# Lecture 7 - OLS review

Thursday, September 09, 2021     3:58 PM

Lecture 7 -
OLS review

---

## OLS Review

Mauricio Romero

1

---

## OLS Review

Linear algebra review

Law of iterated expectations

OLS basics

Conditional expectation function

"Algebraic" properties of OLS

Properties of OLS estimators

Regression (matrix algebra) with a treatment dummy for the experimental case

Frisch–Waugh–Lovell (FWL) theorem

Regression and causality

2

---

## OLS Review

Linear algebra review

Law of iterated expectations

OLS basics

Conditional expectation function

"Algebraic" properties of OLS

Properties of OLS estimators

Regression (matrix algebra) with a treatment dummy for the experimental case

Frisch–Waugh–Lovell (FWL) theorem

Regression and causality

3

---

## Basic matrix operations

- $k \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} ka_1 \\ ka_2 \end{bmatrix}$

- $\begin{bmatrix} a & b \end{bmatrix} \begin{bmatrix} c \\ d \end{bmatrix} = \begin{bmatrix} ac + bd \end{bmatrix}$

- $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} a & b \\ c & d \end{bmatrix} = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$

- $\begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$

4

---

## Matrix multiplication

- Let $A_{n \times m}$ and $B_{m \times k}$, then $(AB)_{n \times k}$

- Let $A_{n \times m}$ and $B_{m \times k}$, then $(BA)$ "conformability error"

$$\left( B_{m \times k} \right) \left( A_{n \times m} \right)$$

5

## Transpose and inverse of a matrix

$$\left(A_{a\times m}\ B_{m\times k}\right)' = \left(AB_{a\times k}\right)' = [k\times a] \longrightarrow B'_{k\times m} A'_{m\times a}\ \checkmark$$

- Transpose of Product $(AB)' = B'A'$ and $(ABC)' = C'B'A'$

- Inverse of Product $(AB)^{-1} = B^{-1}A^{-1}$ and $(ABC)^{-1} = C^{-1}B^{-1}A^{-1}$

- Transpose of an inverse equals inverse of a transpose $(D^{-1})' = (D')^{-1}$

6

---

## OLS Review

Linear algebra review

Law of iterated expectations

OLS basics

Conditional expectation function

"Algebraic" properties of OLS

Properties of OLS estimators

Regression (matrix algebra) with a treatment dummy for the experimental case

Frisch–Waugh–Lovell (FWL) theorem

Regression and causality

7

---

## OLS Review

Linear algebra review

Law of iterated expectations

OLS basics

Conditional expectation function

"Algebraic" properties of OLS

Properties of OLS estimators

Regression (matrix algebra) with a treatment dummy for the experimental case

Frisch–Waugh–Lovell (FWL) theorem

Regression and causality

8

---

## Law of Iterated Expectations (LIE): A useful trick

- **Formally**: The unconditional expectation of a random variable is equal to the expectation of the conditional expectation of the random variable conditional on some other random variable

$$\mathbb{E}(Y) = \mathbb{E}(\mathbb{E}[Y|X])$$

- **Informally**: the weighted average of the conditional averages is the unconditional average

9

---

## Example of LIE

- Say want average wage but only know average wage by education level

- LIE says we get the former by taking conditional expectations by education level and combining them (properly weighted)

$$\begin{aligned}
\mathbb{E}[Wage] &= \mathbb{E}(\mathbb{E}[Wage|Education]) \\
&= \sum_{Education_i} Pr(Education_i) \cdot E[Wage|Education_i]
\end{aligned}$$

10

---

| Person | Gender | IQ |
|--------|--------|-----|
| 1 | M | 120 |
| 2 | M | 115 |
| 3 | M | 110 |
| 4 | F | 130 |
| 5 | F | 125 |
| 6 | F | 120 |

- E[IQ] = 120
- E[IQ | Male] = 115; E[IQ | Female] = 125
- LIE: E ( E [ IQ | Sex ] ) = (0.5)×115 + (0.5)×125 = 120

## LIE: Proof for the discrete case

$$
\begin{aligned}
\mathbb{E}(\mathbb{E}[Y|X]) &= \sum_x \mathbb{E}[Y|X=x]p(x) \\
&= \sum_x \left( \sum_y y p(y|x) \right) p(x) \\
&= \sum_x \sum_y y p(x,y) \\
&= \sum_y y \sum_x p(x,y) \\
&= \sum_y y p(y) \\
&= \mathbb{E}(Y)
\end{aligned}
$$

*(handwritten annotations, red):* $\rightarrow \sum_x \sum_y y\, P(y|x)\, P(x)$

*(handwritten, blue):* $\underbrace{P(y|x)\,P(x)}_{P(x,y)}$ REGLA BAYES

*(handwritten, blue):* $\rightarrow \sum_y \sum_x y\, P(x,y)$

12

## LIE: Proof for the continuous case

$$
\begin{aligned}
\mathbb{E}[\mathbb{E}(Y|X)] &= \int \mathbb{E}(Y|X=u)g_x(u)du \\
&= \int \left[ \int t f_{y|x}(t|X=u)dt \right] g_x(u)du \\
&= \int \int t f_{y|x}(t|X=u)g_x(u)dudt \\
&= \int t \left[ \int f_{y|x}(t|X=u)g_x(u)du \right] dt \\
&= \int t \left[ f_{x,y}du \right] dt \\
&= \int t g_y(t)dt \\
&= \mathbb{E}(y)
\end{aligned}
$$

13

## OLS Review

Linear algebra review

Law of iterated expectations

OLS basics

Conditional expectation function

"Algebraic" properties of OLS

Properties of OLS estimators

Regression (matrix algebra) with a treatment dummy for the experimental case

Frisch–Waugh–Lovell (FWL) theorem

Regression and causality

14

## OLS Review

Linear algebra review

Law of iterated expectations

OLS basics

Conditional expectation function

"Algebraic" properties of OLS

Properties of OLS estimators

Regression (matrix algebra) with a treatment dummy for the experimental case

Frisch–Waugh–Lovell (FWL) theorem

Regression and causality

15

## OLS - As minimizing residuals

- Data with $n$ observations and two variables: $(x_1, ... x_n)$ and $(y_1, ..., y_n)$

- Find the line $(\widehat{\beta}_0 + \widehat{\beta}_1 x)$ that best fits the data

- $\widehat{y}_i = \widehat{\beta}_0 + \widehat{\beta}_1 x_i$ is the fitted value for $i$

- The residual is $\widehat{u}_i = y_i - \widehat{y}_i$

- Goal: minimize residuals or distance from the line (fitted values) to the data

16

## OLS - As minimizing residuals

- We don't care if the residual $\widehat{u}_i$ is positive or negative, we want it to be small

- So we square it: $\widehat{u}_i^2$

- Why not the absolute value? Good statistical reasons + harder to work with $|\cdot|$

- We want all the mistakes to be small, so we really want to minimize $\sum_{i=1}^n \widehat{u}_i^2$

17

## OLS - As minimizing residuals

$$\min_{\widehat{\beta}_0, \widehat{\beta}_1} \sum_{i=1}^n \widehat{u}_i^2 \;\; = \;\; \sum_{i=1}^n (y_i - \widehat{y}_i)^2 = \sum_{i=1}^n \left( y_i - \left( \widehat{\beta}_0 + \widehat{\beta}_1 x_i \right) \right)^2$$

$$\frac{\partial}{\partial \widehat{\beta}_0} = 0$$
$$\frac{\partial}{\partial \widehat{\beta}_1} = 0$$

- Using calculus (deriving with respect to $\widehat{\beta}_0, \widehat{\beta}_1$ and equating to zero):

$$\widehat{\beta}_1^* = \frac{\sum_{i=1}^n (x_i - \overline{x_i})(y_i - \overline{y_i})}{\sum_{i=1}^n (x_i - \overline{x_i})^2} = \frac{\frac{1}{n}\sum_{i=1}^n (x_i - \overline{x_i})(y_i - \overline{y_i})}{\frac{1}{n}\sum_{i=1}^n (x_i - \overline{x_i})^2} = \frac{\text{Sample covariance (x,y)}}{\text{Sample variance (x)}}$$

$$\widehat{\beta}_0^* = \overline{y_i} - \widehat{\beta}_1 \overline{x_i}$$

18

## Visual tour of OLS

- https://ryansafner.shinyapps.io/ols_estimation_by_min_sse/

- https://seeing-theory.brown.edu/regression-analysis/
  index.html#section1

- https://setosa.io/ev/ordinary-least-squares-regression/

- https://mgimond.github.io/Stats-in-R/regression.html

## OLS as an estimator

- There is a population with two random variables $x$ and $y$

- We take a **random sample** of size $n$: $(x_1, x_2, ..x_n)$ and $(y_1, y_2, ..., y_n)$

- We would like to see how $y$ varies with changes in $x$

  - What if y is affected by factors other than x?

  - What is the functional form connecting these two variables?

  - If interested in causal effect of $x$ on $y$, how to distinguish from mere correlation?

