

# Lecture 3: Multiple OLS Regression

*Introduction to Econometrics, Spring 2025*

---

**Zhaopeng Qu**

Business School, Nanjing University

March 12 2025

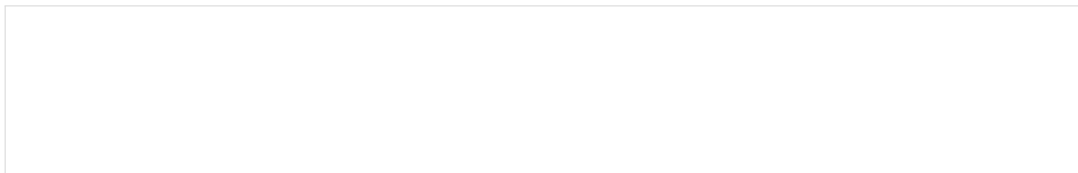


- 1 Review of the Last Lecture**
- 2 Make Comparison Make Sense**
- 3 Multiple OLS Regression: Introduction**
- 4 Multiple OLS Regression: Estimation**
- 5 Multiple OLS Regression Estimators: Partitioned Regression**
- 6 Measures of Fit in Multiple Regression**
- 7 Multiple Regression: Assumption**
- 8 Properties of OLS Estimators in Multiple Regression**

## **Review of the Last Lecture**

# Simple OLS Formula

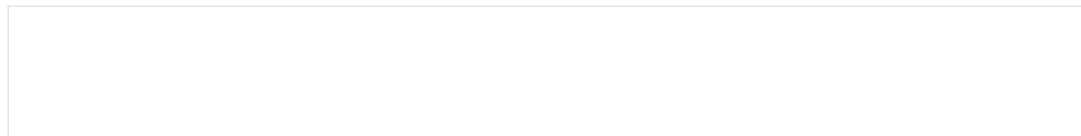
- The linear regression model with one regressor is denoted by



- Where
  - $Y_i$  is the **dependent variable** (Test Score)
  - $X_i$  is the **independent variable** or regressor (Class Size or Student-Teacher Ratio)
  - $u_i$  is the **error term** which contains all the other factors besides  $X$  that determine the value of the dependent variable,  $Y$ , for a specific observation,  $i$ .

# The OLS Estimator

- The estimators of the slope and intercept that **minimize the sum of the squares of  $\hat{u}_i$** , thus



are called the **ordinary least squares (OLS) estimators** of  $\beta_0$  and  $\beta_1$ .

# The OLS Estimator

- The estimators of the slope and intercept that **minimize the sum of the squares of  $\hat{u}_i$** , thus

are called the **ordinary least squares (OLS) estimators** of  $\beta_0$  and  $\beta_1$ .

OLS estimator of  $\beta_1$ :

# Least Squares Assumptions

- Under 3 least squares assumptions,
  1. Assumption 1: ZERO Conditional Mean
  2. Assumption 2: i.i.d. Samples or random sampling
  3. Assumption 3: Without large outliers
- The OLS estimators will be
  1. **unbiased**
  2. **consistent**
  3. **normal sampling distribution**

# Simple OLS Regression v.s. RCT

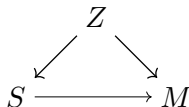
- A simple OLS regression model is a **generalizing continuous version** of RCT assuming three least squares assumptions are held.
- In most observational studies, OLS regression suffers from **selection bias**, which violates the assumption of  $E(u_i|X_i) = 0$ .
- In such cases, OLS estimators are **biased** and **inconsistent**. Therefore the **causal effect** of  $X$  on  $Y$  cannot be identified by simple OLS regression.
- To address the selection bias problem, we have to extend the simple OLS regression model in more general settings.



**Make Comparison Make Sense**

# Case: Smoke and Mortality

- Criticisms from **Ronald A. Fisher**
  - There is no experimental evidence to suggest that smoking is a cause of lung cancer or other serious diseases.
  - Correlation between smoking and mortality may be spurious due to **biased selection** of subjects.



- **Confounder**, *Z*, some other factors, affect on smoking and mortality simultaneously.

## Case: Smoke and Mortality(Cochran 1968)

Table 1: Death rates(死亡率) per 1,000 person-years

Smoking group	Canada	U.K.	U.S.
Non-smokers(不吸烟)	20.2	11.3	13.5
Cigarettes(香烟)	20.5	14.1	13.5
Cigars/pipes(雪茄/烟斗)	35.5	20.7	17.4

- It seems that taking cigars is more hazardous than others to the health.

## Case: Smoke and Mortality(Cochran 1968)

Table 2: Non-smokers and smokers differ in age

Smoking group	Canada	U.K.	U.S.
Non-smokers(不吸烟)	54.9	49.1	57.0
Cigarettes(香烟)	50.5	49.8	53.2
Cigars/pipes(雪茄/烟斗)	65.9	55.7	59.7

- Older people die at a higher rate, and for reasons other than just smoking cigars.
- Perhaps the higher observed death rates among cigar smokers are because **they're older on average.**

## Case: Smoke and Mortality(Cochran 1968)

- The issue is that the ages are **not balanced**; there is a difference in the age distribution between the treatment and control groups.
- let's try to **balance** them, which means to compare mortality rates across the different smoking groups **within age groups** so as to neutralize age imbalances in the observed sample.
- It naturally relates to the concept of **Conditional Expectation Function**.

# Case: Smoke and Mortality(Cochran 1968)

How to balance?

1. Divide the smoking group samples into age groups.
2. For each of the smoking group samples, calculate the mortality rates for the age group.
3. Construct probability weights for each age group as the proportion of the sample with a given age.
4. Compute the **weighted averages** of the age groups mortality rates for each smoking group using the probability weights.

## Case: Smoke and Mortality(Cochran 1968)

	Death rates	Number of	
	Pipe-smokers	Pipe-smokers	Non-smokers
Age 20-50	0.15	11	29
Age 50-70	0.35	13	9
Age +70	0.5	16	2
Total		40	40

- **Question:** What is the average death rate for pipe smokers?

## Case: Smoke and Mortality(Cochran 1968)

	Death rates	Number of	
	Pipe-smokers	Pipe-smokers	Non-smokers
Age 20-50	0.15	11	29
Age 50-70	0.35	13	9
Age +70	0.5	16	2
Total		40	40

- **Question:** What would the average mortality rate be for pipe smokers if they had the same age distribution as the non-smokers?



## Case: Smoke and Mortality(Cochran 1968)

Table 3: Non-smokers and smokers differ in mortality and age

Smoking group	Canada	U.K.	U.S.
Non-smokers(不吸烟)	20.2	11.3	13.5
Cigarettes(香烟)	28.3	12.8	17.7
Cigars/pipes(雪茄/烟斗)	21.2	12.0	14.2

## Case: Smoke and Mortality(Cochran 1968)

Table 3: Non-smokers and smokers differ in mortality and age

Smoking group	Canada	U.K.	U.S.
Non-smokers(不吸烟)	20.2	11.3	13.5
Cigarettes(香烟)	28.3	12.8	17.7
Cigars/pipes(雪茄/烟斗)	21.2	12.0	14.2

- **Conclusion:** It seems that taking cigarettes is most hazardous, and taking pipes is not different from non-smoking.

# Formalization: Covariates

## Definition: Covariates

Variable  $W$  is predetermined with respect to the treatment  $D$  if for each individual  $i$ ,  $W_{0i} = W_{1i}$ , i.e., the value of  $X_i$  does not depend on the value of  $D_i$ . Such characteristics are called *covariates*.

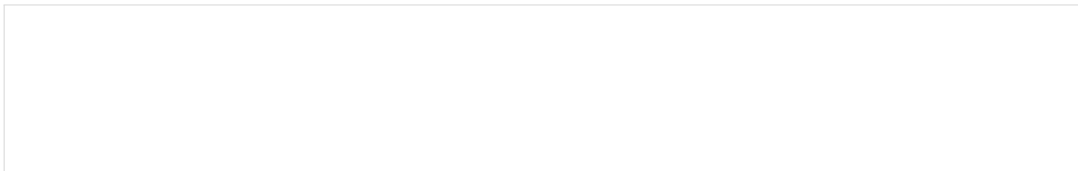
- Covariates are often time invariant (e.g., sex, race), but time invariance is not a necessary condition.

