(\mathbf{P} - \mathbf{P}_1)$ is symmetric and idempotent. Hence

$$\begin{aligned} \text{rank}(\mathbf{P} - \mathbf{P}_1) &= \text{trace}(\mathbf{P} - \mathbf{P}_1) \\ &= \text{trace } \mathbf{P} - \text{trace } \mathbf{P}_1 = K - K_1 = K_2 \end{aligned}$$

Proof.(cont.) Via fact 5.1.19, we conclude $\mathcal{L}(Q_1) = \chi^2(K_2)$, as was to be shown

Now we show that, under the null hypothesis and taking \mathbf{X} as given, Q_1 and Q_2 are independent

Q_1 is a function of $(\mathbf{P} - \mathbf{P}_1)\mathbf{u}$, while Q_2 is a function of \mathbf{Mu}

- we show $(\mathbf{P} - \mathbf{P}_1)\mathbf{u}$ and \mathbf{Mu} are independent given \mathbf{X} by showing their covariance is zero

Observe

$$\begin{aligned}\text{cov}[(\mathbf{P} - \mathbf{P}_1)\mathbf{u}, \mathbf{Mu} \mid \mathbf{X}] \\ = \mathbb{E}[(\mathbf{P} - \mathbf{P}_1)\mathbf{u}(\mathbf{Mu})^\top \mid \mathbf{X}] = \mathbb{E}[(\mathbf{P} - \mathbf{P}_1)\mathbf{u}\mathbf{u}^\top \mathbf{M} \mid \mathbf{X}]\end{aligned}$$

Proof.(cont.) Since \mathbf{P} , \mathbf{P}_1 , and \mathbf{M} are just functions of \mathbf{X} , the above becomes

$$(\mathbf{P} - \mathbf{P}_1)\mathbb{E}[\mathbf{u}\mathbf{u}^\top | \mathbf{X}]\mathbf{M} = \sigma^2(\mathbf{P} - \mathbf{P}_1)\mathbf{M} = \sigma^2(\mathbf{P} - \mathbf{P}_1)(\mathbf{I} - \mathbf{P})$$

Using idempotence and fact 2.2.7, the matrix product on the right is

$$(\mathbf{P} - \mathbf{P}_1)(\mathbf{I} - \mathbf{P}) = \mathbf{P} - \mathbf{P}^2 - \mathbf{P}_1 + \mathbf{P}_1\mathbf{P} = \mathbf{P} - \mathbf{P} - \mathbf{P}_1 + \mathbf{P}_1 = \mathbf{0}$$

This completes the proof of independence, and hence of the theorem \square

Common F-Test: test that all coefficients of nonconstant regressors are zero

Consider:

$$\mathbf{y} = \mathbf{1}\beta_1 + \mathbf{X}_2\beta_2 + \mathbf{u}$$

where β_2 is the vector of coefficients corresponding to the nonconstant regressors. Null hypothesis:

$$\mathbf{y} = \mathbf{1}\beta_1 + \mathbf{u}$$

F statistic in (17) becomes

$$F = \frac{R_c^2}{1 - R_c^2} \frac{N - K}{K_2}$$

(proof is exercise — see page 338)

Why is large F is evidence against the null?

Nonspherical Errors

Heteroskedasticity occurs when the variance of the error term is not constant across observations

- violate assumption 12.1.4

Errors **heteroskedastic** if diagonal terms of $\mathbb{E}[\mathbf{u}\mathbf{u}^T | \mathbf{X}]$ not constant

If the off-diagonal terms are nonzero, then the errors are said to have **serial correlation**

Failure of assumption 12.1.4 does not alone cause bias in the OLS estimator $\hat{\beta}$

However, without assumption 12.1.4, no longer true that $\hat{\sigma}^2$ is unbiased for σ

- expression $\text{var}[\hat{\beta} | \mathbf{X}] = \sigma^2(\mathbf{X}^T \mathbf{X})^{-1}$ from (8) is no longer valid
- Gauss–Markov theorem also breaks down
- results on inference we discussed above from §12.1.4 no longer valid

Replace the assumption $\mathbb{E}[\mathbf{u}\mathbf{u}^\top | \mathbf{X}] = \sigma^2 \mathbf{I}$ with more general assumption

$$\mathbb{E}[\mathbf{u}\mathbf{u}^\top | \mathbf{X}] = \mathbf{\Omega} \quad (18)$$

for some positive definite matrix $\mathbf{\Omega}$

Write the variance of the OLS estimator $\hat{\boldsymbol{\beta}}$ under these conditions as

$$\text{var}[\hat{\boldsymbol{\beta}} | \mathbf{X}] = \mathbb{E}[(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^\top | \mathbf{X}] = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{\Omega} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}$$

