

Lecture 2: Simple OLS Regression Estimation

Introduction to Econometrics, 2025 Spring

Zhaopeng Qu

Nanjing University Business School

March 06 2025



- 1 Review the previous lecture
- 2 OLS Estimation: Simple Regression
- 3 The Least Squares Assumptions
- 4 Properties of the OLS Estimators
- 5 Simple OLS and RCT
- 6 Make Comparison Make Sense
- 7 Appendix

Review the previous lecture

Causal Inference and RCT

- **Causality** is our main goal in the studies of empirical social science.
- The existence of **selection bias** makes social science more difficult than science.
- Although RCTs is a powerful tool for economists, every project or topic can **NOT** be carried on by it.
- This is the reason why modern econometrics exists and develops. The main job of econometrics is using **non-experimental** data to **making convincing causal inference**.

Furious Seven Weapons (七种武器)

- To build a *reasonable counterfactual world* or to find a *proper control group* is the core of econometric methods.
 1. Randomized controlled trial(RCTs)
 2. Regression(回归)
 3. Matching and Propensity Score(匹配与倾向得分)
 4. Instrumental Variable (工具变量)
 5. Regression Discontinuity (断点回归)
 6. Panel Data and Difference in Differences (双差分或倍差法)
 7. Synthetic Control Method (合成控制法)
- The most fundamental of these tools is **regression**. It compares treatment and control subjects with the same observable characteristics **in a generalized manner**.
- It paves the way for the more elaborate tools used in the class that follow.
- **Let's start our exciting journey from it.**

OLS Estimation: Simple Regression

Question: Class Size and Student's Performance

- **Specific Question:**
 - What is the effect on district **test scores** if we would increase district average **class size** by 1 student (or one unit of Student-Teacher's Ratio)
- If we could know the full relationship between two variables which can be summarized by a real value function, $f(\cdot)$
- Unfortunately, the function form is always unknown.

Question: Class Size and Student's Performance

- Two basic methods to describe the function.
 - **non-parametric**: we don't care the specific form of the function, unless we know all the values of two variables, which actually are the *whole distributions* of class size and test scores.
 - **parametric**: we have to suppose the basic form of the function, then to find values of some *unknown parameters* to determine the specific function form.
- Both methods need to use **samples** to inference **populations** in our random and unknown world.

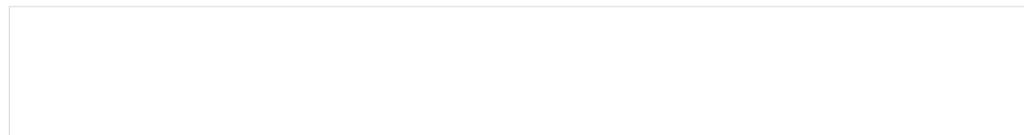
Question: Class Size and Student's Performance

- Suppose we choose *parametric* method, then we just need to know the real value of a **parameter** β_1 to describe the relationship between Class Size and Test Scores
- Next step, we have to suppose specific forms of the function $f(\cdot)$, still two categories: **linear** and **non-linear**
- And we start to use the *simplest* function form: a **linear** equation, which is graphically a **straight line**, to summarize the relationship between two variables.

where β_1 is actually the **the slope** and β_0 is the **intercept** of the straight line.

Class Size and Student's Performance

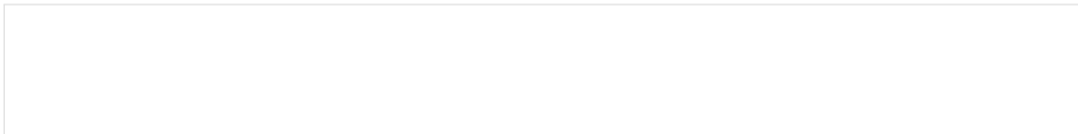
- BUT the average test score in district i does not **only** depend on the average class size
- It also depends on **other factors** such as
 - Student background
 - Quality of the teachers
 - School's facilities
 - Quality of text books
 - Random deviation
- So the equation describing the linear relation between Test score and Class size is **better** written as



where α_i lumps together all other factors that affect average test scores

Terminology for Simple Regression Model

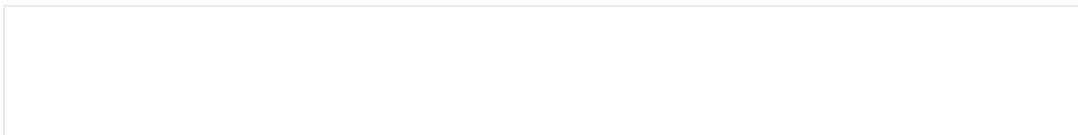
- The linear regression model with one regressor is denoted By



- Where
 - Y_i is the **dependent variable**(Test Score)
 - X_i is the **independent variable** or regressor(Class Size or Student-Teacher Ratio)
 - $\beta_0 + \beta_1 X_i$ is the **population regression line** or the **population regression function**

Population Regression: relationship in average

- The linear regression model with one regressor is denoted by



- Both side to conditional on X , then

$$E[Y_i|X_i] = \beta_0 + \beta_1 X_i + E[u_i|X_i]$$

- Suppose $E[u_i|X_i] = 0$ then

$$E[Y_i|X_i] = \beta_0 + \beta_1 X_i$$

- Population regression function is the relationship that holds between Y and X **on average over the population.**

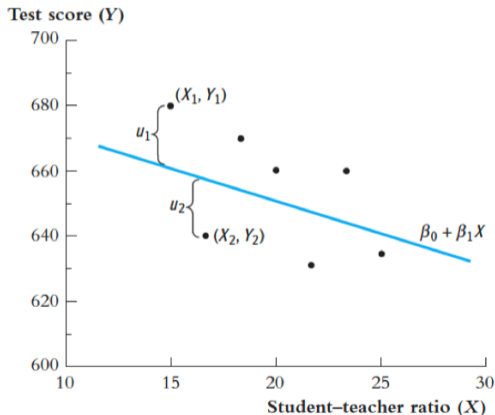
Terminology for Simple Regression Model

- The intercept β_0 and the slope β_1 are the **coefficients** of the **population regression line**, also known as the **parameters** of the population regression line.
- 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 .