# Identification under Independence

- Recall that randomization in RCTs implies

$$(Y_{0i}, Y_{1i}) \perp\!\!\!\perp D$$

and therefore:

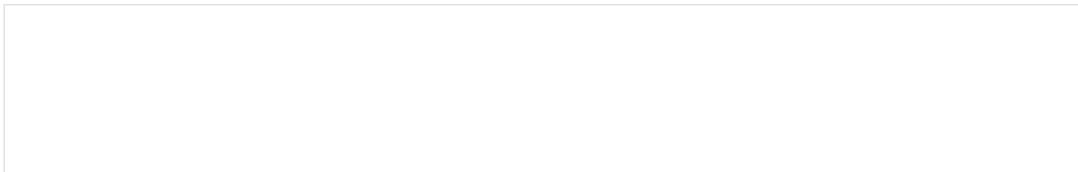


# Identification under Conditional Independence

- **Conditional Independence Assumption(CIA):** which means that if we can “balance” covariates  $X$  then we can take the treatment  $D$  as randomized, thus

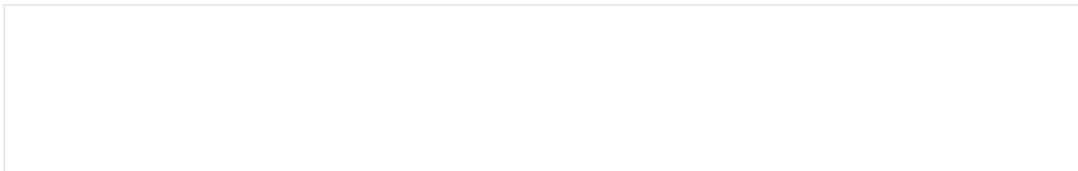
$$(Y_{1i}, Y_{0i}) \perp\!\!\!\perp D | X$$

- Now as  $(Y_{1i}, Y_{0i}) \perp\!\!\!\perp D | X \not\Rightarrow (Y_{1i}, Y_{0i}) \perp\!\!\!\perp D$ ,



# Identification under Conditional Independence(CIA)

- But using the CIA assumption, then



# Curse of Multiple Dimensionality

- Sub-classification in one or two dimensions as Cochran(1968) did in the case of *Smoke and Mortality* is feasible.
- But as the number of covariates we would like to balance grows (like many personal characteristics such as age, gender, education, working experience, married, industries, income, ), then the method become less feasible.
- Assume we have  $k$  covariates and we divide each into 3 coarse categories (e.g., age: young, middle age, old; income: low, medium, high, etc.)
- The number of cells (or groups) is  $3^K$ .
  - If  $k = 10$  then  $3^{10} = 59049$
  - Even if  $k = 6$ , then  $3^6 = 729$ . Assume that we have 1000 observations, then the average number of observations in each cell is less than 2.
- Sub-classification is not a feasible method to balance covariates in high-dimensional space.

# Making Comparison Make Sense

- **Question:** How to make comparison make sense in the presence of covariates?
- *Selection on Observables*
  - Regression
  - Matching
- *Selection on Unobservables*
  - IV, RD, DID, FE and SCM.
- The most fundamental tool among them is **multiple regression**, which compares treatment and control subjects who have the same **observable characteristics in a generalized manner**.



## **Multiple OLS Regression: Introduction**

# Violation of the 1st Least Squares Assumption

- Recall simple OLS regression equation

$$Y_i = \beta_0 + \beta_1 X_i + u_i$$

- **Question:** What does  $u_i$  represent?
  - Answer: contains **all other factors(variables)** which potentially affect  $Y_i$ .
- **Assumption 1**

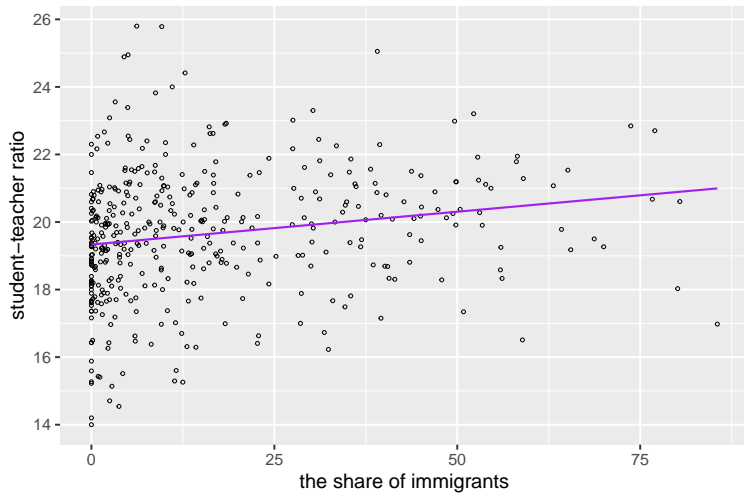
$$E(u_i | X_i) = 0$$

- It states that  $u_i$  are unrelated to  $X_i$  in the sense that, given a value of  $X_i$ , the mean of these other factors equals **zero**.
- But what if  $u_i$  is **correlated** with  $X_i$ ?

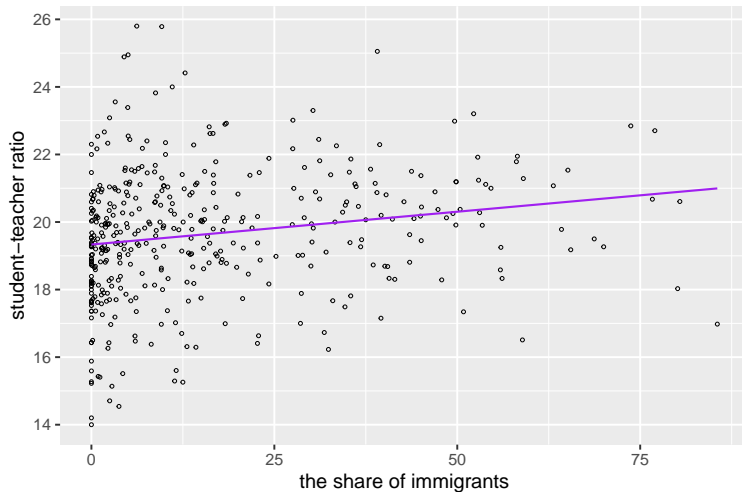
## Example: Class Size and Test Score

- Many other factors can affect student's performance in the school.
- One of factors is **the share of immigrants** in the class. Because immigrant children may have different backgrounds from native children, such as
  - parents' education level
  - family income and wealth
  - parenting style
  - traditional culture