20

## OLS as an estimator of the DGP parameters

- Assume the data generating proces (DGP)s is:

$$y_i = \beta_0 + \beta_1 x_i + u_i$$

- That is, this model holds in the **population**
- Not only $x_i$ affects $y_i$, $u_i$ (called the error term) also does
- Do not confuse $u_i$ with $\widehat{u}_i$
- We assume there is a linear relationship between $y_i$ and $x_i$
- We never observe $\beta_0$ and $\beta_1$

21

## Inference

- Goal: Estimate unknown parameters

- To approximate parameters, we use an estimator, which is a function of the data

- Thus, estimator is a random variable (it is a function of a random variable)

- Infer something about the parameters from the distribution of the estimator

22

## Important notation

Based on this tweet: https://twitter.com/nickchk/status/1272993322395557888

- Greek letters (e.g., $\mu$) are the truth (i.e., parameters of the true DGP)
- Greek letters with hats (e.g., $\widehat{\mu}$) are estimates (i.e., what we *think* the truth is)
- Non-Greek letters (e.g., $X$) denote sample/data
- Non-Greek letters with lines on top (e.g., $\overline{X}$) denote calculations from the data
- We want to estimate the truth, with some calculation from the data ($\widehat{\mu} = \overline{X}$)
- Data $\longrightarrow$ Calculations $\longrightarrow$ Estimate $\underset{\text{Hopefully}}{\longrightarrow}$ Truth

- Example: $X \longrightarrow \overline{X} \longrightarrow \widehat{\mu} \underset{\text{Hopefully}}{\longrightarrow} \mu$

23

- Assume the data generating process is:

$$y_i = \beta_0 + \beta_1 x_i + u_i$$

- Also assume $\mathbb{E}u_i = 0$
  - Without loss of generality
  - We can just change the intercept to force $\mathbb{E}u_i = 0$
  - For example if $\mathbb{E}u_i = \alpha_0$
  - Redefine model to $y_i = \underbrace{\beta_0 + \alpha_0}_{\text{new intercept}} + \beta_1 x_i + \underbrace{u_i - \alpha_0}_{\text{new error term}}$
- Assume mean independence $\mathbb{E}(u_i|x_i) = \mathbb{E}(u_i)$ for all values x
  - This is a non-trivial assumption, but let's take it for granted for now
  - Implies that $\mathbb{E}(u_i|x) = \mathbb{E}(u_i) = 0$
  - Implies that $\mathbb{E}(u_i x_i) = \mathbb{E}(\mathbb{E}(u_i|x_i)) = 0$

24

---

- $\mathbb{E}(y_i|x_i) = \beta_0 + \beta_1 x_i$ 

$$E(Y_i|X_i) = E(\beta_0|X_i) + E(\beta_1 X_i|X_i) + E(U_i|X_i)^{\to 0}$$
$$E(\beta_0 + \beta_1 X_i + U_i|X_i)$$

- $\mathbb{E}(y_i|x_i)$: population regression function or conditional expectation function

- By our assumptions:

  - $\mathbb{E}(u_i|x_i) = \mathbb{E}(y_i - \beta_0 - \beta_1 x_i) = 0$

  - $\mathbb{E}(u_i x_i) = \mathbb{E}(x(y_i - \beta_0 - \beta_1 x_i)) = 0$

- These two conditions determine $\beta_0$ and $\beta_1$

25

---

First equation

$$\begin{aligned}
\mathbb{E}(y_i - \beta_0 - \beta_1 x_i) &= 0 \\
\mathbb{E}y_i - \beta_0 - \beta_1 \mathbb{E}x_i &= 0 \\
\mathbb{E}y_i - \beta_1 \mathbb{E}x_i &= \beta_0
\end{aligned}$$

26

---

Second equation

$$\begin{aligned}
\mathbb{E}x_i(y_i - \beta_0 - \beta_1 x_i) &= 0 \qquad \to \beta_0 \\
\mathbb{E}x_i(y_i - [\mathbb{E}y_i - \beta_1 \mathbb{E}x_i] - \beta_1 x_i) &= 0 \\
\mathbb{E}x_i(y_i - \mathbb{E}y_i - \beta_1(x_i - \mathbb{E}x_i)) &= 0 \\
\mathbb{E}x_i(y_i - \mathbb{E}y_i) &= \mathbb{E}x_i \beta_1(x_i - \mathbb{E}x_i) \\
\mathbb{E}(x_i - \mathbb{E}x_i)(y_i - \mathbb{E}y_i) &= \beta_1 \mathbb{E}(x_i - \mathbb{E}x_i)(x_i - \mathbb{E}x_i) \\
\frac{\mathbb{E}(x_i - \mathbb{E}x_i)(y_i - \mathbb{E}y_i)}{\mathbb{E}(x_i - \mathbb{E}x_i)(x_i - \mathbb{E}x_i)} &= \beta_1 \\
\frac{\text{Population covariance (x,y)}}{\text{Population variance (x)}} &= \beta_1
\end{aligned}$$

$$\beta_1 \; \mathbb{E}\left(x_i - \mathbb{E}(x_i)\right)\left(x_i - \mathbb{E}(x_i)\right)$$
$$\beta_1 \left( \mathbb{E}\left( x_i^2 - 2\mathbb{E}(x_i)x_i + \mathbb{E}^2(x_i) \right) \right)$$
$$\beta_1 \left( \mathbb{E}\left( x_i^2 - 2\mathbb{E}^2(x_i) + \mathbb{E}^2(x_i) \right) \right)$$
$$\beta_1 \left( \mathbb{E}\left( x_i^2 - \mathbb{E}^2(x_i) \right) \right)$$

27

---

- But we don't have x and y, nor do we know $\mathbb{E}y_i$ or $\mathbb{E}x_i$

- We only have a **random sample** of size $n$: $(x_1, ..., x_n)$ and $(y_1, ..., y_n)$

- The sample analogs:

  - $\frac{1}{n}\sum_{i=1}^{n}(y_i - \widehat{\beta}_0 - \widehat{\beta}_1 x_i) = 0$

  - $\frac{1}{n}\sum_{i=1}^{n} x_i(y_i - \widehat{\beta}_0 - \widehat{\beta}_1 x_i) = 0$

28

---

First equation

$$\begin{aligned}
\frac{1}{n}\sum_{i=1}^{n}(y_i - \widehat{\beta}_0 - \widehat{\beta}_1 x_i) &= 0 \\
\frac{1}{n}\sum_{i=1}^{n} y_i - \widehat{\beta}_0 - \widehat{\beta}_1 \frac{1}{n}\sum_{i=1}^{n} x_i &= 0 \\
\overline{y}_i - \widehat{\beta}_0 - \widehat{\beta}_1 \overline{x}_i &= 0 \\
\overline{y}_i - \widehat{\beta}_1 \overline{x}_i &= \widehat{\beta}_0
\end{aligned}$$

29

## OLS as an estimator of the DGP parameters

Second equation

$$\frac{1}{n}\sum_{i=1}^{n} x_i(y_i - \widehat{\beta}_0 - \widehat{\beta}_1 x_i) = 0$$

$$\frac{1}{n}\sum_{i=1}^{n} x_i(y_i - (\overline{y}_i - \widehat{\beta}_1 \overline{x}_i) - \widehat{\beta}_1 x_i) = 0$$

$$\frac{1}{n}\sum_{i=1}^{n} x_i(y_i - \overline{y}_i + \widehat{\beta}_1(\overline{x}_i - x_i)) = 0$$

$$\frac{1}{n}\sum_{i=1}^{n} x_i(y_i - \overline{y}_i) = \frac{1}{n}\sum_{i=1}^{n} x_i(\widehat{\beta}_1(x_i - \overline{x}_i))$$

$$\frac{1}{n}\sum_{i=1}^{n}(x_i - \overline{x}_i)(y_i - \overline{y}_i) = \widehat{\beta}_1\frac{1}{n}\sum_{i=1}^{n}(x_i - \overline{x}_i)(x_i - \overline{x}_i)$$

$$\frac{\frac{1}{n}\sum_{i=1}^{n}(x_i - \overline{x}_i)(y_i - \overline{y}_i)}{\frac{1}{n}\sum_{i=1}^{n}(x_i - \overline{x}_i)^2} = \widehat{\beta}_1$$

$$\frac{\text{Sample covariance (x,y)}}{\text{Sample variance (x)}} = \widehat{\beta}_1$$

30

## OLS as an estimator of the DGP parameters

- Formulas are the same as "minimizing residuals"

- Show the OLS coefficients as estimator of the population parameters ($\beta_0$ and $\beta_1$)

- Some remarks:

  - Can only estimate if the sample variance of $x_i$ is not zero

  - In other words, if $x_i$ is not constant across all values of $i$

  - Intuitively, the variation in x is what permits us to identify its impact in y

31

## Multiple regression – notation

- Consider the multiple linear regression model
$$y_i = x_i'\beta + u_i$$
where $\beta = (\beta_0, \beta_1, ..., \beta_K)'$ and $x_i = (1, ..., x_K)'$
  - $\beta$ is of size $(k \times 1)$
  - $x_i$ is of size $(k \times 1)$
  - $x_i'\beta$ is of size $(1 \times k)(k \times 1) = 1 \times 1$

- Equivalent
$$y = X\beta + u$$
where $\beta = (\beta_0, \beta_1, ..., \beta_K)'$
  - $\beta$ is of size $(k \times 1)$
  - $X$ is of size $(n \times k)$
  - $X\beta$ is of size $(n \times k)(k \times 1) = n \times 1$

32

## Multiple regression

- Consider the multiple linear regression model
$$y_i = x_i'\beta + u_i$$

- $x_i'\beta = \beta_0 + \sum_{k=1}^{K} x_{ik}\beta_k$ is the conditional expectation function $(\mathbb{E}(y_i|x_i))$
- The population regression $\beta$ coefficients solve
$$\beta = \mathbb{E}[x_i x_i']^{-1}\mathbb{E}[x_i y_i] = \arg\min_b \mathbb{E}\left[(y_i - x_i'b)^2\right]$$

- The sample equivalent is
$$\widehat{\beta} = \frac{\frac{1}{n}\sum_{i=1}^{n}(x_i y_i)}{\frac{1}{n}\sum_{i=1}^{n}(x_i x_i')} = \frac{\sum_{i=1}^{n}(x_i y_i)}{\sum_{i=1}^{n}(x_i x_i')} = (X'X)^{-1}X'y$$

33

*(handwritten annotations)*

$\mathbb{E}(Y_i|X_i) = \mathbb{E}(X_i'\beta|X_i) \times$
$+ \mathbb{E}(u|X_i)$

$X_i' = (X_{1i}, ..., X_{Ki})$
$X_i'\beta = (X_{0i}, ..., X_{Ki})\begin{pmatrix}\beta_0 \\ \vdots \\ \beta_K\end{pmatrix}$