Graphics for Simple Regression Model

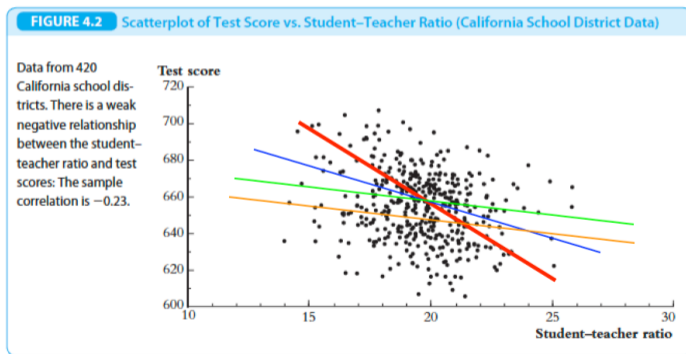
FIGURE 4.1 Scatterplot of Test Score vs. Student-Teacher Ratio (Hypothetical Data)

The scatterplot shows hypothetical observations for seven school districts. The population regression line is $\beta_0 + \beta_1 X$. The vertical distance from the i^{th} point to the population regression line is $Y_i - (\beta_0 + \beta_1 X_i)$, which is the population error term u_i for the i^{th} observation.



How to find the “best” fitting line?

- In general we don't know β_0 and β_1 which are parameters of **population regression function** but have to calculate them using a bunch of data: the **sample**.

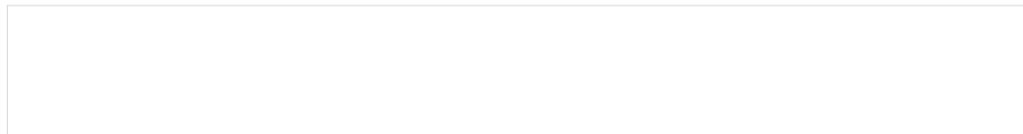


- So how to find the line that fits the data **best**?

The Ordinary Least Squares Estimator (OLS)

The OLS estimator

- Chooses the **best** regression coefficients so that the estimated regression line is **as close as possible** to the observed data, where closeness is measured by **the sum of the squared mistakes** made in predicting Y given X.
- Let b_0 and b_1 be estimators of β_0 and β_1 , thus $b_0 \equiv \hat{\beta}_0, b_1 \equiv \hat{\beta}_1$
- The predicted value of Y_i given X_i using these estimators is $b_0 + b_1 X_i$, or $\hat{\beta}_0 + \hat{\beta}_1 X_i$ formally denotes as \hat{Y}_i , thus



The Ordinary Least Squares Estimator (OLS)

The OLS estimator

- The prediction mistake is **the difference** between Y_i and \hat{Y}_i , which denotes as \hat{u}_i

- 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 β_0 and β_1 .

The Ordinary Least Squares Estimator (OLS)

- OLS minimizes sum of squared prediction mistakes:

- Solve the problem by F.O.C(the first order condition)
 - Step 1 for β_0 :

- Step 2 for β_1 :

Step 1: OLS estimator of β_0

- Recall the sample mean of Y_i is

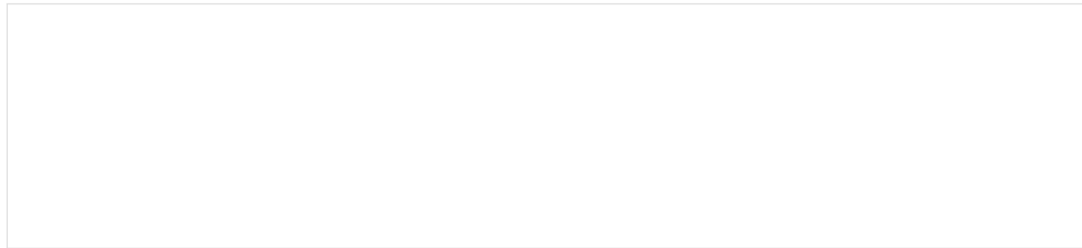
- Optimization

Step 1: OLS estimator of β_0

OLS estimator of β_0 :

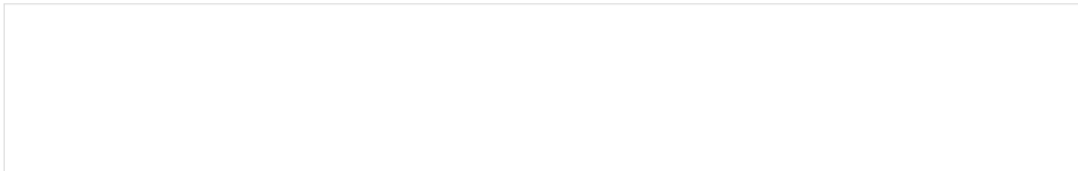
$$b_0 \equiv \hat{\beta}_0 = \bar{Y} - b_1 \bar{X}$$

