Lecture 7: Assessing Regression Studies(I)

Introduction to Econometrics, Spring 2025

Zhaopeng Qu

Business School, Nanjing University

April 10 2025



- **1** Review of Previous Lectures
- 2 Accessing Regression Studies: Introduction
- 3 Omitted Variable Bias(OVB) and Control Variables
- **4** DAGs and Control Variables
- 5 Some practical tips about Control Variables

Review of Previous Lectures

Simple and Multiple OLS Regressions

The simple OLS regression model is

$$Y_i = \beta_0 + \beta_1 X_i + u_i, i = 1, ..., n$$

• The OLS estimator for β_1 is

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

• The OLS regression model is

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

• The OLS estimator for β_j

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

Simple/Multiple OLS Regressions: Assumptions

- If the four least squares assumptions in the multiple regression model hold:
 - Assumption 1: The conditional distribution of u_i given $X_{1i}, ..., X_{ki}$ has mean zero,thus

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

- Assumption 2: $(Y_i, X_{1i}, ..., X_{ki})$ are i.i.d.
- Assumption 3: Large outliers are unlikely.
- Assumption 4: No perfect multicollinearity.(only for Multiple OLS regression)
- Then
 - The OLS estimators $\hat{\beta_0}, \hat{\beta_1}...\hat{\beta_k}$ are unbiased.
 - The OLS estimators $\hat{\beta}_0, \hat{\beta}_1...\hat{\beta}_k$ are consistent.
 - The OLS estimators $\hat{eta_0}, \hat{eta_1}...\hat{eta_k}$ are normally distributed in large samples.

- 1. Nonlinear in Xs
 - Polynomials,Logarithms and Interactions
 - The multiple regression framework can be extended to handle regression functions that are nonlinear in one or more X.
 - the difference from a standard multiple OLS regression is *how to explain estimating coefficients*.

- **2**. Nonlinear in β or Nonlinear in Y
 - Discrete Dependent Variables or Limited Dependent Variables.
 - Linear function in Xs is not a good prediciton function or Y, instead a function which parameters enter nonlinearly, such as logisitic or standard normal.
 - The parameters can not obtained by OLS estimation any more but Maximum Likelihood Estimation
- MLE need more assumptions than OLS, thus the distribution of u_i is known.
- In the large sample, the MLE estimator is also consistent and asymptotically normal

Accessing Regression Studies: Introduction

- The concepts of **internal and external validity** provide a general framework for assessing whether a empirical studies answers a specific question of interest rightly and usefully.
 - **Internal validity**: the statistical inferences about **causal effects are valid** for the population and setting **being studied**.
 - **External validity**: the statistical inferences can be generalized from the population and setting studied to other populations and settings.
- Internal and external validity distinguish between
 - the population and setting being studied
 - the population and setting to which the results are generalized.

Differences between studied and interest

The population and setting studied

- The population studied is the population of entities-people, companies, school districts, and so forth-from which the sample is drawn.
- The setting studied refers to as the institutional, legal, social, and economic environment in which the population studied fits in and the sample is drawn.
- The population and setting of interest
 - The population and setting of interest is the population and setting of entities to which the causal inferences from the study are to be applied(generalized).
- Example: Class size and test score
 - the population studies: elementary schools in CA
 - the population of interest: middle schools in CA
 - different populations and settings: elementary schools in MA

- Internal validity is the top priority in causal inference studies.
- External validity is the secondary focus, but only if internal validity is secured.
- In result, we care about the internal validity over 100 times than the external validity in most studies.
- The following content will focus on the internal validity of regression studies.

Threats to Internal Validity in OLS Regressions

• Suppose we are interested in the causal effect of X_1 on Y and we estimate the following multiple regression model

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

- Internal validity has three components:
- 1. The estimators of β_1 should be **unbiased and consistent**. This is the most critical aspect.
- 2. Hypothesis tests and confidence intervals should have the **desired significance level** (at least 5% significant).
- 3. The value of β_1 should be **large enough** to be meaningful or economically significant.