$\sum_{j=0}^{k} X_{ji}\beta_j$

## Solving for $\widehat{\beta}$

- We let the computer do the calculations, which are tedious even for small $n$

- Good to know what's going on behind the scenes

- But I honestly do not care if you know how invert a matrix

34

## Simulations!

```
alpha=1 #intercept
beta=2 #slope
Nobs=10000 #how many observations?
X=runif(Nobs,-5,5)
#use the DGP to generate data
Y=alpha+beta*X+rnorm(Nobs)
OLS=lm(Y~X)
summary(OLS)
```

35

## Our estimate of the coefficient are pretty close to the truth

```
Call:
lm(formula = Y ~ X)

Residuals:
    Min      1Q  Median      3Q     Max
-3.5501 -0.6739 -0.0023  0.6706  3.7211

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 0.995720   0.010082   98.76   <2e-16 ***
X           1.998857   0.003522  567.59   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.008 on 9998 degrees of freedom
Multiple R-squared:  0.9699,    Adjusted R-squared:  0.9699
F-statistic: 3.222e+05 on 1 and 9998 DF,  p-value: < 2.2e-16
```

$\hat{\alpha}$

$\hat{\beta}$

36

## But how close on average?

```
#Now lets repeat the process and see how close our estimates are
Reps=1000
alpha_estimate=NULL
beta_estimate=NULL
for(r in 1:Reps){
  X=runif(Nobs,-5,5)
  #use the DGP to generate data
  Y=alpha+beta*X+rnorm(Nobs)
  OLS=lm(Y~X)
  Estimates=summary(OLS)$coef[,"Estimate"]
  alpha_estimate=c(alpha_estimate,Estimates[1])
  beta_estimate=c(beta_estimate,Estimates[2])
}
hist(beta_estimate,freq=F,breaks=30,
     main="",las=1,xlab="Estimate of beta")
abline(v=beta,col='red',lwd=2,lty=2)
```

Vector de $\hat{\alpha}$
Vector de $\hat{\beta}$

37

$$\hat{\beta} \sim N\left(\beta, \sigma_{\hat{\beta}}^2\right)$$

ERROR
ESTANDAR
DE
$\hat{\beta}$

## Distribution of estimate of $\hat{\beta}$



38

## OLS Review

Linear algebra review

Law of iterated expectations

OLS basics

Conditional expectation function

"Algebraic" properties of OLS

Properties of OLS estimators

Regression (matrix algebra) with a treatment dummy for the experimental case

Frisch–Waugh–Lovell (FWL) theorem

Regression and causality

39

## OLS Review

Linear algebra review

Law of iterated expectations

OLS basics

Conditional expectation function

"Algebraic" properties of OLS

Properties of OLS estimators

Regression (matrix algebra) with a treatment dummy for the experimental case

Frisch–Waugh–Lovell (FWL) theorem

Regression and causality

40

## Conditional expectation function (CEF)

- Assume we are interested in the returns to schooling

- Summarize the effect of schooling on wages with the CEF ($\mathbb{E}(y_i|x_i)$)

- The CEF is the expectation (i.e, population average) of $y_i$ with $x_i$ held constant

- $\mathbb{E}(y_i|x_i)$ provides a reasonable representation of how $y$ changes with $x$

- Because $x_i$ is random, $\mathbb{E}[y_i \mid x_i]$ is random

- Sometimes work with a particular value of the CEF (e.g., $\mathbb{E}[y_i \mid x_i = 12]$)

41

## Property 1: CEF Decomposition Property

- $y_i = \mathbb{E}(y_i|x_i) + u_i$ where

  1. $u_i$ is mean independent of $x_i$; that is $\mathbb{E}(u_i|x_i) = 0$

  2. $u_i$ is uncorrelated with any function of $x_i$

- In words: any random variable, $y_i$, can be decomposed into two parts: the part that can be explained by $x_i$ and the part left over that cannot be explained by $x_i$

- Proof is in Angrist and Pischke (ch. 3)

42

## Property 2: CEF Prediction Property

- Let $m(x_i)$ be any function of $x_i$

- $\mathbb{E}(y_i|x_i) = \arg\min_{m(x_i)} \mathbb{E}[(y_i - m(x_i))^2]$

- In words: The CEF is the minimum mean squared error predictor of $y_i$ given $x_i$

- Proof is in Angrist and Pischke (ch. 3)

43

## Property 3: Best linear approximation

- The population regression is the best linear approximation to the true nonlinear CEF in a mean squared error sense:

$$\beta = \mathbb{E}[x_i x_i']^{-1}\mathbb{E}[x_i y_i] = \arg\min_b \mathbb{E}[(\mathbb{E}[y_i \mid x_i] - x_i' b)^2]$$

- In words: even if the true CEF is nonlinear (for example, $E[y_i \mid x_i] = \log(x_i)$), regression is still a good approximation to the truth

44

## Why linear regression may be of interest (summary)

- If the CEF is linear. Then the population regression is it

  - Then it makes the most sense to use linear regression to estimate it

- Linear regression may be interesting even if the underlying CEF is not linear

  - $\mathbb{E}(y_i|x_i)$, is the minimum mean squared error predictor of $y_i$ given $x_i$ in the class of all functions of $x_i$

  - The population regression function is the best we can do in the class of all linear functions to approximate $\mathbb{E}(y_i|x_i)$

45

## Big picture

1. Regression provides the best linear predictor for the dependent variable in the same way that the CEF is the best unrestricted predictor of the dependent variable

2. If we prefer to think of approximating $\mathbb{E}(y_i|x_i)$ as opposed to predicting $y_i$, even if the CEF is nonlinear, regression provides the best linear approximation to it

46

## OLS Review

Linear algebra review

Law of iterated expectations

OLS basics

Conditional expectation function

"Algebraic" properties of OLS

Properties of OLS estimators

Regression (matrix algebra) with a treatment dummy for the experimental case

Frisch–Waugh–Lovell (FWL) theorem

Regression and causality

47

## Residuals add up to zero

- Remembering how the **first moment** condition allows us to obtain $\widehat{\beta}_0$ and $\widehat{\beta}_1$:

$$\frac{1}{n}\sum_{i=1}^{n}(y_i - \widehat{\beta}_0 - \widehat{\beta}_1 x_i) = 0$$

*(handwritten:)* $\mathbb{E}(U_i) = 0$   $= \frac{1}{n}\sum \widehat{U}_i = 0$

*(handwritten:)* $\frac{1}{n}\sum U_i = Y_i - \beta_0 - \beta_1 X_i = 0$

- This means the OLS residuals *always* add up to zero, by *construction*,

$$\frac{1}{n}\sum_{i=1}^{n}\widehat{u}_i = 0$$

$$\sum_{i=1}^{n}\widehat{u}_i = 0$$

## The mean of the fitted values is the mean of the data

Because $y_i = \widehat{y}_i + \widehat{u}_i$ by definition,

$$\sum_{i=1}^{n} y_i = \sum_{i=1}^{n} \widehat{y}_i + \sum_{i=1}^{n}\widehat{u}_i$$

*(handwritten: arrow to last term)* $\to 0$

$$\frac{1}{n}\sum_{i=1}^{n} y_i = \frac{1}{n}\sum_{i=1}^{n}\widehat{y}_i$$

$$\overline{y} = \overline{\widehat{y}}$$

*(handwritten at right:)*

$$Y_i = \mathbb{E}(Y_i|X_i) + U_i$$

$$U_i \neq \widehat{U}_i$$

$$\boxed{\widehat{U}_i = Y_i - \widehat{Y}_i}$$

## Sample correlation between $x_i$ and residuals is zero

Similarly the way we obtained our estimates,

$$\frac{1}{n}\sum_{i=1}^{n} x_i(y_i - \widehat{\beta}_0 - \widehat{\beta}_1 x_i) = 0 \tag{1}$$

*(handwritten:)* $\boxed{\mathbb{E}(X_i U_i) = 0}$   $\frac{1}{n}\sum Y_i U_i = 0$

The sample covariance (and therefore the sample correlation) between the explanatory variables and the residuals is always zero:

$$\frac{1}{n}\sum_{i=1}^{n} x_i\widehat{u}_i = 0 \tag{2}$$

*(handwritten:)* $\mathbb{E}(X_i U_i) \neq 0$

51

## Bringing things together

Because the $\widehat{y}_i$ are linear functions of the $x_i$, the fitted values and residuals are uncorrelated, too:

$$\frac{1}{n}\sum_{i=1}^{n}\widehat{y}_i\widehat{u}_i = 0 \tag{3}$$

52

## The point $(\overline{x}, \overline{y})$ is always on the OLS regression line

If we plug in the average for $x$, we predict the sample average for $y$:

$$\overline{y} = \widehat{\beta}_0 + \widehat{\beta}_1\overline{x} \tag{4}$$

(see formula for $\widehat{\beta}_0$)

*(handwritten: "MCO" with a line through points)*

53

```
Nobs=1000 #how many observations?
X=runif(Nobs,-5,5)
#use the DGP to generate data
Y=10+2*X^2+rnorm(Nobs)
OLS=lm(Y~X)
summary(OLS)
plot(X,Y,bty="L")
abline(OLS,col=2,lwd=2,lty=2)
points(mean(X),mean(Y),pch=19,col=4,cex=1.5)
#Not a great fit...yet
#residual add to zero
sum(OLS$residuals)
#mean of fitted values is the mean of true values
mean(OLS$fitted.values)-mean(Y)
#sample covariance between X and residuals is zero
sum(OLS$residuals*X)
#sample covariance between fitted values and residuals is zero
sum(OLS$residuals*OLS$fitted.values)
```

54



## Algebraic properties

```
> #Not a great fit...yet
> #residual add to zero
> sum(OLS$residuals)
[1] 1.097566e-12
> #mean of fitted values is the mean of true values
> mean(OLS$fitted.values)-mean(Y)
[1] 0
> #sample covariance between X and residuals is zero
> sum(OLS$residuals*X)
[1] 1.684319e-12
> #sample covariance between fitted values and residuals is zero
> sum(OLS$residuals*OLS$fitted.values)
[1] 3.242961e-11
```

56

## Big picture

Don't let anyone tell you the model is good because any of the following happens

1. Residuals add to zero

2. Fitted values mean is equal to data mean

3. Residuals are uncorrelated with x

4. If we plug in the average for $x$, we predict the sample average for $y$

These results are mechanical: Unrelated to how appropriate the model is or "causality"

57

## OLS Review

58

## OLS Review

59

## Expected Value of OLS

- Mathematical statistics: How do our estimators behave across different samples of data? On average, would we get the right answer if we could repeatedly sample?