# Scatter Plot: The share of immigrants and STR

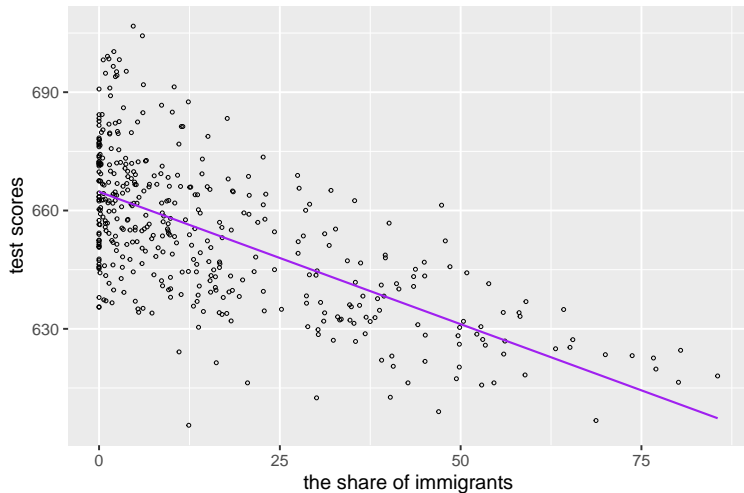


# Scatter Plot: The share of immigrants and STR

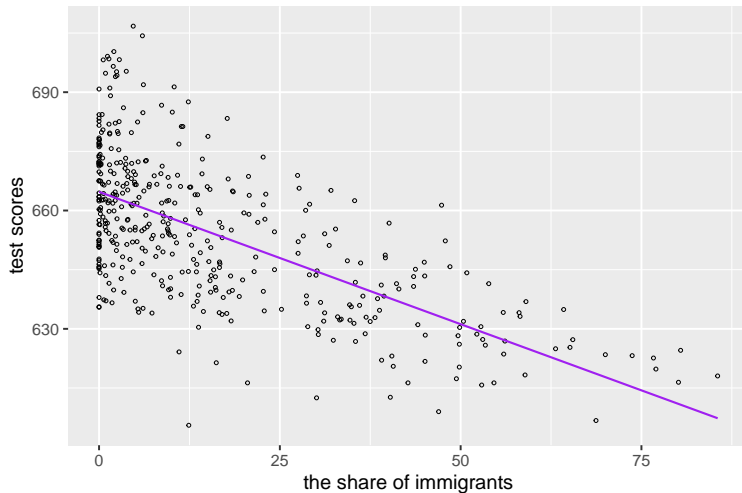


- **higher share of immigrants, bigger class size**

# Scatter Plot: The share of immigrants and STR



# Scatter Plot: The share of immigrants and STR



- **higher share of immigrants, lower test score**

# The share of immigrants as an Omitted Variable

- Class size may be related to percentage of English learners and students who are still learning English likely have lower test scores.
  - In other words, the effect of class size on scores we had obtained in simple OLS may contain *an effect of immigrants on scores*.
- It implies that percentage of English learners is contained in  $u_i$ , in turn that **Assumption 1 is violated**.
  - More precisely, the estimates of  $\hat{\beta}_1$  and  $\hat{\beta}_0$  are **biased and inconsistent**.



# Omitted Variable Bias: Introduction

- As before,  $X_i$  and  $Y_i$  represent **STR** and **Test Score**, respectively.
- Besides,  $W_i$  is the variable which represents **the share of english learners**.
- Suppose that we have no information about it for some reasons, then we have to omit in the regression.
- Thus we have two regressions in mind:
  - **True model**(the Long regression):

$$Y_i = \beta_0 + \beta_1 X_i + \gamma W_i + u_i$$

where  $E(u_i|X_i) = 0$

- **OVB model**(the Short regression):

$$Y_i = \beta_0 + \beta_1 X_i + v_i$$

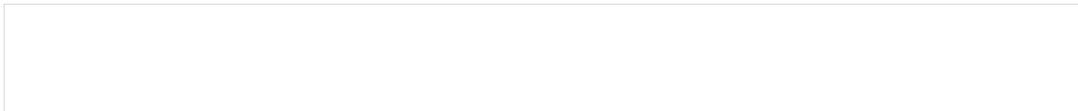
where  $v_i = \gamma W_i + u_i$

# Omitted Variable Bias: Biasedness

- Let us see what is the consequence of OVB

# Omitted Variable Bias: Biasedness

- Let us see what is the consequence of OVB



# Omitted Variable Bias: Biasedness

- Let us see what is the consequence of OVB

- Using the **Law of Iterated Expectation(LIE)** again, we will obtain the following expression(**Skip these steps which are very similar to those for proving unbiasedness of  $\hat{\beta}_1$ , please prove it by yourself**).

# Omitted Variable Bias: Biasedness

- As proving unbiasedness of  $\hat{\beta}_1$ , thus  $E[\hat{\beta}_1] = \beta_1$ , then we need



- Two scenarios:
  1. If  $W_i$  is unrelated to  $X_i$ , then  $E[\hat{\beta}_1] = \beta_1$ .
  2. If  $W_i$  is not determinant of  $Y_i$ , which means that

$$\gamma = 0$$

, then  $E[\hat{\beta}_1] = \beta_1$ , too.

- Only if **both two conditions above are violated simultaneously**, then  $\hat{\beta}_1$  is **biased**, which is normally called **Omitted Variable Bias(OVB)**.

# Omitted Variable Bias(OVB): inconsistency

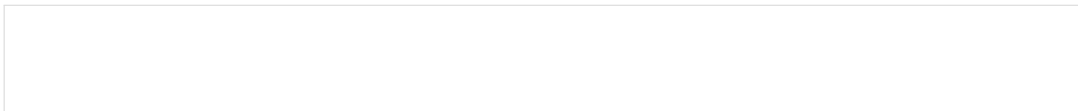
- Recall: simple OLS is consistency when n is large, thus

$$plim\hat{\beta}_1 = \frac{Cov(X_i, Y_i)}{Var(X_i)}$$

# Omitted Variable Bias(OVB): inconsistency

- Recall: simple OLS is consistency when n is large, thus

$$plim\hat{\beta}_1 = \frac{Cov(X_i, Y_i)}{Var(X_i)}$$



# Omitted Variable Bias(OVB): inconsistency

- Thus we obtain



- $\hat{\beta}_1$  is still **consistent**
  - if  $W_i$  is unrelated to  $X$ , thus  $Cov(X_i, W_i) = 0$
  - if  $W_i$  has no effect on  $Y_i$ , thus  $\gamma = 0$
- Only if **both two conditions** above are violated *simultaneously*, then  $\hat{\beta}_1$  is **inconsistent**.



# Omitted Variable Bias(OVB):Directions

- If OVB can be possible in our regressions,then we should guess the **directions** of the bias, in case that we can't eliminate it.
- A summary of the directions of the OVB bias

# Omitted Variable Bias(OVB):Directions

- If OVB can be possible in our regressions, then we should guess the **directions** of the bias, in case that we can't eliminate it.
- A summary of the directions of the OVB bias

---

$$Cov(X_i, W_i) > 0$$

$$Cov(X_i, W_i) < 0$$

---

$$\gamma > 0$$

# Omitted Variable Bias(OVB):Directions

- If OVB can be possible in our regressions, then we should guess the **directions** of the bias, in case that we can't eliminate it.
- A summary of the directions of the OVB bias

---

$$Cov(X_i, W_i) > 0$$

$$Cov(X_i, W_i) < 0$$

---

$$\gamma > 0$$

**Positive bias**

# Omitted Variable Bias(OVB):Directions

- If OVB can be possible in our regressions, then we should guess the **directions** of the bias, in case that we can't eliminate it.
- A summary of the directions of the OVB bias

---

	$Cov(X_i, W_i) > 0$	$Cov(X_i, W_i) < 0$
$\gamma > 0$	Positive bias	Negative bias
$\gamma < 0$		

# Omitted Variable Bias(OVB):Directions

- If OVB can be possible in our regressions, then we should guess the **directions** of the bias, in case that we can't eliminate it.
- A summary of the directions of the OVB bias

---

	$Cov(X_i, W_i) > 0$	$Cov(X_i, W_i) < 0$
$\gamma > 0$	Positive bias	Negative bias
$\gamma < 0$	Negative bias	

# Omitted Variable Bias(OVB):Directions

- If OVB can be possible in our regressions, then we should guess the **directions** of the bias, in case that we can't eliminate it.
- A summary of the directions of the OVB bias

---

	$Cov(X_i, W_i) > 0$	$Cov(X_i, W_i) < 0$
$\gamma > 0$	Positive bias	Negative bias
$\gamma < 0$	Negative bias	Positive bias

---

# Omitted Variable Bias: Examples

- **Question:** If we omit following variables, then what are the directions of these biases? and why?
  1. Time of day of the test[suppose morning(8:00-12:00am) is better,afternoon(13:00-17:00pm) is worse]
  2. The number of dormitories
  3. Teachers' salary
  4. Family income
  5. Percentage of English learners(the share of immigrants)

# Omitted Variable Bias: Examples in R

- Regress *Testscore* on *Class size*

```
#>
#> Call:
#> lm(formula = testscr ~ str, data = ca)
#>
#> Residuals:
#>      Min       1Q   Median       3Q      Max
#> -47.727 -14.251   0.483  12.822  48.540
#>
#> Coefficients:
#>              Estimate Std. Error t value Pr(>|t|)
#> (Intercept)  698.9330     9.4675   73.825 < 2e-16 ***
#> str          -2.2798     0.4798   -4.751 2.78e-06 ***
#> ---
#> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#>
#> Residual standard error: 18.58 on 418 degrees of freedom
#> Multiple R-squared:  0.05124,    Adjusted R-squared:  0.04897
#> F-statistic: 22.58 on 1 and 418 DF,  p-value: 2.783e-06
```