Maintain assumptions 12.1.1–12.1.3

Write the variance of the OLS estimator $\hat{\beta}$ under these conditions as

$$\begin{aligned}\text{var}[\hat{\beta} | \mathbf{X}] &= \mathbb{E}[(\hat{\beta} - \beta)(\hat{\beta} - \beta)^\top | \mathbf{X}] \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{\Omega} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}\end{aligned}$$

If $\mathbf{\Omega}$ is known, then we can estimate using **generalized least squares**

- Cholesky decomposition: there exists a nonsingular matrix \mathbf{C} such that $\mathbf{C}^\top \mathbf{C} = \mathbf{\Omega}^{-1}$
- Use \mathbf{C} to transform the regression model by multiplying both sides of $\mathbf{y} = \mathbf{X}\beta + \mathbf{u}$ by \mathbf{C}

This gives

$$\mathbf{y}_c = \mathbf{X}_c \boldsymbol{\beta} + \mathbf{u}_c \quad \text{where} \quad \mathbf{y}_c := \mathbf{C}\mathbf{y}, \mathbf{X}_c := \mathbf{C}\mathbf{X} \text{ and } \mathbf{u}_c := \mathbf{C}\mathbf{u}$$

Exercise 12.4.17 asks you to show

$$\mathbb{E}[\mathbf{u}_c | \mathbf{X}_c] = \mathbf{0} \quad \text{and} \quad \mathbb{E}[\mathbf{u}_c \mathbf{u}_c^\top | \mathbf{X}_c] = \mathbf{I}$$

All OLS assumptions satisfied for the transformed data

Use the transformed data to estimate $\boldsymbol{\beta}$ via $(\mathbf{X}_c^\top \mathbf{X}_c)^{-1} \mathbf{X}_c^\top \mathbf{y}_c$ —
recover the best linear unbiased property

The new estimator:

$$\hat{\beta}_{\text{GLS}} := (\mathbf{X}^{\top} \mathbf{\Omega}^{-1} \mathbf{X})^{-1} \mathbf{X}^{\top} \mathbf{\Omega}^{-1} \mathbf{y} \quad (19)$$

- estimator is the **generalized least squares (GLS) estimator** of β

In practice, $\mathbf{\Omega}$ is rarely known

Estimation difficult unless additional structure is imposed

- $\mathbf{\Omega}$ is $N \times N$, implying that, absent further restrictions, the number of parameters contained in $\mathbf{\Omega}$ grows *faster* than the sample

Assume Ω is diagonal

- rules out correlation but maintains the possibility of heteroskedasticity

One estimator of Ω :

$$\hat{\Omega} := \text{diag}(\hat{u}_1^2, \dots, \hat{u}_N^2) \quad \text{where} \quad \hat{\mathbf{u}} := \mathbf{M}\mathbf{y}$$

Replace Ω with $\hat{\Omega}$ in expression for $\text{var}[\hat{\beta} | \mathbf{X}]$, get estimate of $\text{var}[\hat{\beta} | \mathbf{X}]$:

$$\hat{\mathbf{V}} := (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \hat{\Omega} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$$

Using this in place of the variance estimate $\hat{\sigma}^2 (\mathbf{X}^T \mathbf{X})^{-1}$ to form standard errors, we get the **heteroskedasticity-consistent standard errors** proposed by White (1980)

Bias

We will look at what happens when Assumptions 12.1.2 and 12.1.3 (linearity and exogeneity) fail

Linearity fundamental to OLS. Without linearity:

- suppose $y_n = f(\mathbf{x}_n) + u_n$ for some arbitrary function f
- cannot even *ask* if the OLS estimator is unbiased because no parameter vector for $\hat{\beta}$ to estimate

Endogeneity Bias

Bias due to failure of exogeneity assumption 12.1.3 is called **endogeneity bias**

Recalling Equation (4) at the start of the lecture,

$$\mathbb{E}[\hat{\beta} | \mathbf{X}] - \beta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}[\mathbf{u} | \mathbf{X}] \quad (20)$$

Without $\mathbb{E}[\mathbf{u} | \mathbf{X}] = \mathbf{0}$, we cannot assert right-hand side is zero, hence $\hat{\beta}$ is biased

Striving for unbiasedness somewhat misguided. Some degree of bias to

- penalize complexity
- recognize empirical distribution not the true distribution

However, good estimators reduce bias as the sample size increases, so that the object they aim to compute can be recovered asymptotically