- Find the expected value of the OLS estimators – the average outcome across all possible random samples – and determine if we are right on average

- Leads to the notion of **unbiasedness**, a "desirable" characteristic for estimators.

$$\mathbb{E}(\widehat{\beta}) = \beta \tag{5}$$

## Don't forget why we're here

- The **population** parameter that describes the relationship between $y$ and $x$ is $\beta$

- Goal: estimate $\beta$ with a sample of data

- $\widehat{\beta}$ is an **estimator** obtained with a *specific* sample from the population

## Uncertainty and sampling variance

- Different samples will generate different estimates ($\widehat{\beta}$) for the "true" $\beta$

- Thus, $\widehat{\beta}$ a random variable (depends on random samples)

- Unbiasedness is the idea that if we could take as many random samples on $y$ as we want from the population, and compute an estimate each time, the average of these estimates would be equal to $\beta$

- But, this also implies that $\widehat{\beta}$ has spread and therefore variance

## Assumption 1 (Linear in Parameters)

- The population model can be written as

$$y = X\beta + u \tag{6}$$

where $\beta$ are the (unknown) population parameters

- We view $X$ and $u$ as outcomes of random variables; thus, $y$ is random

- Our goal is to estimate $\beta$

- $u$ is the unobserved error. It is not the residual that we compute from the data!

## Assumption 2 (Random Sampling)

- We have a random sample of size $n$, $\{(x_i, y_i) : i = 1, ..., n\}$

- We know how to use this data to estimate $\beta$ by OLS

## Assumption 3 (Zero Conditional Mean)

$$Y = X\beta + U$$
$$\text{COVID} \quad \overset{\downarrow}{\text{EDUCACION}}$$

- In the population, the error term has zero mean given any value of $X$:

$$\underline{\mathbb{E}(u|X)} = \mathbb{E}(u) = 0. \tag{7}$$

- This is the key assumption for showing that OLS is unbiased, with the zero value not being important once we assume $\mathbb{E}(u|X)$ does not change with $X$

## Assumption 1-3

- We can compute the OLS estimates whether or not these assumption hold

- But we might not get a "good" estimate

## Assumption 4 (Sample Variation in the Explanatory Variable)

- The sample outcomes on $x_i$ are not all the same value

- Same as saying the sample variance of $\{x_i : i = 1, ..., n\}$ is not zero

- If the $x_i$ are all the same value, we cannot learn how $x$ affects $y$

## Showing OLS is unbiased

- How do we show $\widehat{\beta}$ is unbiased for $\beta$?

- We know $\widehat{\beta} = (X'X)^{-1}X'y$

- And that $y = X\beta + u$ (by assumption 1)

- Therefore: $\widehat{\beta} = (X'X)^{-1}X'(X\beta + u) = \beta + (X'X)^{-1}X'u$

$$(X'X)^{-1}X'X\beta + (X'X)^{-1}X'u$$

- $\mathbb{E}(\widehat{\beta} \mid X) = \beta + (X'X)^{-1}X'\underbrace{\mathbb{E}(u \mid X)}_{=0 \text{ by assumption 3}}$

- $\mathbb{E}(\widehat{\beta} \mid X) = \beta$

$$\mathbb{E}\left(\mathbb{E}\left(\widehat{\beta} \mid X\right)\right) = \beta$$
$$\mathbb{E}\left(\widehat{\beta}\right) = \beta$$

- Each sample leads to a different estimate, $\widehat{\beta}$

- Some will be very close to the true values $\beta$

- Some **could** be very far from those values

- If we repeat the experiment and average the estimates $\to$ very close to $\beta$

- But in a single sample, we can never know whether we are close to $\beta$

- Next: measure of dispersion (spread) in the distribution of the estimators

## Repeat our simulations with different N

```
alpha=1 #intercept
beta=2 #slope
Reps=1000
for(Nobs in c(100,1000,10000)){
  alpha_estimate=NULL
  beta_estimate=NULL
  for(r in 1:Reps){
    X=runif(Nobs,-5,5)
    Y=alpha+beta*X+rnorm(Nobs)
    OLS=lm(Y~X)
    Estimates=summary(OLS)$coef[,"Estimate"]
    alpha_estimate=c(alpha_estimate,Estimates[1])
    beta_estimate=c(beta_estimate,Estimates[2])
  }
  hist(beta_estimate,freq=F,breaks=30,main="",las=1)
  abline(v=beta,col='red',lwd=3,lty=1)
}
```

## Repeat our simulations with different N — Look at the x-axis scale

- **Errors** are the vertical distances between observations and the **unknown** Conditional Expectation Function. Therefore, they are unknown.

- **Residuals** are the vertical distances between observations and the **estimated** regression function. Therefore, they are known.

## Variance of OLS estimators

The correct variance estimation procedure is given by the structure of the data

- It is very unlikely that all observations in a dataset are unrelated, but drawn from identical distributions (**homoskedasticity**)

- For instance, the variance of income is often greater in families belonging to top deciles than among poorer families (**heteroskedasticity**)

- Some phenomena do not affect observations individually, but they do affect groups of observations uniformly within each group (**clustered data**)

## Assumption 5 (Homoskedasticity, or Constant Variance)

The error has the same variance given any value of the explanatory variable $x$:

$$Var(u|X) = \sigma^2 > 0 \qquad (8)$$

where $\sigma^2$ is (virtually always) unknown.

Because $\mathbb{E}(u|x) = 0$ we can also write

$$\mathbb{E}(u^2|x) = \sigma^2 = \mathbb{E}(u^2) \qquad (9)$$

## Assumption 5 (Homoskedasticity, or Constant Variance)

Under the our assumptions

$$
\begin{aligned}
y &= X\beta + u \\
\mathbb{E}(y|x) &= X\beta \\
Var(y|x) &= \sigma^2
\end{aligned}
$$

$$V(y) = V(X\beta) + V(u)$$
$$V(y|x) = V(u|x) = \sigma^2$$

The average or expected value of $y$ is allowed to change with $x$, but the variance does not change with $x$

## Assumption 5 (Homoskedasticity, or Constant Variance)



Figure 2.8
The simple regression model under homoskedasticity.

**Variance of OLS estimators** In matrix form the property that $V(aW) = a^2 V(W)$ where $a$ is constant and $W$ is a random variable is written as:

$$V(AW) = AV(W)A'$$

where $A$ is a constant matrix and $W$ is a random variable

## Variance of OLS estimators

- We know $\widehat{\beta} = (X'X)^{-1}X'y$

- And that $y = X\beta + u$ (by assumption 1)

- Therefore: $\widehat{\beta} = (X'X)^{-1}X'(X\beta + u) = \beta + (X'X)^{-1}X'u$

- $V(\widehat{\beta} \mid X) = \underbrace{V(\beta \mid X)}_{=0 \text{ since it's constant}} + (X'X)^{-1}X' \underbrace{V(u \mid X)}_{=\sigma^2 \text{ by assumption 3}} X(X'X)^{-1}$

- $V(\widehat{\beta} \mid X) = (X'X)^{-1}X'\sigma^2 X(X'X)^{-1} = \sigma^2(X'X)^{-1}$

*(handwritten annotations:)*

$$\left((X'X)^{-1}\right)^t = \left((X'X)^t\right)^{-1} = \left(X^t X\right)^{-1}$$

$$\sigma^2 (X'X)^{-1} X' X (X'X)^{-1}$$

## Estimating the Error Variance

- In the formula
$$V(\widehat{\beta} \mid X) = (X'X)^{-1}X'\sigma^2 X(X'X)^{-1} = \sigma^2(X'X)^{-1}$$
we can compute $(X'X)^{-1}$ but we need to estimate $\sigma^2$

- Recall that
$$\sigma^2 = \mathbb{E}(u^2)$$

## Estimating the Error Variance

- If we could observe the errors $(u_i)$ an unbiased estimator of $\sigma^2$ would be
$$\frac{1}{n}\sum_{i=1}^{n} u_i^2 \tag{10}$$

- But this not a *feasible* estimator because the $u_i$ are unobserved

- How about replacing each $u_i$ with its "estimate", the OLS residual $\widehat{u}_i$?

$$u_i = y_i - x_i'\beta$$
$$\widehat{u}_i = y_i - x_i'\widehat{\beta}$$

## Estimating the Error Variance

$\widehat{u}_i$ can be computed from the data, but $\widehat{u}_i \neq u_i$ for any $i$:

$$\widehat{u}_i = y_i - x_i'\widehat{\beta} = x_i'\beta + u_i - x_i'\widehat{\beta}$$
$$= u_i - (\widehat{\beta} - \beta)x_i$$

$\mathbb{E}(\widehat{\beta}) = \beta$ but the estimators differ from the population values in a given sample

## Estimating the Error Variance

- Now, what about this as an estimator of $\sigma^2$?
$$\frac{1}{n}\sum_{i=1}^{n} \widehat{u}_i^2 \tag{11}$$

- It is a *feasible* estimator and easily computed from the data after OLS

- As it turns out, this estimator is slightly biased

## Estimating the Error Variance

The estimator does not account for the restrictions on the residuals, used to obtain $\widehat{\beta}$

$$\sum_{i=1}^{n} \widehat{u}_i = 0$$
$$\sum_{i=1}^{n} x_{1i}\widehat{u}_i = 0$$
$$\vdots$$
$$\sum_{i=1}^{n} x_{ki}\widehat{u}_i = 0$$

*(handwritten: )* K RESTRICCIONES

There is no such restriction on the unobserved errors

The unbiased estimator of $\sigma^2$ uses a **degrees-of-freedom** adjustment The residuals have only $n - k$ degrees-of-freedom (minus the k restrictions), not $n$