# Omitted Variable Bias: Examples in R

- Regress *Testscore* on *Class size* and *the percentage of English learners*

```
#>
#> Call:
#> lm(formula = testscr ~ str + el_pct, data = ca)
#>
#> Residuals:
#>      Min       1Q   Median       3Q      Max
#> -48.845 -10.240  -0.308   9.815  43.461
#>
#> Coefficients:
#>              Estimate Std. Error t value Pr(>|t|)
#> (Intercept)  686.03225     7.41131   92.566 < 2e-16 ***
#> str          -1.10130     0.38028   -2.896  0.00398 **
#> el_pct       -0.64978     0.03934  -16.516 < 2e-16 ***
#> ---
#> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#>
#> Residual standard error: 14.46 on 417 degrees of freedom
#> Multiple R-squared:  0.4264, Adjusted R-squared:  0.4237
```

# Omitted Variable Bias: Examples in R

<i>Dependent variable:</i>		
testscr		
	(1)	(2)
str	-2.280*** (0.480)	-1.101*** (0.380)
el_pct		-0.650*** (0.039)
Constant	698.933*** (9.467)	686.032*** (7.411)
Observations	420	420
R <sup>2</sup>	0.051	0.426

Note: \* p<0.1; \*\* p<0.05; \*\*\* p<0.01

- OVB is the **most common** bias when we run OLS regressions using nonexperimental data.
- OVB means that there are some variables which should have been included in the regression but actually was not.
- Then the simplest way to overcome OVB: *Put omitted the variable into the right side of the regression*, which means our regression model should be

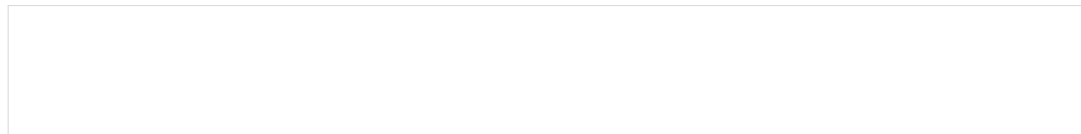
$$Y_i = \beta_0 + \beta_1 X_i + \gamma W_i + u_i$$

- The strategy can be denoted as **controlling** informally, which introduces the more general regression model: **Multiple OLS Regression**.

## **Multiple OLS Regression: Estimation**

# Multiple regression model with k regressors

- The multiple regression model is



where

- $Y_i$  is the **dependent variable**
- $X_1, X_2, \dots, X_k$  are the **independent variables**(includes one is our of interest and **some control variables**)
- $\beta_j, j = 1 \dots k$  are slope coefficients on  $X_j$  corresponding.
- $\beta_0$  is the estimate *intercept*, the value of Y when all  $X_j = 0, j = 1 \dots k$
- $u_i$  is the regression *error term*, still all other factors affect outcomes.

# Interpretation of coefficients $\beta_j, j = 1 \dots k$

- $\beta_j$  is **partial (marginal) effect** of  $X_j$  on  $Y$ .

- $\beta_j$  is also partial (marginal) effect of  $E[Y_i|X_1 \dots X_k]$ .

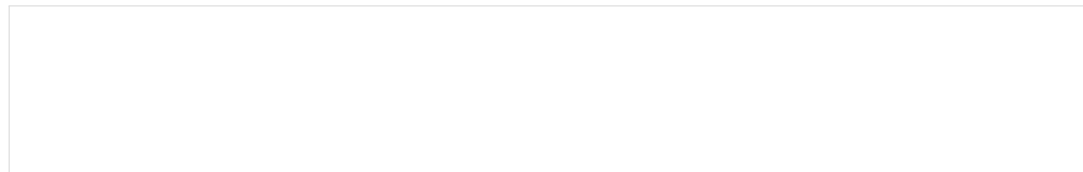
- it does mean that we are estimate the effect of  $X$  on  $Y$  when “**other things equal**”, thus the concept of **ceteris paribus**.

# OLS Estimation in Multiple Regressors

- As in a **Simple OLS Regression**, the estimators of **Multiple OLS Regression** is just a minimize the following question

# OLS Estimation in Multiple Regressors

- As in a **Simple OLS Regression**, the estimators of **Multiple OLS Regression** is just a minimize the following question



where  $b_0 = \hat{\beta}_1, b_1 = \hat{\beta}_2, \dots, b_k = \hat{\beta}_k$  are estimators.

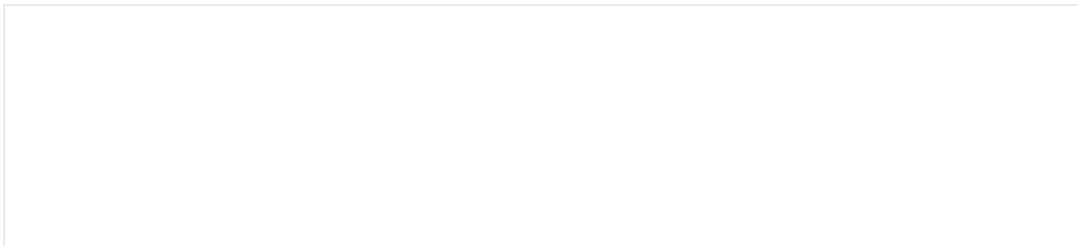


# OLS Estimation in Multiple Regressors

- Similarly in Simple OLS, based on F.O.C, the multiple OLS estimators  $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$  are obtained by solving the following **system of normal equations**

# OLS Estimation in Multiple Regressors

- Similarly in Simple OLS, based on F.O.C, the multiple OLS estimators  $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$  are obtained by solving the following **system of normal equations**



# OLS Estimation in Multiple Regressors

- Similar to in Simple OLS, the fitted residuals are

- Therefore, the normal equations also can be written as

## Multiple OLS Regression Estimators: Partitioned Regression

# Partitioned regression: OLS estimators

- The Multiple OLS estimators of  $\hat{\beta}_j$  could be obtained by running **partitioned regression**, which can let us exempt from the matrix algebra.
- We can obtain OLS estimators of  $\beta_j$ ;  $j = 1, 2 \dots k$  in following 3 steps:

- The last step implies that the OLS estimator of  $\beta_j$  can be expressed as follows

# Partitioned regression: OLS estimators

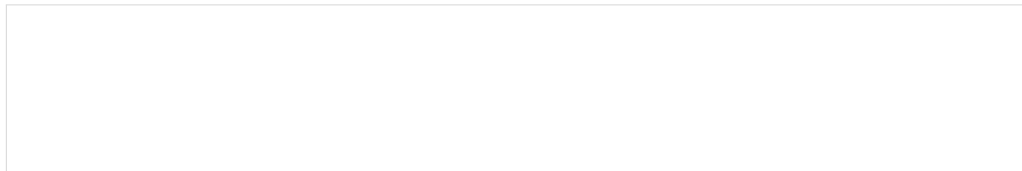
- Suppose we want to obtain an expression for  $\hat{\beta}_1$ .
  - the first step: regress  $X_{1,i}$  on other regressors  $X_s$ , thus

$$X_{1,i} = \gamma_0 + \gamma_2 X_{2,i} + \dots + \gamma_k X_{k,i} + v_i$$

- the second step: obtain the residuals from the regression above, denoted as  $\tilde{X}_{1,i} = \hat{v}_{1i}$ , thus

$$X_{1,i} = \hat{\gamma}_0 + \hat{\gamma}_2 X_{2,i} + \dots + \hat{\gamma}_k X_{k,i} + \tilde{X}_{1,i}$$

- Then we could prove that



# Proof of Partitioned regression result(1)

- Recall  $u_i$  are the residuals for the Multiple OLS regression equation, thus

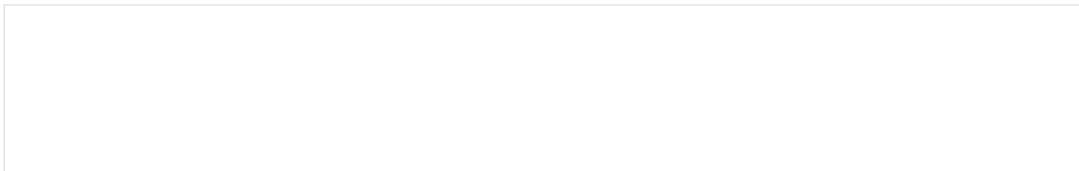
$$Y_i = \hat{\beta}_0 + \hat{\beta}_1 X_{1,i} + \hat{\beta}_2 X_{2,i} + \dots + \hat{\beta}_k X_{k,i} + \hat{u}_i$$