Step 2: OLS estimator of β_1



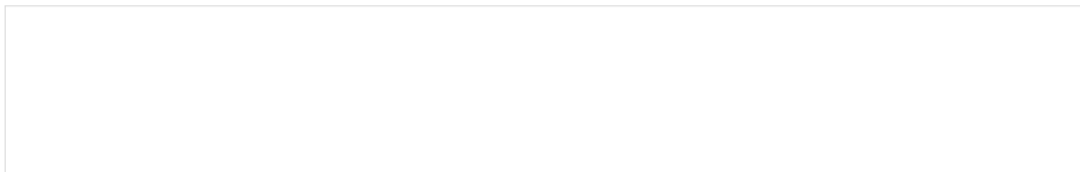
Step 2: OLS estimator of β_1

- Some Algebraic Facts

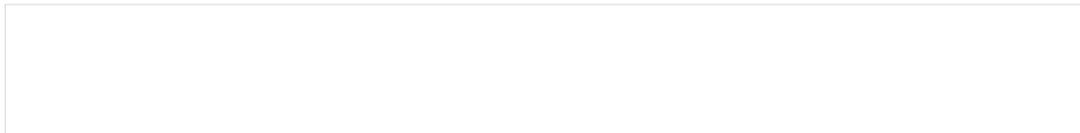


Step 2: OLS estimator of β_1

- Some Algebraic Facts

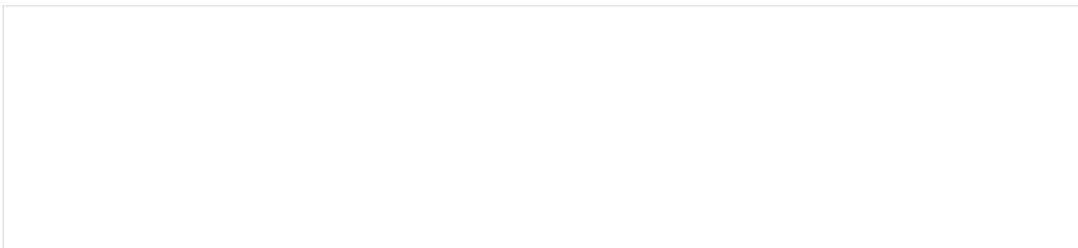


- By a similar reasoning, we could obtain



Step 2: OLS estimator of β_1

Step 2: OLS estimator of β_1

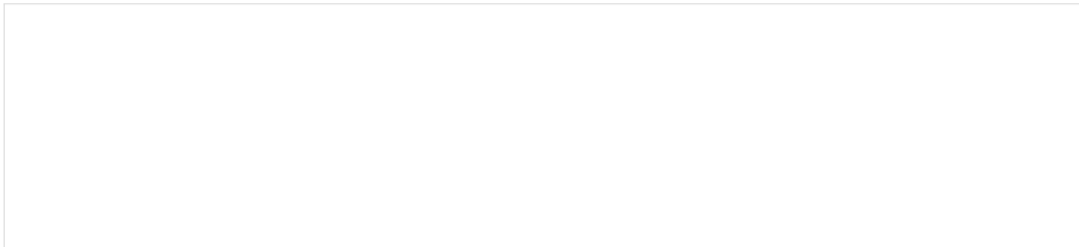


OLS estimator of β_1 :

$$b_1 \equiv \hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})}$$

Some Algebraic of \hat{u}_i

- Recall the F.O.C

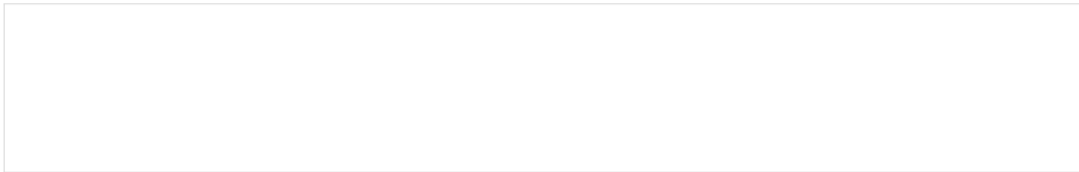


- We obtain two intermediate formulas

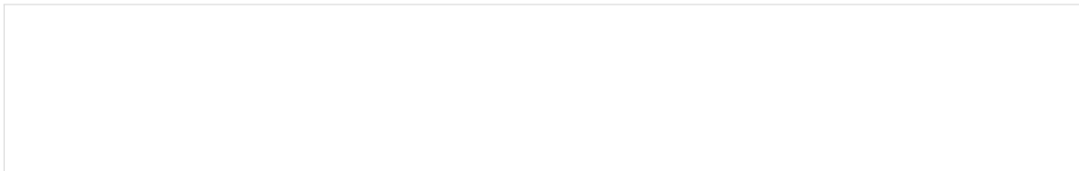


Some Algebraic of \hat{u}_i

- Recall the OLS predicted values \hat{Y}_i and residuals \hat{u}_i are:



- Then we have (*prove them by yourself, Appendix 4.3 in SW, pp184-185*)



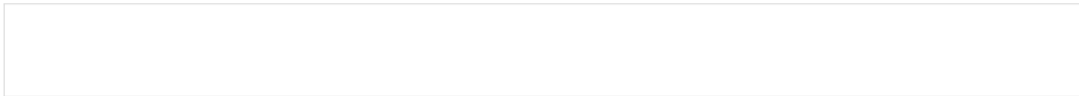
The Estimated Regression Line

- Obtain the values of OLS estimator for a certain data,

- Then the regression line is

The Estimated Regression Line

- Obtain the values of OLS estimator for a certain data,



- Then the regression line is



Measures of Fit: The R^2

- Because the variation of Y can be summarized by a statistic: **Variance**, so the total variation of Y_i , which are also called as the **total sum of squares (TSS)**, is:

-
- Because Y_i can be decomposed into the fitted value plus the residual:
 $Y_i = \hat{Y}_i + \hat{u}_i$, then likewise Y_i , we can obtain

-
- The **explained sum of squares (ESS)**:

Measures of Fit: The R^2

Proof of $TSS = ESS + SSR$

Measures of Fit: The R^2

R^2 or the coefficient of determination

R^2 or the coefficient of determination, is the fraction of the sample variance of Y_i explained/predicted by X_i

- So $0 \leq R^2 \leq 1$, it measures that how much can the variations of Y be explained by the variations of X_i in share.
- **Question:** *If R-squares is bigger, is the regression better?*
- **Answer:** **Not necessarily**, especially when we make **causal inference** in cross-sectional data.

The Least Squares Assumptions

The Linear Regression Model

- In order to investigate the statistical properties of OLS, we need to make some statistical assumptions

The Linear Regression Model

- In order to investigate the statistical properties of OLS, we need to make some statistical assumptions

Linear Regression Model

Two random variables Y_i and X_i , their relationship can satisfy the linear regression equation, thus

- This is not a required assumption. We will extend the model to be nonlinear later on.

Assumption 1: Conditional Mean is Zero

Assumption 1: Zero conditional mean of the errors given X

The error, u_i has expected value of 0 given any value of the independent variable

Assumption 1: Conditional Mean is Zero

Assumption 1: Zero conditional mean of the errors given X

The error, u_i has expected value of 0 given any value of the independent variable

Implications of Assumption 1

With the Iterated Expectation Law, we can obtain an extra implicit assumption about u_i , thus

Assumption 1: Conditional Mean is Zero

- An *weaker* condition that u_i and X_i are uncorrelated:

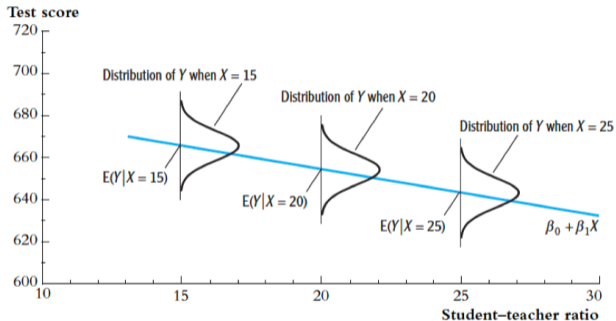
Covariance and Conditional Mean

Although $Cov[u_i, X_i] = 0 \not\Rightarrow E[Y_i|X_i]$, we have

- if u_i and X_i are correlated, then **Assumption 1 is violated**.