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^{n} \hat{u}_i^2}{(n-k)}$$

**THEOREM: Unbiased Estimator of** $\sigma^2$

Under Assumptions 1-5,

$$\mathbb{E}(\hat{\sigma}^2) = \sigma^2$$

---

- Given $\hat{\sigma}$, we can now estimate $V(\hat{\beta})$

- $V(\hat{\beta})$ is a variance-covariance matrix (size $k \times k$)

- The diagonal elements of $V(\hat{\beta})$ give us the variance of the estimators $\hat{\beta}$

- $\hat{\sigma}_{\hat{\beta}}$: The square root of the diagonal elements of the estimator of $V(\hat{\beta})$ is usually called the **standard errors** (i.e., estimate of the standard deviation of the estimator)

*[handwritten: ERROR STD $\widehat{\sigma}_\beta$]*

*[handwritten: DESV. ESTIMADOR $\sigma_\beta$]*

---

**Bringing the central limit theorem to play**

- By some version of the central limit theorem:

$$\frac{\hat{\beta} - \beta}{\sigma_{\hat{\beta}}} \to_d N(0,1)$$
$$\hat{\beta} \to_d \sigma_{\hat{\beta}} N(0,1) + \beta$$
$$\hat{\beta} \to_d N(\beta, \sigma_{\hat{\beta}}^2)$$

- $\sigma_\beta = \sigma^2 (X'X)^{-1}$
- Since we do not know $\sigma^2$, we estimate it
- $\hat{\sigma}_{\hat{\beta}} = \hat{\sigma}^2 (X'X)^{-1}$
- By some version of the central limit theorem + some statistical properties

$$\frac{\hat{\beta} - \beta}{\hat{\sigma}_{\hat{\beta}}} \to_d t_{n-k}$$
$$\hat{\beta} \to_d \hat{\sigma}_{\hat{\beta}} t_{n-k} + \beta$$

*[handwritten: $V(aX) = a^2 V(X)$]*

---

**To keep things simple**

- $t_{n-k} \to N(0,1)$ as $(n-k) \to \infty$

- So as long as your sample is large, we can keep thinking of normal distributions

$$\hat{\beta} \approx N(\beta, \hat{\sigma}_{\hat{\beta}})$$

---

**31.7% of estimates will be more than $\hat{\sigma}_{\hat{\beta}}$ away from $\beta$**

---

**4.55% of estimates will be more than $2\hat{\sigma}_{\hat{\beta}}$ away from $\beta$**

## We can know learn something about the true $\beta$

- We know $\widehat{\beta} \sim N(\beta, \widehat{\sigma_{\widehat{\beta}}})$
- We want to find some interval on which $\beta$ is likely to live:

$$P\left(a \leq \beta \leq b\right) = 1 - \alpha$$

- $P\left(-a \geq -\beta \geq b\right) = 1 - \alpha$

- $P\left(\frac{\widehat{\beta} - a}{\sigma_{\widehat{\beta}}} \geq \frac{\widehat{\beta} - \beta}{\sigma_{\widehat{\beta}}} \geq \frac{\widehat{\beta} - b}{\sigma_{\widehat{\beta}}}\right) = 1 - \alpha$

  standard normal
  ($t_{n-k}$ to be exact)

- Assuming we want symmetry (so $\frac{\alpha}{2}$ on each side), then:
  - $\Phi\left(\frac{\widehat{\beta} - a}{\sigma_{\widehat{\beta}}}\right) = 1 - \frac{\alpha}{2}$
  - $\Phi\left(\frac{\widehat{\beta} - b}{\sigma_{\widehat{\beta}}}\right) = \frac{\alpha}{2}$



## Confidence interval

- Thus:
  - $\Phi^{-1}\left(\frac{\alpha}{2}\right) = \frac{\widehat{\beta} - b}{\sigma_{\widehat{\beta}}}$
  - $\Phi^{-1}\left(1 - \frac{\alpha}{2}\right) = \frac{\widehat{\beta} - a}{\sigma_{\widehat{\beta}}}$
  - $b = \widehat{\beta} - \Phi^{-1}\left(\frac{\alpha}{2}\right)\widehat{\sigma_{\widehat{\beta}}}$
  - $a = \widehat{\beta} - \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)\widehat{\sigma_{\widehat{\beta}}}$

- $\beta$ is between $\widehat{\beta} - \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)\widehat{\sigma_{\widehat{\beta}}}$ and $\widehat{\beta} - \Phi^{-1}\left(\frac{\alpha}{2}\right)\widehat{\sigma_{\widehat{\beta}}}$ with probability $1 - \alpha$

91

## Confidence interval

- Say $\alpha = 5\%$, then $\Phi^{-1}\left(\frac{\alpha}{2}\right) = -1.96$ and $\Phi^{-1}\left(1 - \frac{\alpha}{2}\right) = 1.96$

- Then we know $\beta$ is between with probability 95%:
  - $\widehat{\beta} - 1.96\widehat{\sigma_{\widehat{\beta}}}$
  - $\widehat{\beta} + 1.96\widehat{\sigma_{\widehat{\beta}}}$

- Generally speaking, confidence intervals are wider, the smaller $\alpha$ is

*QNORM(0.025)*

*QNORM(1 - 0.025)*

92

## Simulations!

```
beta0=1 #intercept
beta1=2 #slope
#Now lets repeat the process and see how close our estimates are
Reps=1000
Nobs=100 #number of obs
beta0.estimate=NULL #vector to store estimates of beta0
beta1.estimate=NULL #vector to store estimates of beta1
EstimateInCI=NULL #vector to store whether estimate is in CI
LowerCI=NULL #vector to store lower bound CI
UpperCI=NULL #vector to store upper bound CI
Confidence.level=0.05 #alpha
for(r in 1:Reps){
  X=runif(Nobs,-5,5)
  Y=beta0+beta1*X+rnorm(Nobs)  #use the DGP to generate data
  OLS=lm(Y~X)
  Estimates=summary(OLS)$coef[,"Estimate"] #estimate
  SE=summary(OLS)$coef[,"Std. Error"] #estimate sigma_beta
  beta0.estimate=c(beta0.estimate,Estimates[1])
  beta1.estimate=c(beta1.estimate,Estimates[2])
  CI.Beta.lowerbound=Estimates[2]+qnorm(Confidence.level/2)*SE[2]  #Confidence intervals
  LowerCI=c(LowerCI,CI.Beta.lowerbound)
  CI.Beta.upperbound=Estimates[2]+qnorm(1-Confidence.level/2)*SE[2]
  UpperCI=c(UpperCI,CI.Beta.upperbound)
  DummyInCI=c(CI.Beta.lowerbound<beta1 & CI.Beta.upperbound>beta1)  #Is the true value in CI?
  EstimateInCI=c(EstimateInCI,DummyInCI)
}
mean(EstimateInCI)
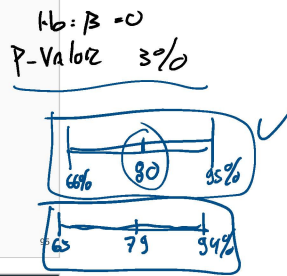```

$\Phi^{-1}\left(\frac{\alpha}{2}\right)$

93

## First ten simulations (red line is true $\beta$)



94

## Test hypothesis

- Is $\beta \neq \beta_0$?
- Usually posed as testing $H_0 : \beta = \beta_0$ vs $H_a : \beta \neq \beta_0$
- Different way to look at this: is $\beta_0$ in the confidence interval of $\beta$?
- Confidence interval depends on our choice of $\alpha$
- Pick largest $\alpha$ for which $\beta_0$ is not in the confidence interval
- This is called the $p - value$
- Largest probability of obtaining results at least as extreme as those actually observed, under the assumption that the null hypothesis is correct

*H0: $\beta$ = 0*
*P-Value  3%*



## Simulations!

$$V(\widehat{\beta}) = V\begin{pmatrix} \widehat{\beta}_{2020} \\ \widehat{\beta}_{2021} \end{pmatrix} = \widehat{\sigma}_{\widehat{\beta}}^2 (X'X)^{-1}$$

$$\widehat{D} = \widehat{\beta}_{2020} - \widehat{\beta}_{2021}$$

$$V(\widehat{D}) = \widehat{\sigma}_{\widehat{\beta}_{2020}}^2 + \widehat{\sigma}_{\widehat{\beta}_{2021}}^2 - 2\,cov(\widehat{\beta}_{2020}, \widehat{\beta}_{2021})$$

$$V(\hat{\delta}) = \hat{\sigma}^2_{\hat{\beta}_{2020}} + \hat{\sigma}^2_{\hat{\beta}_{2021}} - 2\left(\text{cov}\left(\hat{\beta}_{2020}, \hat{\beta}_{2021}\right)\right)$$

## Simulations!