- not true for the OLS estimator when exogeneity fails

Recall from theorem 11.1.2 that:

$$\hat{\boldsymbol{\beta}} \xrightarrow{p} \mathbb{E}[\mathbf{x}\mathbf{x}^T]^{-1} \mathbb{E}[\mathbf{x}y]$$

as $N \rightarrow \infty$

The right-hand side is the vector of coefficients in the best linear predictor

If $y = \mathbf{x}^T \boldsymbol{\beta} + u$, then

$$\hat{\boldsymbol{\beta}} \xrightarrow{p} \mathbb{E}[\mathbf{x}\mathbf{x}^T]^{-1} \mathbb{E}[\mathbf{x}(\mathbf{x}^T \boldsymbol{\beta} + u)] = \boldsymbol{\beta} + \mathbb{E}[\mathbf{x}\mathbf{x}^T]^{-1} \mathbb{E}[\mathbf{x}u]$$

With exogeneity we get $\mathbb{E}[\mathbf{x}u] = \mathbb{E}[\mathbf{x} \mathbb{E}[u | \mathbf{x}]] = 0$ and the OLS estimator is consistent — cannot assert without exogeneity

Examples of Exogeneity

Consider again the Cobb–Douglas production function which yields the regression model

$$\ln y_n = \gamma + a \ln k_n + \delta \ln \ell_n + u_n$$

- y is output, k is capital and ℓ is labor
- subscript n indicates observation on the n th firm and
- term u_n is a firm specific productivity shock

Source of endogeneity bias: firm chooses higher levels of capital and labor when it anticipates high productivity in the current period

To illustrate, suppose firm n received and observed $u_{n,-1}$ last period

Productivity follows random walk, with $u_n = u_{n,-1} + \eta_n$, where η_n is zero-mean white noise

Firms

- forecast period n productivity as $\mathbb{E}[u_n | u_{n,-1}] = u_{n,-1}$
- increase labor input when productivity is anticipated to be high, with the specific relationship $\ell_n = a + b\mathbb{E}[u_n | u_{n,-1}]$ for $b > 0$

With zero-mean shocks:

$$\mathbb{E}[\ell_n u_n] = \mathbb{E}[(a + bu_{n,-1})(u_{n,-1} + \eta_n)] = \mathbb{E}[bu_{n,-1}^2]$$

Strictly positive whenever $u_{n,-1}$ has positive variance

- exogeneity fails as conditions of fact 12.1.2 fail

Now consider data generated according to AR(1) model

$$x_0 = 0 \quad \text{and} \quad x_n = \beta x_{n-1} + u_n \quad \text{for } n = 1, \dots, N \quad (21)$$

Let $\{u_n\}_{n=1}^N$ be IID with distribution $N(0, \sigma^2)$

The unknown parameters are β and σ^2

Let $\mathbf{y} := (x_1, \dots, x_N)$, $\mathbf{x} := (x_0, \dots, x_{N-1})$, and $\mathbf{u} := (u_1, \dots, u_N)$,

Write N equations in (21) as $\mathbf{y} = \beta \mathbf{x} + \mathbf{u}$

The OLS estimate of β is $\hat{\beta} := (\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top \mathbf{y}$

If $n \geq m$, then we can write x_n as

$$x_n = \sum_{j=0}^{n-1} \beta^j u_{n-j}$$

Hence $\mathbb{E}[x_n u_m] = \sum_{j=0}^{n-1} \beta^j \mathbb{E}[u_{n-j} u_m] = \beta^{n-m} \sigma^2$

Once again exogeneity fails since conditions of fact 12.1.2 fail

```
import numpy as np; import matplotlib.pyplot as plt

N = 25; x = np.zeros(N); beta = 0.9
num_reps = 10000; betahat_obs = np.empty(num_reps)

for j in range(num_reps):
    u = np.random.randn(N)
    for t in range(N-1):
        x[t+1] = beta * x[t] + u[t+1]
    y = x[1:]          #  $x_1, \dots, x_N$ 
    x_vec = x[:-1]     #  $x_0, \dots, x_{N-1}$ 
    betahat_obs[j] = np.sum(x_vec * y) \
                    / np.sum(x_vec**2)

plt.hist(betahat_obs, bins=50, alpha=0.6, normed=True)
plt.show()
```

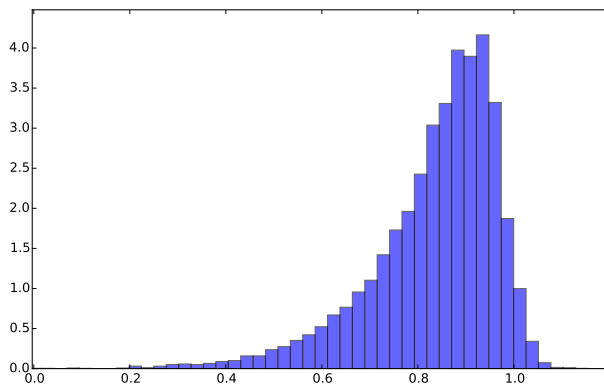


Figure: Observations of $\hat{\beta}$ when $\beta = 0.9$

Instrumental Variables

When regressors are endogenous, we can find consistent estimators of β when extra information is available in the form of “exogenous” variables called **instruments**