Assumption 1: Conditional Mean is Zero

FIGURE 4.4 The Conditional Probability Distributions and the Population Regression Line



The figure shows the conditional probability of test scores for districts with class sizes of 15, 20, and 25 students. The mean of the conditional distribution of test scores, given the student-teacher ratio, $E(Y|X)$, is the population regression line. At a given value of X , Y is distributed around the regression line and the error, $u = Y - (\beta_0 + \beta_1 X)$, has a conditional mean of zero for all values of X .

Assumption 2: Random Sample

Assumption 2: Random Sample

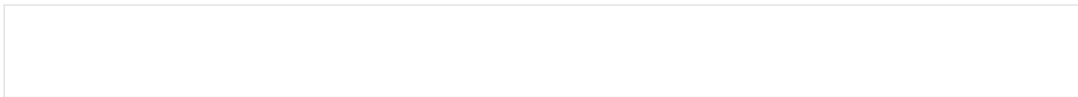
We have a i.i.d random sample of size , $\{(X_i, Y_i), i = 1, \dots, n\}$ from the population regression model above.

Assumption 2: Random Sample

Assumption 2: Random Sample

We have a i.i.d random sample of size , $\{(X_i, Y_i), i = 1, \dots, n\}$ from the population regression model above.

- This is an implication of random sampling. Then we have such as



- And it generally won't hold in other data structures.
 - time-series, cluster samples and spatial data.

Assumption 3: Large outliers are unlikely

Assumption 3: Large outliers are unlikely

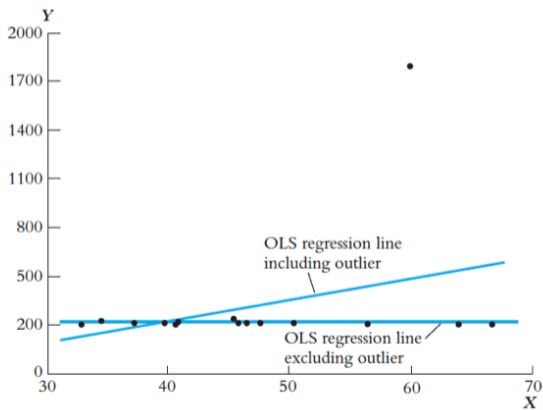
It states that observations with values of X_i , Y_i or both that are far outside the usual range of the data (Outlier) are unlikely. Mathematically, it assumes that X and Y have nonzero finite fourth moments.

- Large outliers can make OLS regression results misleading.
- One source of large outliers is data entry errors, such as a typographical error or incorrectly using different units for different observations.
- Data entry errors aside, the assumption of finite kurtosis is a plausible one in many applications with economic data.

Assumption 3: Large outliers are unlikely

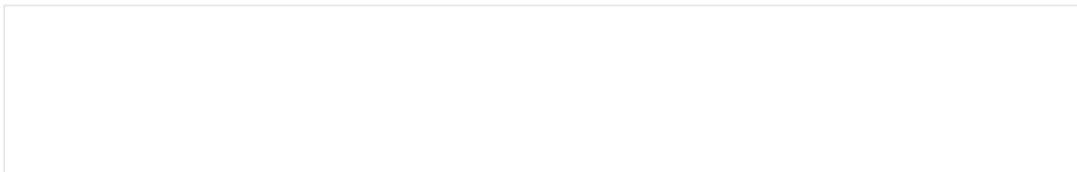
FIGURE 4.5 The Sensitivity of OLS to Large Outliers

This hypothetical data set has one outlier. The OLS regression line estimated with the outlier shows a strong positive relationship between X and Y , but the OLS regression line estimated without the outlier shows no relationship.



Underlying Assumptions of OLS

- The OLS estimator is **unbiased, consistent and has asymptotically normal sampling distribution** if



Underlying assumptions of OLS

- OLS is an **estimator**: it is a machine that we plug data into and we get out estimates.
- It has a **sampling distribution**, with a sampling variance/standard error, etc. like the sample mean, sample difference in means, or the sample variance.
- Let's discuss these characteristics of OLS in the next section.

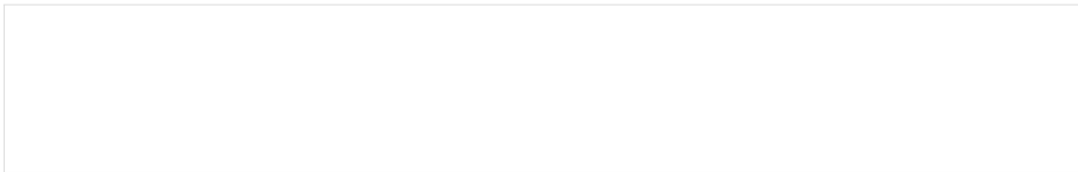
Properties of the OLS Estimators

The OLS estimators

- Question of interest: What is the effect of a change in X_i (Class Size) on Y_i (Test Score)

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

- We derived the OLS estimators of β_0 and β_1 :



Least Squares Assumptions

1. Assumption 1: Conditional Mean is Zero
 2. Assumption 2: Random Sample
 3. Assumption 3: Large outliers are unlikely
- If the 3 least squares assumptions hold the OLS estimators will be
 - **unbiased**
 - **consistent**
 - **normal sampling distribution**

Properties of the OLS estimator: unbiasedness

- Recall:

- take expectation to β_0 :

- Then we have: if β_1 is unbiased, then β_0 is also unbiased.