- Then we have

- Likewise,  $\tilde{X}_{1i}$  are the residuals for the partitioned regression equation on  $X_{2i}, \dots, X_{ki}$ , then we have

- Additionally, because  $\tilde{X}_{1,i} = X_{1,i} - \hat{\gamma}_0 - \hat{\gamma}_2 X_{2,i} - \dots - \hat{\gamma}_k X_{k,i}$ , then

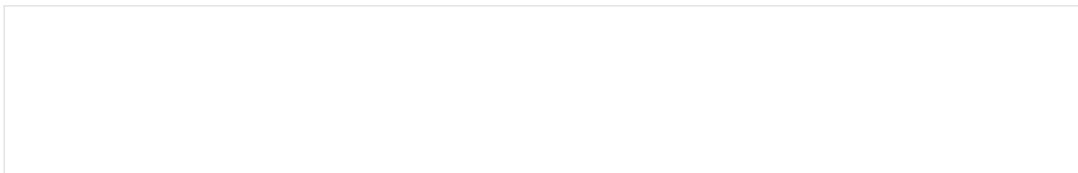
# Proof of Partitioned regression result(2)





# Proof of Partitioned regression result(3)

- We will see



- Then

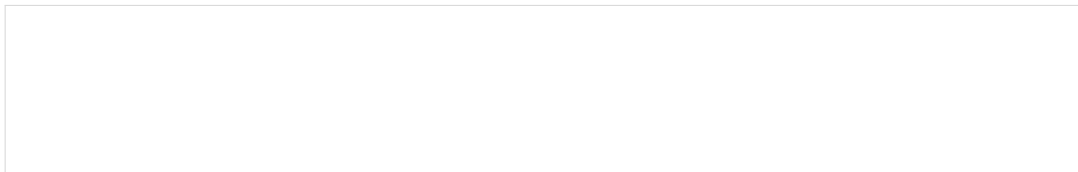


# Frisch-Waugh-Lowell Theorem

## FWL Theorem

# The Intuition of FWL Theorem

## Partialling Out



- FWL Theorem provides a new and important perspective to understand the multiple OLS estimator.

# Test Scores and Student-Teacher Ratios(1)

- Now we put one additional control variables into our OLS regression model

$$Testscore = \beta_0 + \beta_1 STR + \beta_2 elpct + u_i$$

- `elpct`: the share of English learners as an indicator for the share of immigrants.
- We want to know the effect of `STR` on `testscr` after controlling for `elpct`.
- Two steps:
  - First, we regress `str` on `elpct` and keep the residuals  $\widetilde{STR}$ , thus

$$STR = \hat{\gamma}_0 + \hat{\gamma}_1 elpct + \widetilde{STR}$$

- Second, we regress `testscr` on  $\widetilde{STR}$  to get the effect of `STR` after controlling for `elpct`.

## Test Scores and Student-Teacher Ratios(2)

- The residuals of the regression of `str` on `elpct` are

$$\widetilde{STR} = STR - \widehat{STR} = STR - (\hat{\gamma}_0 + \hat{\gamma}_1 \text{elpct})$$

- Check whether the sum of  $\widetilde{STR}$ ,  $\widetilde{STR} \times \text{elpct}$  and  $\widetilde{testscr} \times \widetilde{STR}$  are zero.

# Test Scores and Student-Teacher Ratios(2)

- The residuals of the regression of `str` on `el_pct` are

$$\widetilde{STR} = STR - \widehat{STR} = STR - (\hat{\gamma}_0 + \hat{\gamma}_1 elpct)$$

- Check whether the sum of  $\widetilde{STR}$ ,  $\widetilde{STR} \times elpct$  and  $\widetilde{testscr} \times \widetilde{STR}$  are zero.

```
tilde.str <- residuals(lm(str ~ el_pct, data=ca))
tilde.score <- residuals(lm(testscr ~ str+el_pct, data=ca))
sum(tilde.str) # also is zero
```

```
#> [1] -8.104628e-15
```

```
sum(tilde.str*ca$el_pct) # also should be zero
```

```
#> [1] -3.896883e-13
```

```
sum(tilde.score*tilde.str) # also should be zero
```

```
#> [1] 1.275424e-12
```

# Test Scores and Student-Teacher Ratios(3)

- Multiple OLS estimator in a partitioned way

$$\hat{\beta}_j = \frac{\sum_{i=1}^n \tilde{X}_{j,i} Y_i}{\sum_{i=1}^n \tilde{X}_{j,i}^2} \text{ for } j = 1, 2, \dots, k$$

# Test Scores and Student-Teacher Ratios(3)

- Multiple OLS estimator in a partitioned way

$$\hat{\beta}_j = \frac{\sum_{i=1}^n \tilde{X}_{j,i} Y_i}{\sum_{i=1}^n \tilde{X}_{j,i}^2} \text{ for } j = 1, 2, \dots, k$$

```
sum(tilde.str*ca$testscr) / sum(tilde.str^2)
```

```
#> [1] -1.101296
```



## Test Scores and Student-Teacher Ratios(4)

```
reg3 <- lm(testscr ~ tilde.str,data = ca)
summary(reg3)
```

```
#>
#> Call:
#> lm(formula = testscr ~ tilde.str, data = ca)
#>
#> Residuals:
#>      Min       1Q   Median       3Q      Max
#> -48.693 -14.124   0.988  13.209  50.872
#>
#> Coefficients:
#>              Estimate Std. Error t value Pr(>|t|)
#> (Intercept)  654.1565     0.9254  706.864  <2e-16 ***
#> tilde.str    -1.1013     0.4986   -2.209   0.0277 *
#> ---
#> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

# Test Scores and Student-Teacher Ratios(5)