Threats to Internal Validity

- Threats to internal validity:
 - 1. Omitted Variables
 - 2. Misspecification
 - 3. Measurement Error
 - 4. Simultaneous Causality
 - 5. Missing Data and Sample Selection
 - 6. Heteroskedasticity and/or Correlated error terms
 - 7. Significant coefficients or marginal effects
- In a narrow sense,
 - **Internal Invalidity = endogeneity in the estimation** which is caused by the above 1-5 threats.
- In a broad sense,
 - Internal Invalidity = 1-5 threats + 6-7 threats

Omitted Variable Bias(OVB) and Control Variables

OVB Review

- Suppose we want to estimate the causal effect of X_i on Y_i , which represent STR and Test Score, respectively.
- Besides, W_i is the share of English learners which is **omitted** in the regression.
- Two models are as follows:
 - True model:

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

where $E(u_i|X_i, W_i) = 0$

• Observed model:

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

where $v_i = \gamma W_i + u_i$

- Then we have the OLS estimator of β_1 as follows

$$plim\hat{\beta}_1 = \beta_1 + \gamma \frac{Cov(X_i, W_i)}{VarX_i}$$

- An omitted variable W_i leads to an inconsistent OLS estimate of the causal effect of X_i if **both**
 - W_i is related to X, thus $Cov(X_i, W_i) \neq 0$
 - W_i has some effect on Y_i , thus $\gamma \neq 0$
- The OLS estimator does not provide a unbiased and consistent estimate of the causal effect of *X_i*, in other words, the OLS regression is not **internally valid**.

Wrap Up

- OVB bias is the most possible bias when we run OLS regression using nonexperimental data.
- OVB bias 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: Control
 - Putting omitted variables into the right side of the regression as **control variables**, which are independent variables that are not the variable of interest.
- A critical question, often overlooked by many students and even experienced researchers, is:
 - Should we control as many variables as possible to avoid omitted variable bias (OVB)?
 - What kinds of variables can serve as control variables?
- Let us dig deeper into the control variables and OVB in the following sections.

OLS Regression Estimators in Partitioned Regression

Recall: OLS Regression Estimators in Multiple OLS

• OLS estimator in Multiple OLS

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

• The OLS estimator of β_j is

$$\hat{\beta}_j = \frac{\sum_{i=1}^n \tilde{X}_{ij}, Y_i}{\sum_{i=1}^n (\tilde{X}_{ij})^2}$$

- The asymptotic OLS estimator of β_j

$$plim\hat{\beta}_j = \frac{Cov(\tilde{X}_{ij}, Y_i)}{Var(\tilde{X}_{ij})}$$

• Where \tilde{X}_{ij} is the fitted OLS residual of regressing X_{ij} on other regressors, thus

$$X_{ij} = \hat{\gamma}_0 + \hat{\gamma}_1 X_{i1} + \hat{\gamma}_2 X_{i2} + \dots + \hat{\gamma}_{j-1} X_{i,j-1} + \hat{\gamma}_{j+1} X_{i,j+1} + \dots + \hat{\gamma}_k X_{ik} + \tilde{X}_{ij}$$

$$\hat{\boldsymbol{\beta}}_j = \frac{\sum_{i=1}^n \tilde{X}_{ij} Y_i}{\sum_{i=1}^n \tilde{X}_{ij}^2}$$

$$\hat{\boldsymbol{\beta}}_{j} = \frac{\sum_{i=1}^{n} \tilde{X}_{ij} Y_{i}}{\sum_{i=1}^{n} \tilde{X}_{ij}^{2}} = \frac{\sum_{i=1}^{n} \tilde{X}_{ij} (\boldsymbol{\beta}_{0} + \boldsymbol{\beta}_{1} X_{i1} + \boldsymbol{\beta}_{2} X_{i2} + \dots + \boldsymbol{\beta}_{k} X_{ik} + u_{i})}{\sum_{i=1}^{n} \tilde{X}_{ij}^{2}}$$

$$\hat{\boldsymbol{\beta}}_{j} = \frac{\sum_{i=1}^{n} \tilde{X}_{ij} Y_{i}}{\sum_{i=1}^{n} \tilde{X}_{ij}^{2}} = \frac{\sum_{i=1}^{n} \tilde{X}_{ij} (\boldsymbol{\beta}_{0} + \boldsymbol{\beta}_{1} X_{i1} + \boldsymbol{\beta}_{2} X_{i2} + \dots + \boldsymbol{\beta}_{k} X_{ik} + u_{i})}{\sum_{i=1}^{n} \tilde{X}_{ij}^{2}}$$
$$= \frac{\boldsymbol{\beta}_{j} \sum_{i=1}^{n} (\tilde{X}_{ij} X_{ij} + \tilde{X}_{ij} u_{i})}{\sum_{i=1}^{n} \tilde{X}_{ij}^{2}}$$