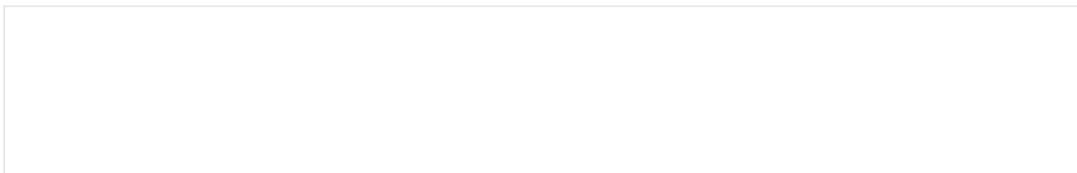
Properties of the OLS estimator: unbiasedness

- Remind we have

- So take expectation to β_1 :

Properties of the OLS estimator: unbiasedness

- Continued



Review: Conditional Expectation Function(CEF)

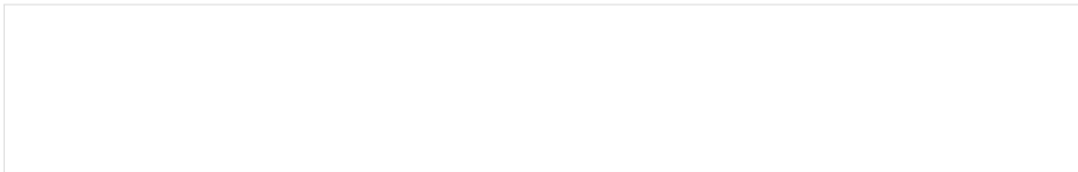
- Expectation(for a continuous r.v.)

- Conditional probability density function

- Conditional Expectation Function: the Expectation of Y conditional on X is

Review: Properties of CEF

Let X, Y, Z are random variables; $a, b \in \mathbb{R}$; $g(\cdot)$ is a real valued function, then we have



Review: the Law of Iterated Expectations(LIE)

the Law of Iterated Expectations

It states that an unconditional expectation can be written as the unconditional average of conditional expectation function :

Proof: the Law of Iterated Expectation(LIE)

- Prove it by a continuous variable way

Proof

Conditional Expectation and Covariance

- Please prove if $E(Y|X) = 0 \Rightarrow Cov(X, Y) = 0$ `\begin{block}{Proof}`

`\end{block}`

Properties of the OLS estimator: unbiasedness

- Because $\sum (X_i - \bar{X})(u_i - \bar{u}) = \sum (X_i - \bar{X})u_i$, so

Properties of the OLS estimator: unbiasedness

- Because $\sum(X_i - \bar{X})(u_i - \bar{u}) = \sum(X_i - \bar{X})u_i$, so

- Then we can obtain

$$E[\hat{\beta}_1] = \beta_1 \text{ if } E[u_i|X_i] = 0$$

- Both β_0 and β_1 are unbiased on the condition of **Assumption 1**.

Properties of the OLS estimator: Consistency

- **Notation:** $\hat{\beta}_1 \xrightarrow{p} \beta_1$ or $plim\hat{\beta}_1 = \beta_1$, so

- Then we could obtain

where s_{xy} and s_x^2 are sample covariance and sample variance.

Recall: Sample Variance and Sample Covariance

Math Review: Continuous Mapping Theorem

- **Continuous Mapping Theorem:** For every continuous function $g(t)$ and random variable X :



- **Example:**

$$plim(X + Y) = plim(X) + plim(Y)$$

$$plim\left(\frac{X}{Y}\right) = \frac{plim(X)}{plim(Y)} \text{ if } plim(Y) \neq 0$$

Properties of the OLS estimator: Consistency

- 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 OLS estimator $\hat{\beta}_1$,when $n \rightarrow \infty$

Properties of the OLS estimator: Consistency

- Then we could obtain

$$plim \hat{\beta}_1 = \beta_1 \text{ if } E[u_i | X_i] = 0$$

Wrap Up: Unbiasedness vs Consistency

- **Unbiasedness & Consistency** both rely on $E[u_i|X_i] = 0$
- **Unbiasedness** implies that $E[\hat{\beta}_1] = \beta_1$ for a certain sample size n . (“small sample”)
- **Consistency** implies that the distribution of $\hat{\beta}_1$ becomes more and more **tightly distributed** around β_1 if the sample size n becomes larger and larger. (“large sample”)
- Additionally, you could prove that $\hat{\beta}_0$ is likewise **Unbiased and Consistent** on the condition of **Assumption 1**.

Sampling Distribution of $\hat{\beta}_0$ and $\hat{\beta}_1$: Recall of \bar{Y}

- Firstly, Let's recall: Sampling Distribution of \bar{Y}
- Because Y_1, \dots, Y_n are i.i.d. and μ_Y is the mean of the population, then for L.L.N, we have

$$E(\bar{Y}) = \mu_Y$$

- Based on the Central Limit theorem (C.L.T) and the σ_Y^2 is the variance of the population, the sample distribution in a large sample can *approximates to a normal distribution*, thus

$$\bar{Y} \sim N\left(\mu_Y, \frac{\sigma_Y^2}{n}\right)$$

- Therefore, the OLS estimators $\hat{\beta}_0$ and $\hat{\beta}_1$ could have similar sample distributions *when three least squares assumptions hold.*

Sampling Distribution of $\hat{\beta}_0$ and $\hat{\beta}_1$: Expectation

- Likewise as \bar{Y} , the sample distribution of β_1 or β_0 in a large sample can also *approximates to a normal distribution* based on the **Central Limit theorem(C.L.T)**

- The **expectation** of the OLS estimators is by the **unbiasedness** of the OLS estimators. It implies that

Sampling Distribution of $\hat{\beta}_1$ and $\hat{\beta}_0$: Variance

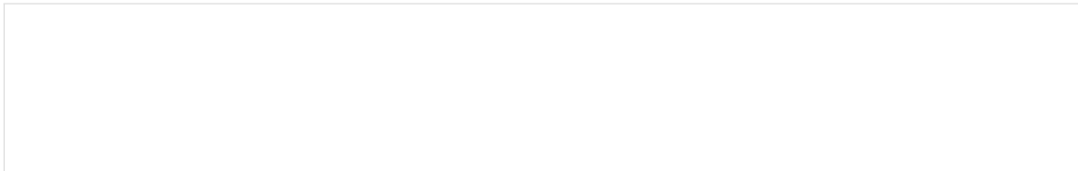
- Likewise as \bar{Y} , the sample distribution of β_1 or β_0 in a large sample can also *approximates to a normal distribution* based on the **Central Limit theorem(C.L.T)**