```
reg4 <- lm(testscr ~ str+el_pct,data = ca)
summary(reg4)
```

```
#>
#> Call:
#> lm(formula = testscr ~ str + el_pct, data = ca)
#>
#> Residuals:
#>      Min       1Q   Median       3Q      Max
#> -48.845 -10.240  -0.308   9.815  43.461
#>
#> Coefficients:
#>              Estimate Std. Error t value Pr(>|t|)
#> (Intercept)  686.03225    7.41131   92.566 < 2e-16 ***
#> str          -1.10130    0.38028   -2.896  0.00398 **
#> el_pct       -0.64978    0.03934  -16.516 < 2e-16 ***
#> ---
```

# Test Scores and Student-Teacher Ratios(6)

<i>Dependent variable:</i>		
testscr		
	(1)	(2)
tilde.str	-1.101** (0.499)	
str		-1.101*** (0.380)
el_pct		-0.650*** (0.039)
Constant	654.157*** (0.925)	686.032*** (7.411)
Observations	420	420
Adjusted R <sup>2</sup>	0.009	0.424

Note: \* p<0.1; \*\* p<0.05; \*\*\* p<0.01

## Measures of Fit in Multiple Regression

## Recall: Measures of Fit: The $R^2$

- Decompose  $Y_i$  into the fitted value plus the residual  $Y_i = \hat{Y}_i + \hat{u}_i$
- The **total sum of squares (TSS)**:  $TSS = \sum_{i=1}^n (Y_i - \bar{Y})^2$
- The **explained sum of squares (ESS)**:  $\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$
- The **sum of squared residuals (SSR)**:  $\sum_{i=1}^n (\hat{Y}_i - Y_i)^2 = \sum_{i=1}^n \hat{u}_i^2$
- And

$$TSS = ESS + SSR$$

- The regression  $R^2$  is the fraction of the sample variance of  $Y_i$  explained by (or predicted by) the regressors.

$$R^2 = \frac{ESS}{TSS} = 1 - \frac{SSR}{TSS}$$

# Measures of Fit in Multiple Regression

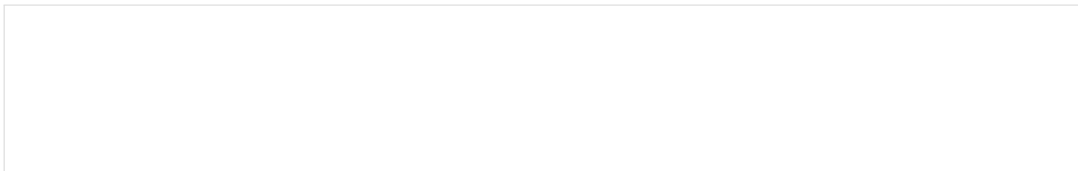
- When you put more variables into the regression, then  $R^2$  **always increases** when you *add another regressor*. Because in general the SSR will decrease.

- Consider two models

$$Y_i = \beta_0 + \beta_1 X_{1i} + u_i$$

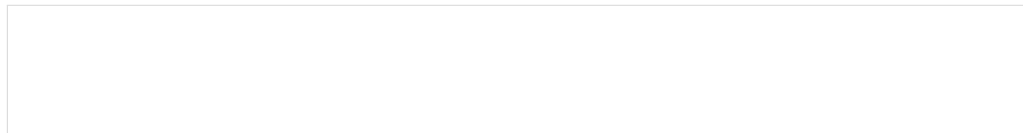
$$Y_i = \tilde{\beta}_0 + \tilde{\beta}_1 X_{1i} + \tilde{\beta}_2 X_{2i} + v_i$$

- Recall: about two residuals  $\hat{u}_i$  and  $\hat{v}_i$ , we have



# Measures of Fit in Multiple Regression

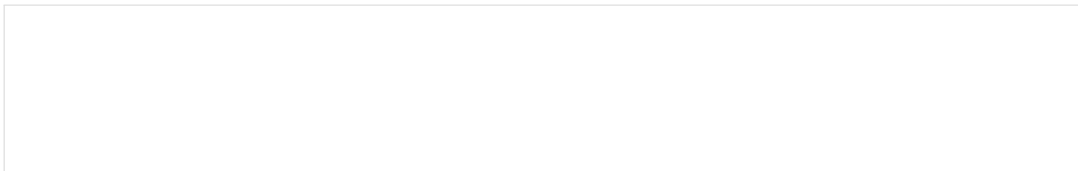
- we will show that



- therefore  $R_v^2 \geq R_u^2$ , thus  $R^2$  that corresponds to the regression with one regressor is **less or equal** than  $R^2$  that corresponds to the regression with two regressors.
- This conclusion can be generalized to the case of  $k + 1$  regressors.

# Measures of Fit in Multiple Regression

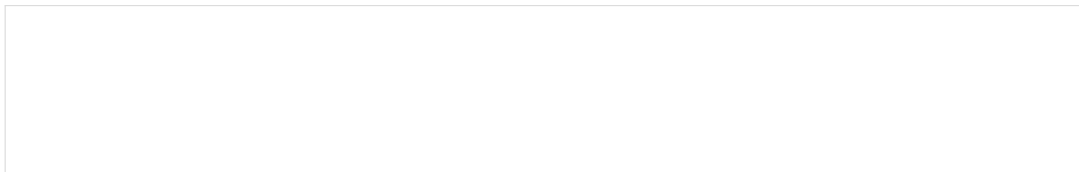
- At first we would like to know  $\sum_{i=1}^n \hat{u}_i \hat{v}_i$





# Measures of Fit in Multiple Regression

- Then we can obtain



- Therefore  $R_v^2 \geq R_u^2$ , thus  $R^2$  the regression with one regressor is **less or equal** than  $R^2$  that corresponds to the regression with two regressors.

# Measures of Fit: The Adjusted $R^2$

- the Adjusted  $R^2$ , is a modified version of the  $R^2$  that does not necessarily increase when a new regressor is added.

$$\overline{R^2} = 1 - \frac{n-1}{n-k-1} \frac{SSR}{TSS} = 1 - \frac{s_{\hat{u}}^2}{s_Y^2}$$

- because  $\frac{n-1}{n-k-1}$  is always greater than 1, so  $\overline{R^2} < R^2$
- adding a regressor has two opposite effects on the  $\overline{R^2}$ .
- $\overline{R^2}$  can be negative.
- **Remind:** *neither  $R^2$  nor  $\overline{R^2}$  is not the golden criterion for good or bad OLS estimation.*

# Example: Test scores and Student Teacher Ratios

```
1 . reg testscr str el_pct
```

Source	SS	df	MS	Number of obs	=	420
Model	64864.3011	2	32432.1506	F(2, 417)	=	155.01
Residual	87245.2925	417	209.221325	Prob > F	=	0.0000
Total	152109.594	419	363.030056	R-squared	=	0.4264
				Adj R-squared	=	0.4237
				Root MSE	=	14.464

testscr	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
str	-1.101296	.3802783	-2.90	0.004	-1.848797	-.3537945
el_pct	-.6497768	.0393425	-16.52	0.000	-.7271112	-.5724423
_cons	686.0322	7.411312	92.57	0.000	671.4641	700.6004

## Multiple Regression: Assumption

# Multiple Regression: Assumption

- Assumption 1: The conditional distribution of  $u_i$  given  $X_{1i}, \dots, X_{ki}$  has mean zero, thus

$$E[u_i | X_{1i}, \dots, X_{ki}] = 0$$

- which is a very strong assumption, which means  $u_i$  is uncorrelated with all the independent variables.
- Assumption 2:  $(Y_i, X_{1i}, \dots, X_{ki})$  are i.i.d.
- Assumption 3: Large outliers are unlikely.
- At last, we have to add one more assumption for multiple regression.
  - Assumption 4: **No perfect multicollinearity.**

# Perfect multicollinearity

- **Perfect multicollinearity** arises when one of the regressors is a **perfect** linear combination of the other regressors.
- If you include a full set of binary variables (a complete and mutually exclusive categorization) and an intercept in the regression, you will have perfect multicollinearity.
  - eg. female and male = 1-female
- This is called the **dummy variable trap**.
- Solutions to the dummy variable trap:
  - Omit one of the groups or the intercept

# Categorized Variable as Independent Variables

- Recall if  $X$  is a dummy variable, then we can put it into regression equation straightly.
- What if  $X$  is a categorical variable?
  - **Question:** What is a categorical variable?
- For example, we may define  $D_i$  as follows:

# Categorized Variable as Independent Variables

- Recall if  $X$  is a dummy variable, then we can put it into regression equation straightly.
- What if  $X$  is a categorical variable?
  - **Question:** What is a categorical variable?
- For example, we may define  $D_i$  as follows:

$$D_i = \begin{cases} 1 & \text{small-size class if } STR \text{ in } i^{th} \text{ school district} < 18 \\ 2 & \text{middle-size class if } 18 \leq STR \text{ in } i^{th} \text{ school district} < 22 \\ 3 & \text{large-size class if } STR \text{ in } i^{th} \text{ school district} \geq 22 \end{cases} \quad (4.5)$$



# A Special Case: Categorical Variable as $X$

- Naive Solution: a simple OLS regression model

$$TestScore_i = \beta_0 + \beta_1 D_i + u_i$$

- **Question:** Can you explain the meaning of estimate coefficient  $\beta_1$ ?
- **Answer:** It does not make sense that the coefficient of  $\beta_1$  can be explained as continuous variables.

## A Special Case: Categorical Variables as $X$

- The first step: turn a categorical variable ( $D_i$ ) into multiple dummy variables ( $D_{1i}, D_{2i}, D_{3i}$ )

## A Special Case: Categorical Variables as $X$

- We put these dummies into a multiple regression

$$TestScore_i = \beta_0 + \beta_1 D_{1i} + \beta_2 D_{2i} + \beta_3 D_{3i} + u_i \quad (4.6)$$

- Then as a dummy variable as the independent variable in a simple regression  
The coefficients  $(\beta_1, \beta_2, \beta_3)$  represent the effect of every categorical class on *testscore* respectively.

## A Special Case: Categorical Variables as $X$

- In practice, we can't put all dummies into the regression, but only have  $n - 1$  dummies unless we will suffer **perfect multi-collinearity**.
- The regression may be like as

$$TestScore_i = \beta_0 + \beta_1 D_{1i} + \beta_2 D_{2i} + u_i \quad (4.6)$$

- The default intercept term,  $\beta_0$ , represents the large-sized class. Then, the coefficients  $(\beta_1, \beta_2)$  represent *testscore* gaps between small\_sized, middle-sized class and large-sized class, respectively.

# Perfect multicollinearity

- regress *Testscore* on *Class size* and the *percentage of English learners*

```
#>
#> Call:
#> lm(formula = testscr ~ str + el_pct, data = ca)
#>
#> Residuals:
#>      Min       1Q   Median       3Q      Max
#> -48.845 -10.240  -0.308   9.815  43.461
#>
#> Coefficients:
#>              Estimate Std. Error t value Pr(>|t|)
#> (Intercept)  686.03225     7.41131   92.566 < 2e-16 ***
#> str          -1.10130     0.38028  -2.896  0.00398 **
```

# Perfect multicollinearity

- add a new variable `nel=1-el_pct` into the regression