```
#### p-values
beta0=1 #intercept
beta1=2 #slope
#Now lets repeat the process and see how close our estimates are
Reps=1000
Nobs=1000 #number of obs
pvalue_beta0=NULL #vector to store whether 0 is in CI
pvalue_beta2=NULL #vector to store whether 2 is in CI
Confidence_level=0.05 #alpha
for(r in 1:Reps){
  X=runif(Nobs,-5,5)
  #use the DGP to generate data
  Y=beta0+beta1*X+rnorm(Nobs)
  OLS=lm(Y~X)
  Estimates=summary(OLS)$coef[,"Estimate"] #estimate
  SE=summary(OLS)$coef[,"Std. Error"] #estimate sigma_beta
  beta0_estimate=c(beta0_estimate,Estimates[1])
  beta1_estimate=c(beta1_estimate,Estimates[2])
  #Confidence intervals
  CI_Beta_lowerbound=Estimates[2]+qnorm(Confidence_level/2)*SE[2]
  CI_Beta_upperbound=Estimates[2]+qnorm(1-Confidence_level/2)*SE[2]
  #Is 0 in CI?
  pvalue_beta0=c(pvalue_beta0,(CI_Beta_lowerbound<0 & CI_Beta_upperbound>0))
  #Is 2 in CI?
  pvalue_beta2=c(pvalue_beta2,(CI_Beta_lowerbound<2 & CI_Beta_upperbound>2))
}
mean(pvalue_beta0)
mean(pvalue_beta2)
```

## OLS Review

Linear algebra review

Law of iterated expectations

OLS basics

Conditional expectation function

"Algebraic" properties of OLS

Properties of OLS estimators

Regression (matrix algebra) with a treatment dummy for the experimental case

Frisch–Waugh–Lovell (FWL) theorem

Regression and causality

## OLS Review

Linear algebra review

Law of iterated expectations

OLS basics

Conditional expectation function

"Algebraic" properties of OLS

Properties of OLS estimators

Regression (matrix algebra) with a treatment dummy for the experimental case

Frisch–Waugh–Lovell (FWL) theorem

Regression and causality

## OLS

- $\hat{\beta} = (X'X)^{-1}X'y$

What's going on behind the scenes?

## Simple case

- Relationship between outcome $Y_i$ and treatment indicator $T_i$
- Regress the outcome on the treatment indicator, and a constant
- $X_i = \begin{pmatrix} 1 & T_i \end{pmatrix}$
- Assume first $N_T$ units are treated ($N_C = N - N_T$ units are untreated)
- $X = \begin{pmatrix} 1 & T_1 \\ 1 & T_2 \\ \vdots & \\ 1 & T_{N_T} \\ 1 & T_{N_T+1} \\ \vdots & \\ 1 & T_N \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 1 & 1 \\ \vdots & \vdots \\ 1 & 1 \\ 1 & 0 \\ \vdots & \vdots \\ 1 & 0 \end{pmatrix} \begin{array}{l} \left. \rule{0pt}{2em} \right\} N_T \\ \left. \rule{0pt}{2em} \right\} N_C \end{array}$

## Simple case

$$\bullet \; (X'X) = \begin{pmatrix} 1 & 1\cdots1 & 1 & \cdots1 \\ 1 & 1\cdots1 & 0 & \cdots0 \end{pmatrix} \underbrace{\phantom{xx}}_{N_T}\underbrace{\phantom{xx}}_{N_C} \begin{pmatrix} 1 & 1 \\ 1 & 1 \\ \vdots & \vdots \\ 1 & 1 \\ 1 & 0 \\ \vdots & \vdots \\ 1 & 0 \end{pmatrix} \begin{matrix} N_T \end{matrix} = \begin{pmatrix} N & N_T \\ N_T & N_T \end{pmatrix}$$

$$\bullet \; (X'X)^{-1} = \frac{1}{N_T(N-N_T)}\begin{pmatrix} N_T & -N_T \\ -N_T & N \end{pmatrix} = \frac{1}{N_C}\begin{pmatrix} 1 & -1 \\ -1 & \frac{N}{N_T} \end{pmatrix}$$

*(handwritten annotations)*

$$\frac{N_T}{N_T(N-N_T)} = \frac{1}{N_C} \qquad -\frac{1}{N_C}$$

$$\frac{-N_T}{N_T(N-N_T)} = -\frac{1}{N_C} \qquad \frac{N}{N_T(N-N_T)} = \frac{N}{N_T N_C}$$

101

## Simple case

$$\bullet \; X'y = \begin{pmatrix} 1 & 1\cdots1 & 1 & \cdots1 \\ 1 & 1\cdots1 & 0 & \cdots0 \end{pmatrix} \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_{N_T} \\ Y_{N_T+1} \\ \vdots \\ Y_N \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^N Y_i \\ \sum_{i=1}^{N_T} Y_i \end{pmatrix}$$

102

## Simple case

$$(X'X)^{-1}X'y = \frac{1}{N_C}\begin{pmatrix} 1 & -1 \\ -1 & \frac{N}{N_T} \end{pmatrix}\begin{pmatrix} \sum_{i=1}^N Y_i \\ \sum_{i=1}^{N_T} Y_i \end{pmatrix}$$

103

## Simple case

$$(X'X)^{-1}X'y = \frac{1}{N_C}\begin{pmatrix} 1 & -1 \\ -1 & \frac{N}{N_T} \end{pmatrix}\begin{pmatrix} \sum_{i=1}^N Y_i \\ \sum_{i=1}^{N_T} Y_i \end{pmatrix}$$

$$= \frac{1}{N_C}\begin{pmatrix} \sum_{i=1}^N Y_i - \sum_{i=1}^{N_T} Y_i \\ \frac{N}{N_T}\sum_{i=1}^{N_T} Y_i - \sum_{i=1}^N Y_i \end{pmatrix} \rightarrow$$

103

## Simple case

$$(X'X)^{-1}X'y = \frac{1}{N_C}\begin{pmatrix} 1 & -1 \\ -1 & \frac{N}{N_T} \end{pmatrix}\begin{pmatrix} \sum_{i=1}^N Y_i \\ \sum_{i=1}^{N_T} Y_i \end{pmatrix}$$

$$= \frac{1}{N_C}\begin{pmatrix} \sum_{i=1}^N Y_i - \sum_{i=1}^{N_T} Y_i \\ \frac{N}{N_T}\sum_{i=1}^{N_T} Y_i - \sum_{i=1}^N Y_i \end{pmatrix}$$

$$= \frac{1}{N_C}\begin{pmatrix} \sum_C Y_i \\ \frac{N}{N_T}\sum_T Y_i - \sum_T Y_i - \sum_C Y_i \end{pmatrix}$$

$$= \frac{1}{N_C}\begin{pmatrix} \sum_C Y_i \\ \frac{N-N_T}{N_T}(\sum_T Y_i) - \sum_C Y_i \end{pmatrix}$$

103

## Simple case

$$(X'X)^{-1}X'y = \frac{1}{N_C}\begin{pmatrix} 1 & -1 \\ -1 & \frac{N}{N_T} \end{pmatrix}\begin{pmatrix} \sum_{i=1}^N Y_i \\ \sum_{i=1}^{N_T} Y_i \end{pmatrix}$$

$$= \frac{1}{N_C}\begin{pmatrix} \sum_{i=1}^N Y_i - \sum_{i=1}^{N_T} Y_i \\ \frac{N}{N_T}\sum_{i=1}^{N_T} Y_i - \sum_{i=1}^N Y_i \end{pmatrix}$$

$$= \frac{1}{N_C}\begin{pmatrix} \sum_C Y_i \\ \frac{N}{N_T}\sum_T Y_i - \sum_T Y_i - \sum_C Y_i \end{pmatrix}$$

$$= \frac{1}{N_C}\begin{pmatrix} \sum_C Y_i \\ \frac{N-N_T}{N}(\sum_T Y_i) - \sum_C Y_i \end{pmatrix}$$

$$\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = \begin{pmatrix} \bar{Y}_C \\ \bar{Y}_T - \bar{Y}_C \end{pmatrix}$$

*(handwritten annotations)*

$$\frac{1}{N_C}\sum_C Y_i = \bar{Y}_C$$

$$\frac{N_C}{N_T N_C}\sum_T Y - \frac{1}{N_C}\sum_C Y_i$$

$$\bar{Y}_T - \bar{Y}_C$$

103

## Simple case

## Simple case

$$(X'X)^{-1}X'y = \frac{1}{N_C}\begin{pmatrix} 1 & -1 \\ -1 & \frac{N}{N_T} \end{pmatrix}\begin{pmatrix} \sum_{i=1}^{N} Y_i \\ \sum_{i=1}^{N_T} Y_i \end{pmatrix}$$

$$= \frac{1}{N_C}\begin{pmatrix} \sum_{i=1}^{N} Y_i - \sum_{i=1}^{N_T} Y_i \\ \frac{N}{N_T}\sum_{i=1}^{N_T} Y_i - \sum_{i=1}^{N} Y_i \end{pmatrix}$$

$$= \frac{1}{N_C}\begin{pmatrix} \sum_C Y_i \\ \frac{N}{N_T}\sum_T Y_i - \sum_T Y_i - \sum_C Y_i \end{pmatrix}$$

$$= \frac{1}{N_C}\begin{pmatrix} \sum_C Y_i \\ \frac{N-N_T}{N_T}(\sum_T Y_i) - \sum_C Y_i \end{pmatrix}$$

$$\widehat{\beta} = \begin{pmatrix} \overline{Y_C} \\ \overline{Y_T} - \overline{Y_C} \end{pmatrix}$$

103

## Simple case

$$(X'X)^{-1}X'y = \begin{pmatrix} \overline{Y_C} \\ \overline{Y_T} - \overline{Y_C} \end{pmatrix}$$

104

## Simple case

$$(X'X)^{-1}X'y = \begin{pmatrix} \boxed{\overline{Y_C}} \\ \overline{Y_T} - \overline{Y_C} \end{pmatrix}$$

- The OLS estimate of the intercept is $\overline{Y_C}$

- The coefficient of the treatment dummy is $\overline{Y_T} - \overline{Y_C}$

104

## How precise are these estimates?

- What is the variance of $\widehat{\beta} = (X'X)^{-1}X'y$
- Recall $Y = X\beta + u$
- $\widehat{\beta} = (X'X)^{-1}X'(X\beta + u)$
- $\widehat{\beta} = (X'X)^{-1}X'X\beta + (X'X)^{-1}X'u$
- $\widehat{\beta} = \beta + (X'X)^{-1}X'u$
- If $\mathbb{E}(uX) = 0$
- $V(\widehat{\beta}) = (X'X)^{-1}X'V(u)X(X'X)^{-1}$ [matrix version of $V(b+aY) = a^2Y$]
- If $V(\varepsilon) = \sigma^2 I$ [Homoskedasticity] then
- $V(\widehat{\beta}) = (X'X)^{-1}X'\sigma^2 IX(X'X)^{-1} = \sigma^2(X'X)^{-1}X'X(X'X)^{-1} = \sigma^2(X'X)^{-1}$
- $V(\widehat{\beta}) = \sigma^2\frac{1}{N_T(N-N_T)}\begin{pmatrix} N_T & -N_T \\ -N_T & N \end{pmatrix} = \frac{\sigma^2}{N_C}\begin{pmatrix} 1 & -1 \\ -1 & \frac{N}{N_T} \end{pmatrix}$

105

## How precise are these estimates?

$$V(\widehat{\beta}) = \sigma^2\begin{pmatrix} \frac{1}{N_C} & -\frac{1}{N_C} \\ -\frac{1}{N_C} & \frac{N}{N_T N_C} \end{pmatrix}$$

- Let $N_T = \kappa N$ and $N_C = (1-\kappa)N$

- Since we don't know $\sigma^2$, use $\frac{1}{N-1}(Y - \widehat{Y})^2 = \frac{1}{N-1}(\widehat{u})^2$ as an estimator

106

## Simulations!