- The **variance** of the OLS estimators can be shown as follows:

Where $H_i = 1 - \left[\frac{\mu_x}{E[X_i]} \right] X_i$

Sampling Distribution $\hat{\beta}_1$ in large-sample

- We have shown that

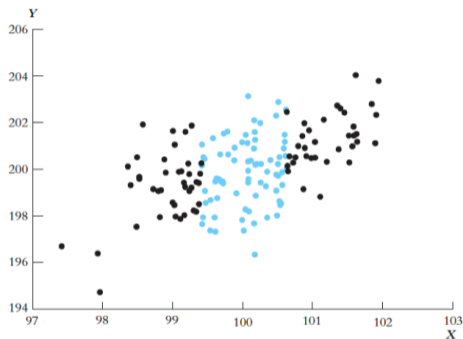


- **An intuition:**The **variation** of X_i is very important.
 - Because if $Var(X_i)$ is *small*, it is difficult to obtain an accurate estimate of the effect of X on Y which implies that $Var(\hat{\beta}_1)$ is *large*.

Variation of X

FIGURE 4.6 The Variance of $\hat{\beta}_1$ and the Variance of X

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



- When more **variation** in X_i , then there is more information in the data that you can use to fit the regression line.

In a Summary

Under 3 least squares assumptions, the OLS estimators will be

- **unbiased**
- **consistent**
- **normal sampling distribution**
- *more variation in X , more accurate estimation*

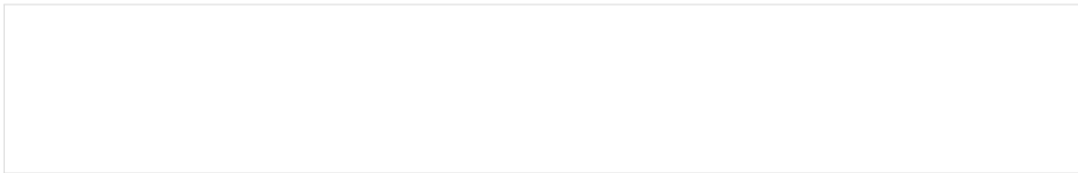
Simple OLS and RCT

OLS Regression and RCT

- We learned RCT is the “**golden standard**” for causal inference. Because it can naturally eliminate **selection bias**.
- So far, we did not discuss the relationship between RCT and OLS regression, which means that we can not be sure that the result from an OLS regression can be explained as “causal”.
- Instead of using a continuous regressor X , the regression where D_i is a binary variable, a so-called **dummy variable**, will help us to unveil the relationship between RCT and OLS regression.

Regression when X is a Binary Variable

- For example, we may define D_i as follows:



- The regression can be written as

$$Y_i = \beta_0 + \beta_1 D_i + u_i \quad (4.1)$$

Regression when X is a Binary Variable

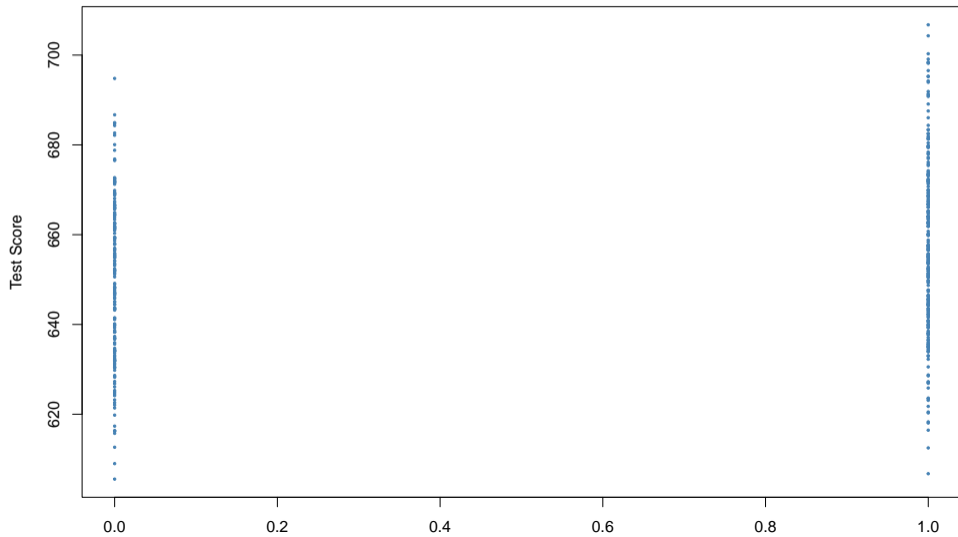
- More precisely, the regression model now is



- With D as the regressor, it is not useful to think of β
 - Since $D_i \in \{0, 1\}$, i.e., we only observe two discrete
-
- There is no continuous line depicting the conditional expectation function $E(\text{TestScore}_i | D_i)$ since this function is solely defined for x -positions 0 and 1.