```
#>
#> Call:
#> lm(formula = testscr ~ str + nel_pct + el_pct, data = ca)
#>
#> Residuals:
#>      Min       1Q   Median       3Q      Max
#> -48.845 -10.240  -0.308   9.815  43.461
#>
#> Coefficients: (1 not defined because of singularities)
#>              Estimate Std. Error t value Pr(>|t|)
#> (Intercept)  685.38247    7.41556   92.425 < 2e-16 ***
#> str          -1.10130    0.38028  -2.896  0.00398 **
#> nel_pct       0.64978    0.03934  16.516 < 2e-16 ***
#> el_pct                NA           NA      NA      NA
#> ---
#> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#>
#> Residual standard error: 14.46 on 417 degrees of freedom
```

# Perfect multicollinearity

<i>Dependent variable:</i>		
testscr		
	(1)	(2)
str	-1.101*** (0.380)	-1.101*** (0.380)
nel_pct		0.650*** (0.039)
el_pct	-0.650*** (0.039)	
Constant	686.032*** (7.411)	685.382*** (7.416)
Observations	420	420
Adjusted R <sup>2</sup>	0.424	0.424

Note: \* p<0.1; \*\* p<0.05; \*\*\* p<0.01

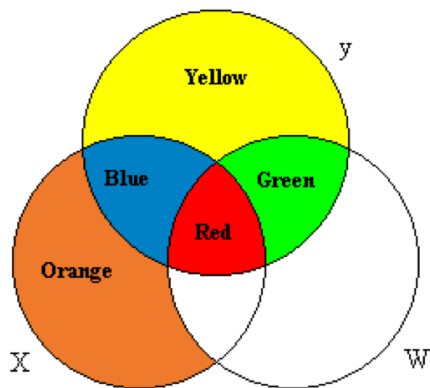
# Multicollinearity

**Multicollinearity** means that two or more regressors are **highly** correlated, but one regressor is **NOT** a perfect linear function of one or more of the other regressors.

- **multicollinearity** is **NOT** a violation of OLS assumptions.
  - It does not impose theoretical problem for the calculation of OLS estimators.
- But if two regressors are highly correlated, then the the coefficient on at least one of the regressors is imprecisely estimated (high variance).
- To what extent two correlated variables can be seen as “highly correlated”?
  - **rule of thumb:** correlation coefficient is over **0.8**.

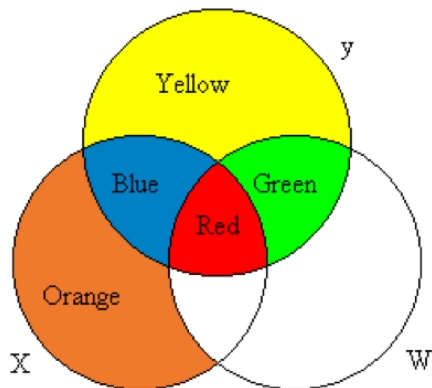


# Venn Diagrams for Multiple Regression Model

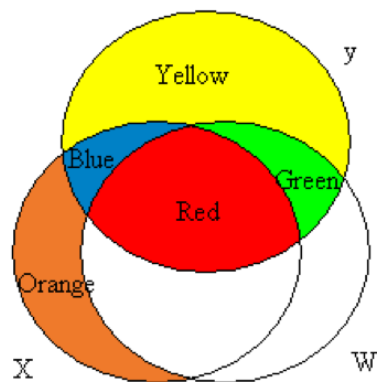


- In a simple model (y on X), OLS uses 'Blue' + 'Red' to estimate  $\beta$ .
- When y is regressed on X and W: OLS throws away the red area and just uses blue to estimate  $\beta$ .
- Idea: Red area is contaminated (we do not know if the movements in y are due to X or to W).

# Venn Diagrams for Multicollinearity



**Figure 3a Modest collinearity**



**Figure 3b Considerable collinearity**

# Multiple Regression: Test Scores and Class Size

- New variable: expense per student `expn_stu`