```
N=100 #Number of individuals
mu0=1
s.sq=1
beta=0.2
#Let's create potential outcomes
Y0 <-  rnorm(n=N, mean=mu0, sd=s.sq)  # control potential outcome
Y1 <- Y0 + beta   # treatment potential outcome
#Lets randomly assign people to treatment
Z.sim <- rbinom(n=N, size=1, prob=.5) # Do a random assignment
Y.sim <- Y1*Z.sim + Y0*(1-Z.sim) # Reveal outcomes according to assignment

OLS=lm(Y.sim~Z.sim)
summary(OLS)
```

107

$\overline{Y_C} = 80\%$

$\overline{Y_T} - \overline{Y_C} = -0.5 = \widehat{\beta}$

$$V(\widehat{\beta_T}) = \left(\frac{N}{N_T N_C}\right)\sigma^2$$

$$\left(\frac{N_T + N_C}{N_T N_C}\right)\sigma^2 \qquad N_T + N_C = N$$

$$N_T = N_C = \frac{1}{2}N$$

## OLS estimator

---

## How precise are these estimates?

- Standard error of the intercept is: $\sqrt{\hat{\sigma}\dfrac{1}{N_C}}$

- Standard error of the slope is: $\sqrt{\hat{\sigma}\dfrac{N}{N_T N_C}}$

- This should tell us how much our estimates vary on different samples

```
sqrt(sum(OLS$residuals^2)/(N-2)*(1/sum(Z.sim==0)))
sqrt(sum(OLS$residuals^2)/(N-2)*(N/(sum(Z.sim==0)*sum(Z.sim==1))))
```

$$N_C \qquad N_T$$

---

## Simulations!

```
N=100 #Number of individuals
mu0=1
s.sq=1
beta=0.2
Reps=1000
estimate_vector=NULL
for(r in 1:Reps){
    Y0 <-  rnorm(n=N, mean=mu0, sd=s.sq) # control potential outcome
    Y1 <- Y0 + beta # treatment potential outcome
    Z.sim <- rbinom(n=N, size=1, prob=.5) # Do a random assignment
    Y.sim <- Y1*Z.sim + Y0*(1-Z.sim) # Reveal outcomes according to assignment
    OLS=lm(Y.sim~Z.sim)
    beta_estimate=summary(OLS)$coef[2,1]
    estimate_vector=c(estimate_vector,beta_estimate)
}
sd(estimate_vector)
```

---

## Big picture

- We let the computer do the calculations, which are tedious even for small $n$
- Good to know what's going on behind the scenes
- But I honestly do not care if you know how invert a matrix
- Important things in life to understand:
  - What $\hat{\beta}$ is (an estimator of a parameter we do not observe)
  - What the standard error is (the standard deviation of the estimator)
  - What a confidence interval is (an interval where we know with some probability the true estimate lives)
  - What a p-value is (largest probability of obtaining results at least as extreme as those actually observed, under the assumption that the null hypothesis is correct)

---

## OLS Review

Linear algebra review

Law of iterated expectations

OLS basics

Conditional expectation function

"Algebraic" properties of OLS

Properties of OLS estimators

Regression (matrix algebra) with a treatment dummy for the experimental case

Frisch–Waugh–Lovell (FWL) theorem

Regression and causality

---

## OLS Review

Linear algebra review

Law of iterated expectations

OLS basics

Conditional expectation function

"Algebraic" properties of OLS

Properties of OLS estimators

Regression (matrix algebra) with a treatment dummy for the experimental case

Frisch–Waugh–Lovell (FWL) theorem

Regression and causality

## Regression Anatomy Theorem – Frisch–Waugh–Lovell (FWL) theorem

- Assume your main multiple regression model of interest:

$$y_i = \beta_0 + \beta_1 x_{1i} + \cdots + \beta_k x_{ki} + \cdots + \beta_K x_{Ki} + e_i$$

- An auxiliary regression in which the variable $x_{1i}$ is regressed on all the remaining independent variables

$$x_{1i} = \gamma_0 + \gamma_{k-1} x_{k-1i} + \gamma_{k+1} x_{k+1i} + \cdots + \gamma_K x_{Ki} + f_i$$

- $\tilde{x}_{1i} = x_{1i} - \hat{x}_{1i}$ is the residual from the auxiliary regression
- The parameter $\beta_1$ can be rewritten as

$$\beta_1 = \frac{Cov(y_i, x_{1i})}{Var(x_{1i})} = \frac{Cov(y_i, \tilde{x}_{1i})}{Var(\tilde{x}_{1i})}$$

- $\hat{\beta}_1$ is a scaled covariance with the actual data *xor* with the $\tilde{x}_1$ residual

114

## Regression Anatomy Theorem – Frisch–Waugh–Lovell (FWL) theorem II

- Assume your main multiple regression model of interest:

$$y_i = \beta_0 + \beta_1 x_{1i} + \cdots + \beta_k x_{ki} + \cdots + \beta_K x_{Ki} + e_i$$

- Two auxiliary regressions
  - $x_{1i}$ is regressed on all the remaining independent variables

$$x_{1i} = \gamma_0 + \gamma_{k-1} x_{k-1i} + \gamma_{k+1} x_{k+1i} + \cdots + \gamma_K x_{Ki} + f_i$$

  - $y_i i$ is regressed on all the remaining independent variables

$$y_i = \gamma_0 + \alpha_{k-1} x_{k-1i} + \alpha_{k+1} x_{k+1i} + \cdots + \alpha_K x_{Ki} + g_i$$

- $\tilde{x}_{1i} = x_{1i} - \hat{x}_{1i}$ and $\tilde{y}_i = y_i - \hat{y}_i$ residuals from auxiliary regressions
- The parameter $\beta_1$ can be rewritten as

$$\beta_1 = \frac{Cov(y_i, x_{1i})}{Var(x_{1i})} = \frac{Cov(\tilde{y}_i, \tilde{x}_{1i})}{Var(\tilde{x}_{1i})}$$

- $\hat{\beta}_1$ is a scaled covariance with the actual data or with the residuals

115

## Big picture

- Regression anatomy theorem helps us interpret a single slope coefficient in a multiple regression model by the aforementioned decomposition

- Also, help us understand "OLS" as a "matching estimator" (try to compare observations that are alike in the Xs)

116

## OLS Review

117

## OLS Review

118

## Regression and causality

- When is regression causal? Whenever the CEF that regression approximates (or equals if the truth is linear) is causal

- Next: discuss one assumption under which the CEF has a causal interpretation

119

**Potential outcomes - reminder**

- A treatment ($T$) induces two "potential outcomes" for individual $i$
  - The untreated outcome $Y_{0i}$
  - The treated outcome $Y_{1i}$

---

**Potential outcomes - reminder**

- A treatment ($T$) induces two "potential outcomes" for individual $i$
  - The untreated outcome $Y_{0i}$
  - The treated outcome $Y_{1i}$
- The observed outcome

$$Y_i = \begin{cases} Y_{1i} & \text{if } T_i = 1 \\ Y_{0i} & \text{if } T_i = 0 \end{cases}$$
$$= Y_{0i} + (Y_{1i} - Y_{0i})T_i$$

---

**Potential outcomes - reminder**

- A treatment ($T$) induces two "potential outcomes" for individual $i$
  - The untreated outcome $Y_{0i}$
  - The treated outcome $Y_{1i}$
- The observed outcome

$$Y_i = \begin{cases} Y_{1i} & \text{if } T_i = 1 \\ Y_{0i} & \text{if } T_i = 0 \end{cases}$$
$$= Y_{0i} + (Y_{1i} - Y_{0i})T_i$$

- The impact for any individual is $\delta_i = Y_{1i} - Y_{0i}$

---

**Potential outcomes - reminder**

- A treatment ($T$) induces two "potential outcomes" for individual $i$
  - The untreated outcome $Y_{0i}$
  - The treated outcome $Y_{1i}$
- The observed outcome

$$Y_i = \begin{cases} Y_{1i} & \text{if } T_i = 1 \\ Y_{0i} & \text{if } T_i = 0 \end{cases}$$
$$= Y_{0i} + (Y_{1i} - Y_{0i})T_i$$

- The impact for any individual is $\delta_i = Y_{1i} - Y_{0i}$
- Fundamental problem: **Never observe both potential outcomes for the same individual**

---

**We can't just compared treated/untreated individuals**

- We observe $Y_i = Y_{0i} + \underbrace{(Y_{1i} - Y_{0i})}_{\delta_i = \text{impact}} T_i$

- If we compare the outcomes of treated and untreated individuals:

$$\underbrace{\mathbb{E}(Y_i|T_i = 1) - \mathbb{E}(Y_i|T_i = 0)}_{\text{Observed difference}} =$$

---

**We can't just compared treated/untreated individuals**

- We observe $Y_i = Y_{0i} + \underbrace{(Y_{1i} - Y_{0i})}_{\delta_i = \text{impact}} T_i$

- If we compare the outcomes of treated and untreated individuals:

$$\underbrace{\mathbb{E}(Y_i|T_i = 1) - \mathbb{E}(Y_i|T_i = 0)}_{\text{Observed difference}} = \mathbb{E}(Y_{1i}|T_i = 1) - \mathbb{E}(Y_{0i}|T_i = 1) +$$
$$\mathbb{E}(Y_{0i}|T_i = 1) - \mathbb{E}(Y_{0i}|T_i = 0)$$

## We can't just compared treated/untreated individuals

- We observe $Y_i = Y_{0i} + \underbrace{(Y_{1i} - Y_{0i})}_{\delta_i = \text{impact}} T_i$

- If we compare the outcomes of treated and untreated individuals:

$$
\underbrace{\mathbb{E}(Y_i | T_i = 1) - \mathbb{E}(Y_i | T_i = 0)}_{\text{Observed difference}} = \mathbb{E}(Y_{1i} | T_i = 1) - \mathbb{E}(Y_{0i} | T_i = 1) +
$$

$$
\mathbb{E}(Y_{0i} | T_i = 1) - \mathbb{E}(Y_{0i} | T_i = 0)
$$

$$
= \underbrace{\mathbb{E}(Y_{1i} | T_i = 1) - \mathbb{E}(Y_{0i} | T_i = 1)}_{\text{average treatment effect on the treated}} +
$$

$$
\underbrace{\mathbb{E}(Y_{0i} | T_i = 1) - \mathbb{E}(Y_{0i} | T_i = 0)}_{\text{selection bias}}
$$

---

## Unconfoundedness

**Assumption (Unconfoundedness)**
$(Y_{1i}, Y_{0i}) \coprod T_i \mid X_i$

In words:

1. Once we condition on observable characteristics $X_i$, the treatment $T_i$ is as good as randomly assigned

2. Put differently, within the group of individuals with the same characteristics $x_i$, we have a randomized experiment

3. Yet another way of saying it is that conditional on $x_i$, the selection bias disappears

Uncounfoundedness is **fundamentally untestable** and should always be discussed!

---

## Overlap

- In order to exploit the unconfoundedness assumption, for all values of $x_i$ we need to have both treated and untreated units
- Otherwise, either no treatment or no control group for some values of $x_i$
- *Propensity score*, which gives us the probability of $T_i = 1$ given $X_i = x$

$$
p(x) = P(T_i = 1 \mid x_i = x)
$$

- $p(x) = 1$ means that there are no control units (everyone is treated)
- $p(x) = 0$ means that there are no treated units (no one is treated)

**Assumption (Overlap)**
$0 < p(x) < 1$ for all $x$

- In contrast to unconfoundedness, overlap is testable since we can compute $p(x)$ from the data

---

## Identification under unconfoundedness

- How can we identify treatment effects under unconfoundedness?
- Define the conditional mean difference as

$$
\delta_x = E[Y_i \mid T_i = 1, x_i = x] - E[Y_i \mid T_i = 0, x_i = x]
$$

- Conditional on $x_i = x$, we can use the same arguments as the experimental case:

$$
\begin{aligned}
\delta_x &= E[Y_i \mid T_i = 1, x_i = x] - E[Y_i \mid T_i = 0, x_i = x] \\
&= E[Y_{1i} \mid T_i = 1, x_i = x] - E[Y_{0i} \mid T_i = 0, x_i = x] \\
&= E[Y_{1i} \mid x_i = x] - E[Y_{0i} \mid x_i = x] \\
&= E[Y_{1i} - Y_{0i} \mid x_i = x]
\end{aligned}
$$

- The second equality is well-defined for every $x$ by the overlap assumption
- The third equality is by unconfoundedness

---

## Identification under unconfoundedness

- $\delta_x$ is the ATE for individuals with charateristics $x_i = x$

- We can get the (unconditional) ATE as

$$
\begin{aligned}
E[\delta_{x_i}] &= E[E[Y_{1i} - Y_{0i} \mid x_i]] \\
&= E[Y_{1i} - Y_{0i}]
\end{aligned}
$$

---

## Discrete covariates

- The results so far are rather abstract.
- It is easier to understand the results with discrete covariates $x_i$
- In this case,

$$
\text{ATE} = E[Y_{1i} - Y_{0i}] = \sum_x \delta_x P(x_i = x) = E(\delta_x)
$$

- Suppose $x_i$ is binary. In this case the formula becomes:

$$
\begin{aligned}
E[Y_{1i} - Y_{0i}] = \quad & \underbrace{\delta_{x_i=1}}_{\text{mean diff. in group with } x_i = 1} \cdot \underbrace{P(x_i = 1)}_{\text{fract. with } x_i = 1} \\
+ \; & \underbrace{\delta_{x_i=0}}_{\text{mean diff. in group with } x_i = 0} \cdot \underbrace{P(x_i = 0)}_{\text{fract. with } x_i = 0}
\end{aligned}
$$

## An example: causal effect of gender on admissions

| Major | Admissions | Admit | Deny | Total |
|-------|-----------|-------|------|-------|
| A | Men | 400 | 200 | 600 |
| B | Men | 100 | 300 | 400 |
| A | Women | 50 | 50 | 100 |
| B | Women | 300 | 100 | 400 |

- $T_i$ is gender ($T_i = 1$ if male and $T_i = 0$ if female)
- $x_i = M_i$ is choice of major
- Unconfoundedness: gender is independent of admission outcomes conditional on major

## An example: causal effect of gender on admissions

$50\% - 75\% \to \Delta 25PP$
$d\ 50\% +$

$$\delta_A = 400/600 - 50/100 = 0.166$$
$$\delta_B = 100/400 - 300/400 = -0.5$$
$$P(M_i = A) = (600 + 100)/(1500) = 0.466$$
$$P(M_i = B) = (400 + 400)/1500 = 0.533$$
$$ATE = 0.167 \cdot 0.47 + (-0.5) \cdot 0.533 = -0.19$$

## Regression and causality

- Under unconfoundedness, the CEF $E[Y_i \mid T_i, X_i]$ has a causal interpretation.

- Thus, a linear regression model has an (approximate) causal interpretation under unconfoundedness.

## Regression and causality

- If the population regression model is:

$$Y_i = \theta T_i + X_i'\beta + u_i$$

- Then $\delta_x = \theta$ is constant across $x$ and thus $ATE = \theta$

## What if $\delta_x$ is not constant?

| Major | Admissions | Admit | Deny | Total |
|-------|-----------|-------|------|-------|
| A | Men | 400 | 200 | 600 |
| B | Men | 100 | 300 | 400 |
| A | Women | 50 | 50 | 100 |
| B | Women | 300 | 100 | 400 |

## Simulations!