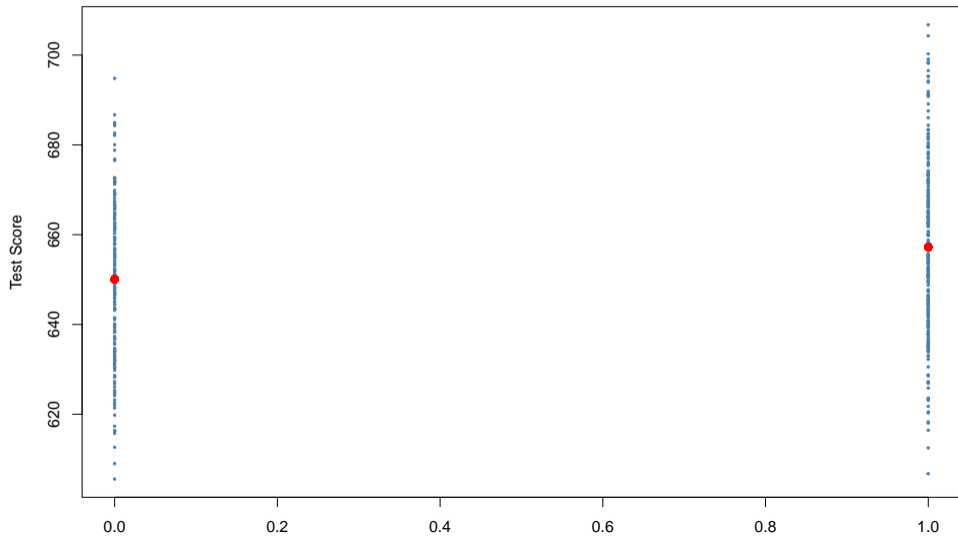
Class Size and STR

Dummy Regression



Class Size and STR

Dummy Regression

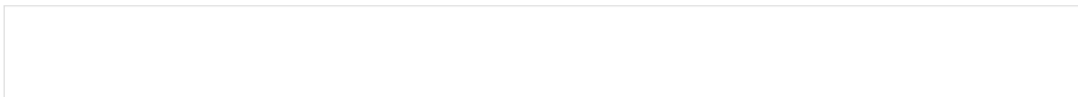


Regression when X is a Binary Variable

- Therefore, the interpretation of the coefficients in this regression model is as follows:
 - $E(Y_i|D_i = 0) = \beta_0$, so β_0 is the expected test score in districts where $D_i = 0$ where STR is below 20.
 - $E(Y_i|D_i = 1) = \beta_0 + \beta_1$ where STR is above 20
- Thus, β_1 is the difference in group specific expectations, i.e., the difference in expected test score between districts with $STR < 20$ and those with $STR \geq 20$,

Causality and OLS

- Let us recall, the individual treatment effect



- The ATE is the average of the ICE and ATT is the average of the ICE for the treated group.



- Either way, the treatment effect is a constant, i.e., it does not depend on the individual.
- Our OLS regression function is to estimate a constant treatment effect ρ , thus

Causality and OLS

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

- Take the difference

Causality and OLS

- Again, our estimate of the **treatment effect** (ρ) is only going to be as good as our ability to shut down the **selection bias**.
- *Selection bias in regression model:* $E[\eta_i | \mathbf{D}_i = 1] - E[\eta_i | \mathbf{D}_i = 0]$
- There is something in our disturbance η_i that is affecting Y_i and is also correlated with D_i .

Simple OLS Regression v.s. RCT

- In a simple regression model, OLS estimators are just a generalizing continuous version of RCT when least squares assumptions are hold.
- Ideally, regression is a way to control observable confounding factors, which assume the source of selection bias is only from the difference in observed characteristics.

Simple OLS Regression v.s. RCT

- But in contrast to RCT, in observational studies, researchers cannot control the assignment of treatment into a treatment group versus a control group, which means that the two groups are **incomparable**.
- To make two groups comparable, we need to keep treatment and control group “**other thing equal**” in observed characteristics and unobserved characteristics.
- OLS regression is valid only when least squares assumptions are hold.
- In most cases, it is not easy to obtain. We have to know how to make a convincing causal inference when these assumptions are not hold.

Make Comparison Make Sense

Case: Smoke and Mortality

- Criticisms from **Ronald A. Fisher**
 - No experimental evidence to incriminate smoking as a cause of lung cancer or other serious disease.
 - Correlation between smoking and mortality may be spurious due to **biased selection** of subjects.

- **Confounder, Z**, creates backdoor path between smoking and mortality

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

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

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.
- Maybe cigar smokers higher observed death rates is because **they're older on average.**

Case: Smoke and Mortality(Cochran 1968)

- The problem is that the age are *not balanced*, thus their mean values differ for treatment and control group.
- 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:**

Formalization: Covariates

Definition: Covariates

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

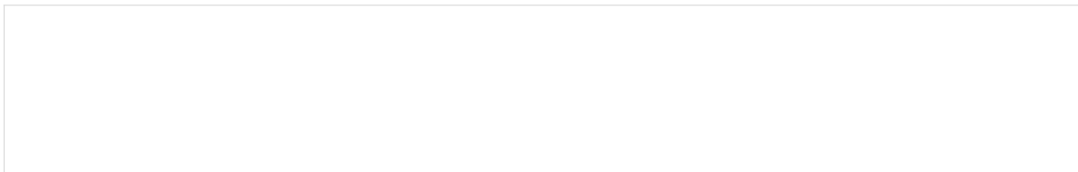
- Covariates are often time invariant (e.g., sex, race), but time invariance is not a necessary condition.

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

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

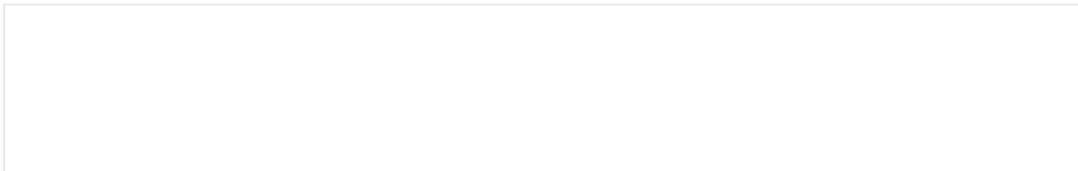
$$(Y_{1i}, Y_{0i}) \perp\!\!\!\perp D | X$$

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

$$E[Y_{1i}|D = 1] - E[Y_{0i}|D = 0] \neq E[Y_{1i}|D = 1] - E[Y_{0i}|D = 1]$$

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 $3^{10} = 59049$

Making Comparison Make Sense

- *Selection on Observables*
 - Regression
 - Matching
- *Selection on Unobservables*
 - IV, RD, DID, FE and SCM.
- The most fundamental tool among them is **regression**, which compares treatment and control subjects who have the same **observable** characteristics **in a generalized manner**.

Extending Reading

- Ho, Chong and Xia(2017),“Yellow taxis have fewer accidents than blue taxis because yellow is more visible than blue”,PNAS,Vol.114(12),pp3074-3078.
- Cochran, W. G. (1968). The effectiveness of adjustment by subclassification in removing bias in observational studies. Biometrics, 24(2), pp295–313.