These collected into a $N \times J$ matrix \mathbf{Z} , where each column is observations on one exogenous variable

Assumptions:

1. $\mathbf{y} = \mathbf{X}\beta + \mathbf{u}$
2. $\mathbb{E}[\mathbf{u} | \mathbf{Z}] = \mathbf{0}$
3. $\mathbb{E}[\mathbf{u}\mathbf{u}^\top | \mathbf{Z}] = \sigma^2 \mathbf{I}$ for some positive constant σ
4. $\mathbf{Z}^\top \mathbf{X}$ has full column rank

Assumption 4. implies $J \geq K$

- If $J > K$, the model is said to be **overidentified**
- If $J = K$, the model is called **exactly identified**
- If $J < K$, the model is called **underidentified**

Multiply $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$ by \mathbf{Z}^T to produce

$$\mathbf{Z}^T \mathbf{y} = \mathbf{Z}^T \mathbf{X} \boldsymbol{\beta} + \mathbf{w} \quad \text{where} \quad \mathbf{w} := \mathbf{Z}^T \mathbf{u} \quad (22)$$

Recall the GLS estimator

$$\hat{\beta}_{\text{GLS}} := (\mathbf{X}^{\top} \mathbf{\Omega}^{-1} \mathbf{X})^{-1} \mathbf{X}^{\top} \mathbf{\Omega}^{-1} \mathbf{y}$$

Substitute \mathbf{X} for $\mathbf{Z}^{\top} \mathbf{X}$, \mathbf{y} for $\mathbf{Z}^{\top} \mathbf{y}$, and $\mathbf{\Omega}$ for $\sigma^2 \mathbf{Z}^{\top} \mathbf{Z}$

The **instrumental variable least squares (IVLS)** estimator

$$\hat{\beta}_{\text{IVLS}} := (\mathbf{X}^{\top} \mathbf{Z} (\mathbf{Z}^{\top} \mathbf{Z})^{-1} \mathbf{Z}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \mathbf{Z} (\mathbf{Z}^{\top} \mathbf{Z})^{-1} \mathbf{Z}^{\top} \mathbf{y}$$

If $J = K$, simplify to

$$\hat{\beta}_{\text{IVLS}} := (\mathbf{Z}^{\top} \mathbf{X})^{-1} \mathbf{Z}^{\top} \mathbf{y}$$

Cannot claim unbiasedness and other properties by appealing to the GLS theory — $\hat{\beta}_{IVLS}$ biased in general (see ex. 12.4.20)

- because direct application of GLS to (22) requires $\mathbb{E}[\mathbf{w} | \mathbf{Z}^T \mathbf{X}] = \mathbf{0}$ rather than $\mathbb{E}[\mathbf{w} | \mathbf{Z}] = \mathbf{0}$

However, bias typically smaller than OLS and vanishes asymptotically

Figure: OLS and IVLS estimates of β in $y = x\beta + u$ when $\beta = 10$:

- z drawn independently of u
- $x = \alpha z + (1 - \alpha)u$
- α is constant

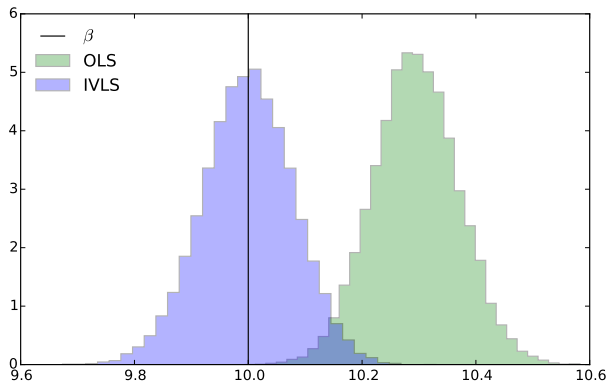


Figure: Simulated draws of OLS and IVLS estimators of β

Causality

Interpret R^2 as measuring the “explanatory power”?

- common terminology: “the total variation in \mathbf{y} is the sum of explained and unexplained variation”
- value of R^2 claimed to be the fraction of the variation in \mathbf{y} “explained” by the regressors

Above terminology misleading: R^2 says nothing about causation *per se*

- R^2 better thought of as a measure of correlation

Understanding causality requires either good experiments or good theory — no magic econometric procedure to reliably extract causality from observational