```
cor(ca$str, ca$el_pct)
```

```
[1] 0.1876424
```

```
cor(ca$str, ca$expn_stu)
```

```
[1] -0.6199821
```

# Multiple Regression: Test Scores and Class Size

Table 5: Class Size and Test Score

	testscr		
	(1)	(2)	(3)
str	-2.280*** (0.480)	-1.101*** (0.380)	-0.286 (0.481)
el_pct		-0.650*** (0.039)	-0.656*** (0.039)
expn_stu			0.004*** (0.001)
Constant	698.933*** (9.467)	686.032*** (7.411)	649.578*** (15.206)
N	420	420	420
Adjusted R <sup>2</sup>	0.049	0.424	0.433

Notes:

\*\*\* Significant at the 1 percent level.

\*\* Significant at the 5 percent level.

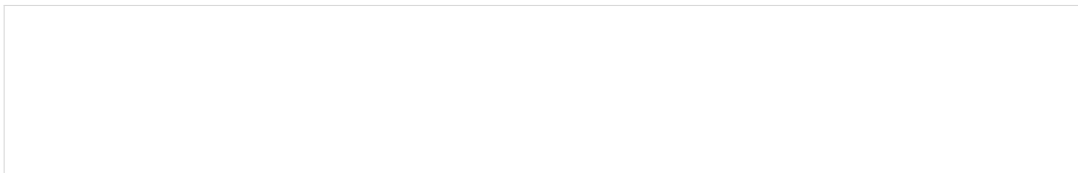
## Properties of OLS Estimators in Multiple Regression

# Properties of OLS estimators: Unbiasedness(1)

- Use partitioned regression formula

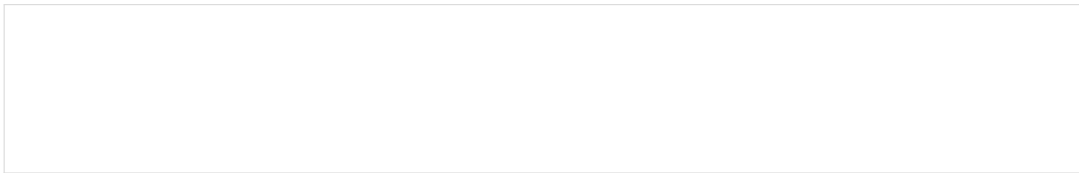
$$\hat{\beta}_1 = \frac{\sum_{i=1}^n \tilde{X}_{1,i} Y_i}{\sum_{i=1}^n \tilde{X}_{1,i}^2}$$

- Substitute  $Y_i = \beta_0 + \beta_1 X_{1,i} + \beta_2 X_{2,i} + \dots + \beta_k X_{k,i} + u_i, i = 1, \dots, n$ , then



# Properties of OLS estimators: Unbiasedness(2)

- Because



- Therefore

$$\hat{\beta}_1 = \beta_1 + \frac{\sum_{i=1}^n \tilde{X}_{1,i} u_i}{\sum_{i=1}^n \tilde{X}_{1,i}^2}$$

# Properties of OLS estimators: Unbiasedness(3)

- Recall **Assumption 1**:  $E[u_i | X_{1i}, X_{2i} \dots X_{ki}] = 0$  and  $\tilde{X}_{1i}$  is a function of  $X_{2i} \dots X_{ki}$
- Then take expectations of  $\hat{\beta}_1$  and **The Law of Iterated Expectations** again

- Identical argument works for  $\beta_2, \dots, \beta_k$ , thus

$$E[\hat{\beta}_j] = \beta_j \text{ where } j = 1, 2, \dots, k$$



# Properties of OLS estimators: Consistency(1)

- Recall

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n \tilde{X}_{1,i} Y_i}{\sum_{i=1}^n \tilde{X}_{1,i}^2}$$

- Similar to the proof in the Simple OLS Regression, thus

# Properties of OLS estimators: Consistency(1)

- Recall

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n \tilde{X}_{1,i} Y_i}{\sum_{i=1}^n \tilde{X}_{1,i}^2}$$

- Similar to the proof in the Simple OLS Regression, thus

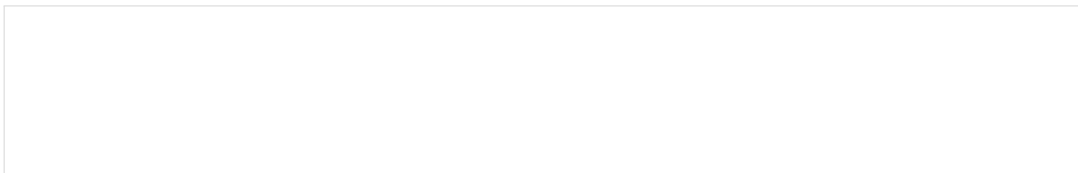
where  $s_{\tilde{X}_1 Y}$  and  $s_{\tilde{X}_1}^2$  are the sample covariance of  $\tilde{X}_1$  and  $Y$  and the sample variance of  $\tilde{X}_1$ .

## Properties of OLS estimators: Consistency(2)

- Base on L.L.N(the law of large numbers) and random sample(i.i.d)

## Properties of OLS estimators: Consistency(2)

- Base on L.L.N(the law of large numbers) and random sample(i.i.d)



- Combining with *Continuous Mapping Theorem*,then we obtain the partitioned multiple OLS estimator  $\hat{\beta}_1$ ,when  $n \rightarrow \infty$

$$plim\hat{\beta}_1 =$$

## Properties of OLS estimators: Consistency(2)

- Base on L.L.N(the law of large numbers) and random sample(i.i.d)

- Combining with *Continuous Mapping Theorem*,then we obtain the partitioned multiple OLS estimator  $\hat{\beta}_1$ ,when  $n \rightarrow \infty$

$$plim \hat{\beta}_1 = plim \left( \frac{s_{\tilde{X}_1 Y}}{s_{\tilde{X}_1}^2} \right) =$$

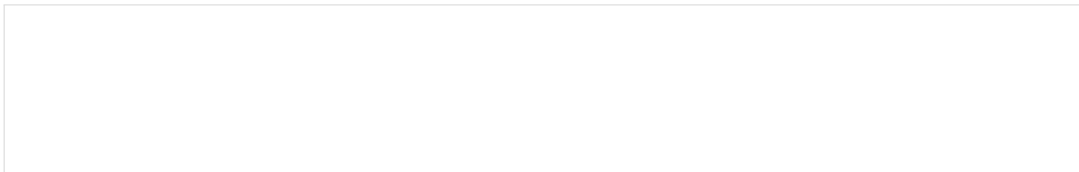
## Properties of OLS estimators: Consistency(2)

- Base on L.L.N(the law of large numbers) and random sample(i.i.d)

- Combining with *Continuous Mapping Theorem*,then we obtain the partitioned multiple OLS estimator  $\hat{\beta}_1$ ,when  $n \rightarrow \infty$

$$plim\hat{\beta}_1 = plim\left(\frac{s_{\tilde{X}_1 Y}}{s_{\tilde{X}_1}^2}\right) = \frac{Cov(\tilde{X}_1, Y)}{Var(\tilde{X}_1)}$$

## Properties of OLS estimators: Consistency(3)



# Properties of OLS estimators: Consistency(4)

- Based on *Assumption 1*:  $E[u_i | X_{1i}, X_{2i} \dots X_{ki}] = 0$
- And  $\tilde{X}_{1i}$  is a function of  $X_{2i} \dots X_{ki}$

- Then

$$Cov(\tilde{X}_1, u_i) = 0$$

- Then we can obtain

$$plim \hat{\beta}_1 = \beta_1$$

- Identical argument works for  $\beta_2, \dots, \beta_k$ , thus

$$plim \hat{\beta}_j = \beta_j \text{ where } j = 1, 2, \dots, k$$



# Recall: The Distribution of Simple OLS Estimators

- Under the least squares assumptions, the Simple OLS estimators  $\hat{\beta}_1$  and  $\hat{\beta}_0$ , are **unbiased and consistent** estimators of  $\beta_1$  and  $\beta_0$ .
- In large samples, the sampling distribution of  $\hat{\beta}_1$  and  $\hat{\beta}_0$  is well approximated by a bivariate normal distribution.
- Specifically, the sampling distribution of  $\hat{\beta}_1$  is

$$\hat{\beta}_1 \xrightarrow{d} N(\beta_1, \sigma_{\hat{\beta}_1}^2)$$

where

$$\sigma_{\hat{\beta}_1}^2 = \frac{\text{Var}[(X_i - \mu_x)u_i]}{n[\text{Var}(X_i)]^2}$$

# The Distribution of Multiple OLS Estimators

- Similarly as in the simple OLS, the multiple OLS estimators are averages of the randomly sampled data, and if the sample size is sufficiently large, the sampling distribution of those averages becomes normal.

$$\hat{\beta}_j = \beta_j + \frac{\left(\sum_{i=1}^n \tilde{X}_{ij} u_i\right)}{\left(\sum_{i=1}^n \tilde{X}_{ij}^2\right)}$$

- Then we have

## Wrap up

- Under the four least squares assumptions, the multiple OLS estimators are **unbiased and consistent** estimators and **approximately normal** in large samples.
- The FWL Theorem provides a simple but important perspective to understand the multiple OLS estimator, which is the **partialling out** or **netting out** is a key concept in the multiple regression analysis.

## Multiple OLS Regression and Causality

# Independent Variable v.s Control Variables

- Generally, we would like to pay more attention to **only one** independent variable (thus we would like to call it **treatment variable**), though there could be many independent variables.
- Because  $\beta_j$  is **partial (marginal) effect** of  $X_j$  on  $Y$ .

$$\beta_j = \frac{\partial Y_i}{\partial X_{j,i}}$$

which means that we are estimate the effect of  $X$  on  $Y$  when “**other things equal**”, thus the concept of **ceteris paribus**.

- Therefore, other variables in the right hand of equation, we call them **control variables**, which we would like to explicitly **hold fixed** when studying the effect of  $X_1$  or  $D$  on  $Y$ .

# Independent Variable v.s Control Variables

- In a multiple regression, OLS is a way to **control observable confounding factors**, which assume the source of selection bias is only from the difference in observed characteristics(Selection-on-Observables)
- If the multiple regression model is

$$Y_i = \beta_0 + \beta_1 X_{1,i} + \beta_2 X_{2,i} + \dots + \beta_k X_{k,i} + u_i, i = 1, \dots, n$$

- Generally, we would like to pay more attention to **only one** independent variable(thus we would like to call it **treatment variable**), though there could be many independent variables.
- Other variables in the right hand of equation, we call them **control variables**, which we would like to explicitly hold fixed when studying the effect of  $X_1$  on  $Y$ .

# OLS Regression, Covariates and RCT

- More specifically, our multiple regression model turns into

$$Y_i = \beta_0 + \beta_1 D_i + \gamma_2 C_{2,i} + \dots + \gamma_k C_{k,i} + u_i, i = 1, \dots, n$$

- We could transform it into as follows

$$Y_i = \alpha + \rho D_i + C_i' \Gamma + u_i$$

where  $\alpha = \beta_0, \rho = \beta_1, \Gamma = (\gamma_2, \dots, \gamma_k), C_i = (C_{2i}, \dots, C_{ki})$

# OLS Regression, Covariates and RCT

- Now write out the conditional expectation of  $Y_i$  for both levels of  $D_i$  conditional on C

- Taking the difference



# OLS Regression, Covariates and RCT

- Again, our estimate of the **treatment effect** ( $\rho$ ) is only going to be as good as our ability to eliminate the **selection bias**, thus

$$E[u_{1i} | \mathbf{D}_i = 1, C] - E[u_{0i} | \mathbf{D}_i = 0, C] \neq 0$$

# OLS Regression, Covariates and RCT

- This is the equivalence of the **CIA** assumption, which is also equivalent to the **1st assumption** of Multiple OLS

- Then we can eliminate the **selection bias**, thus making

$$E[u_{1i} | \mathbf{D}_i = 1, C] = E[u_{0i} | \mathbf{D}_i = 0, C]$$

- Thus

# Wrap up

- OLS regression is valid or can obtain a causal explanation only when least squares assumptions are held.
- The most critical assumption is either

or

# Picking Control Variables

- **Questions:** Are “more controls” always better (or at least never worse)?
- **Answer:** *It depends on.*

- *We will come back soon to discuss the topic in details(in lecture 7 or 8).*