Appendix

Sampling Distribution of $\hat{\beta}_1$

- $\hat{\beta}_1$ in terms of regression and errors in following equation

$$\hat{\beta}_1 = \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})}$$

Sampling Distribution of $\hat{\beta}_1$

- $\hat{\beta}_1$ in terms of regression and errors in following equation

$$\begin{aligned}\hat{\beta}_1 &= \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})} \\ &= \beta_1 + \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(u_i - \bar{u})}{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})}\end{aligned}$$

Sampling Distribution of $\hat{\beta}_1$: the numerator

- The numerator: $\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(u_i - \bar{u})$
- Because \bar{X} is consistent, thus $\bar{X} \xrightarrow{p} \mu_x$, then combine with Continuous Mapping Theorem

$$\sum_{i=1}^n (X_i - \bar{X})(u_i - \bar{u}) = \sum_{i=1}^n (X_i - \bar{X})u_i$$

Sampling Distribution of $\hat{\beta}_1$: the numerator

- The numerator: $\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(u_i - \bar{u})$
- Because \bar{X} is consistent, thus $\bar{X} \xrightarrow{p} \mu_x$, then combine with Continuous Mapping Theorem

$$\begin{aligned} \sum_{i=1}^n (X_i - \bar{X})(u_i - \bar{u}) &= \sum_{i=1}^n (X_i - \bar{X})u_i \\ \implies \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(u_i - \bar{u}) &\xrightarrow{p} \frac{1}{n} \sum_{i=1}^n (X_i - \mu_x)u_i \end{aligned}$$

Sampling Distribution of $\hat{\beta}_1$: the numerator

- Let $v_i = (X_i - \mu_x)u_i$
 - Based on **Assumption 1**, then $E(v_i) = 0$
 - Based on **Assumption 2**, $\sigma_v^2 = Var[(X_i - \mu_x)u_i]$

- Then

$$\frac{1}{n} \sum_{i=1}^n (X_i - \mu_x)u_i = \frac{1}{n} \sum_{i=1}^n v_i = \bar{v}$$

Sampling Distribution of $\hat{\beta}_1$: the numerator

- Recall: \bar{Y} to Y_i and based on C.L.T,

$$\frac{\bar{Y} - 0}{\sigma_{\bar{Y}}} \xrightarrow{d} N(0,1) \text{ or } \bar{Y} \xrightarrow{d} N\left(0, \frac{\sigma_Y^2}{n}\right)$$

- The \bar{v} is the *sample mean* of v_i , based on C.L.T,

$$\frac{\bar{v} - 0}{\sigma_{\bar{v}}} \xrightarrow{d} N(0,1) \text{ or } \bar{v} \xrightarrow{d} N\left(0, \frac{\sigma_v^2}{n}\right)$$

Sampling Distribution of $\hat{\beta}_1$: the denominator

- Recall the sample variance of X_i is $s_{X_i}^2$

$$s_{X_i}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

- Then the denominator, is a variation of **sample variance** of X (except dividing by n rather than $n-1$, which is *inconsequential if n is large*)

$$\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})$$

- Based on discussion of the *sample variance* is a **consistent** estimator of the *population variance*, thus

$$s_{X_i}^2 \xrightarrow{p} \text{Var}[X_i] = \sigma_{X_i}^2$$

Sampling Distribution of $\hat{\beta}_1$

- $\hat{\beta}_1$ in terms of regression and errors

$$\hat{\beta}_1 = \beta_1 + \frac{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(u_i - \bar{u})}{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})}$$

- the numerator is \bar{v} and $\bar{v} \xrightarrow{d} N(0, \frac{\sigma_v^2}{n})$
- the denominator is

$$\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X}) \xrightarrow{p} \text{Var}[X_i] = \sigma_{X_i}^2$$

- Combining these two results, we have that, *in large samples*

$$\hat{\beta}_1 - \beta_1 \xrightarrow{p} \frac{\bar{v}}{\text{Var}[X_i]}$$

Slutsky's Theorem

- It combines consistency and convergence in distribution.

Slutsky's Theorem

Suppose that $a_n \xrightarrow{p} a$, where a is a constant, and $S_n \xrightarrow{d} S$. Then

$$a_n + S_n \xrightarrow{d} a + S$$

$$a_n S_n \xrightarrow{d} aS$$

$$\frac{S_n}{a_n} \xrightarrow{d} \frac{S}{a} \text{ if } a \neq 0$$

Sampling Distribution of $\hat{\beta}_1$

- Based on \bar{v} follow a normal distribution, in large samples, thus

$$\bar{v} \xrightarrow{d} N\left(0, \frac{\sigma_v^2}{n}\right)$$

Sampling Distribution of $\hat{\beta}_1$

- Based on \bar{v} follow a normal distribution, in large samples, thus

$$\bar{v} \xrightarrow{d} N\left(0, \frac{\sigma_v^2}{n}\right)$$
$$\Rightarrow \frac{\bar{v}}{\text{Var}[X_i]} \xrightarrow{d} N\left(0, \frac{\sigma_v^2}{n[\text{Var}(X_i)]^2}\right)$$

Sampling Distribution of $\hat{\beta}_1$

- Based on \bar{v} follow a normal distribution, in large samples, thus

$$\begin{aligned}\bar{v} &\xrightarrow{d} N\left(0, \frac{\sigma_v^2}{n}\right) \\ \Rightarrow \frac{\bar{v}}{\text{Var}[X_i]} &\xrightarrow{d} N\left(0, \frac{\sigma_v^2}{n[\text{Var}(X_i)]^2}\right) \\ \Rightarrow \hat{\beta}_1 - \beta_1 &\xrightarrow{d} N\left(0, \frac{\sigma_v^2}{n[\text{Var}(X_i)]^2}\right)\end{aligned}$$

Sampling Distribution of $\hat{\beta}_1$

- Based on \bar{v} follow a normal distribution, in large samples, thus

$$\begin{aligned}\bar{v} &\xrightarrow{d} N\left(0, \frac{\sigma_v^2}{n}\right) \\ \Rightarrow \frac{\bar{v}}{\text{Var}[X_i]} &\xrightarrow{d} N\left(0, \frac{\sigma_v^2}{n[\text{Var}(X_i)]^2}\right) \\ \Rightarrow \hat{\beta}_1 - \beta_1 &\xrightarrow{d} N\left(0, \frac{\sigma_v^2}{n[\text{Var}(X_i)]^2}\right)\end{aligned}$$

- Then the sampling distribution of $\hat{\beta}_1$ is

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

where

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