```
#Lets create the data
#600 men in A,  400 men in B, 100 women in A,  400 women in B
#code major A as 1 and B as zeroq()
Data=data.frame(AdminExact=c(rep(1,400),rep(0,200),
                             rep(1,100),rep(0,300),
                             rep(1,50),rep(0,50),
                             rep(1,300),rep(0,100)),
             Major=c(rep(1,600),rep(0,400),rep(1,100),rep(0,400)),
             Treatment=c(rep(1,600),rep(1,400),rep(0,100),rep(0,400)))


summary(lm(AdminExact~Treatment+Major,subset=Major==1,data=Data))
summary(lm(AdminExact~Treatment+Major,subset=Major==0,data=Data))
summary(lm(AdminExact~Treatment+Major,data=Data))
```

## What if $\delta_x$ is not constant?

- Estimate
$$Y_i = \theta T_i + X_i'\beta + u_i$$
- Regression yields $\widehat{\theta} = -0.3 \neq ATE = -0.19$
- In general, we get the following weighted average
$$\theta = \frac{\mathbb{E}\left(\sigma^2_{T_i|X}\delta_x\right)}{\mathbb{E}\left(\sigma^2_{T_i|X_i}\right)}$$
- Regression produces a treatment-variance weighted average of $\delta_x$ (proof in Angrist and Pischke MHE 3.3.1)
- In our case $\sigma^2_{T_i|X_i} = P(T_i = 1|X_i)(1 - P(T_i = 1|X_i))$
  - $\sigma^2_{T_i|Major=A} = \frac{600}{700}\frac{100}{700}$
  - $\sigma^2_{T_i|Major=B} = \frac{400}{800}\frac{400}{800}$

133

## What if $\delta_x$ is not constant?

- Therefore:
$$\widehat{\theta} = \frac{\overbrace{\delta_A}^{0.166}\ \overbrace{\sigma^2_{T_i|Major=A}}^{\frac{600}{700}\frac{100}{700}}\ \overbrace{P(Major = A)}^{\frac{700}{1500}} + \overbrace{\delta_B}^{-0.5}\ \overbrace{\sigma_{T_i|Major=B}}^{\frac{400}{800}\frac{400}{800}}\ \overbrace{P(Major = B)}^{\frac{800}{1500}}}{\sigma^2_{T_i|Major=A}P(Major = A) + \sigma^2_{T_i|Major=B}P(Major = B)}$$
$$= -0.3$$

134

## Big picture

- Beware of what OLS gives you

- Still causal interpretation, even if $\delta_x$ is not constant

- Weighted average of different $\delta_x$

- Weights depend on the variance!

135

## Beyond regression

- Regression is only one method to obtain causal effects under uncounfoundedness

- Other popular methods are: matching and inverse probability weighting

  - Assumption are the same, they generally yield similar results (but implicit weights are different)

- A great review is: *Recent Developments in the Econometrics of Program Evaluation* by Imbens and Wooldrige (2009)

- Check this out: `http://www.nber.org/minicourse3.html`

136

## Some important remarks
## (based on Cyrus Samii's lecture notes)

*For most researchers, the math obscures the assumptions. Without an experiment, a natural experiment, a discontinuity, or some other strong design, no amount of econometric or statistical modeling can make the move from correlation to causation persuasive. (Sekhon, 2009, p. 503)*

- At the end of the day, OLS (and other matching/weighting estimators) "mop up" imbalances that makes CIA plausible
- Thought experiment necessary to test CIA:
  - How could it be that two units that are identical with respect to all meaningful background factors nonetheless receive different treatment?
- Your answer to this question is your source of identification

137