$$\begin{split} \hat{\boldsymbol{\beta}}_{j} &= \frac{\sum_{i=1}^{n} \tilde{X}_{ij} Y_{i}}{\sum_{i=1}^{n} \tilde{X}_{ij}^{2}} = \frac{\sum_{i=1}^{n} \tilde{X}_{ij} (\boldsymbol{\beta}_{0} + \boldsymbol{\beta}_{1} X_{i1} + \boldsymbol{\beta}_{2} X_{i2} + \ldots + \boldsymbol{\beta}_{k} X_{ik} + u_{i})}{\sum_{i=1}^{n} \tilde{X}_{ij}^{2}} \\ &= \frac{\boldsymbol{\beta}_{j} \sum_{i=1}^{n} (\tilde{X}_{ij} X_{ij} + \tilde{X}_{ij} u_{i})}{\sum_{i=1}^{n} \tilde{X}_{ij}^{2}} \\ &= \boldsymbol{\beta}_{j} + \frac{\sum_{i=1}^{n} \tilde{X}_{ij} u_{i}}{\sum_{i=1}^{n} \tilde{X}_{ij}^{2}} \end{split}$$

• Based on the multiple OLS estimators in ,we have

$$\begin{split} \hat{\boldsymbol{\beta}}_{j} &= \frac{\sum_{i=1}^{n} \tilde{X}_{ij} Y_{i}}{\sum_{i=1}^{n} \tilde{X}_{ij}^{2}} = \frac{\sum_{i=1}^{n} \tilde{X}_{ij} (\boldsymbol{\beta}_{0} + \boldsymbol{\beta}_{1} X_{i1} + \boldsymbol{\beta}_{2} X_{i2} + \ldots + \boldsymbol{\beta}_{k} X_{ik} + u_{i})}{\sum_{i=1}^{n} \tilde{X}_{ij}^{2}} \\ &= \frac{\boldsymbol{\beta}_{j} \sum_{i=1}^{n} (\tilde{X}_{ij} X_{ij} + \tilde{X}_{ij} u_{i})}{\sum_{i=1}^{n} \tilde{X}_{ij}^{2}} \\ &= \boldsymbol{\beta}_{j} + \frac{\sum_{i=1}^{n} \tilde{X}_{ij} u_{i}}{\sum_{i=1}^{n} \tilde{X}_{ij}^{2}} \end{split}$$

• For simplicity, under the 5th assumption of multiple OLS regression: homoskedasticity, thus

$$Var(u_i|X_{1i}, X_{2i}..., X_{ki}) = Var(u_i) = \sigma_u^2$$

• Then we have the variance of $\hat{\beta}_j$ as follows

 $Var(\hat{\beta}_j)$

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

$$\begin{split} Var(\hat{\beta}_j) &= Var\left(\boldsymbol{\beta}_j + \frac{\left(\sum_{i=1}^n \tilde{X}_{ij} u_i\right)}{\left(\sum_{i=1}^n \tilde{X}_{ij}^2\right)}\right) \\ &= \frac{\left(\sum_{i=1}^n \tilde{X}_{ij}^2 Var(u_i)\right)}{\left(\sum_{i=1}^n \tilde{X}_{i1}^2\right)^2} \text{ for i.i.d sample} \end{split}$$

$$\begin{split} Var(\hat{\beta}_j) &= Var\left(\beta_j + \frac{\left(\sum_{i=1}^n \tilde{X}_{ij}u_i\right)}{\left(\sum_{i=1}^n \tilde{X}_{ij}^2\right)}\right) \\ &= \frac{\left(\sum_{i=1}^n \tilde{X}_{ij}^2 Var(u_i)\right)}{\left(\sum_{i=1}^n \tilde{X}_{i1}^2\right)^2} \text{ for i.i.d sample} \\ &= \frac{\left(\sum_{i=1}^n \tilde{X}_{ij}^2 \sigma_u^2\right)}{\left(\sum_{i=1}^n \tilde{X}_{ij}^2\right)^2} \text{ for homoskedasticity} \end{split}$$

$$\begin{aligned} Var(\hat{\beta}_j) &= Var\left(\beta_j + \frac{\left(\sum_{i=1}^n \tilde{X}_{ij} u_i\right)}{\left(\sum_{i=1}^n \tilde{X}_{ij}^2\right)}\right) \\ &= \frac{\left(\sum_{i=1}^n \tilde{X}_{ij}^2 Var(u_i)\right)}{\left(\sum_{i=1}^n \tilde{X}_{i1}^2\right)^2} \text{ for i.i.d sample} \\ &= \frac{\left(\sum_{i=1}^n \tilde{X}_{ij}^2 \sigma_u^2\right)}{\left(\sum_{i=1}^n \tilde{X}_{ij}^2\right)^2} \text{ for homoskedasticity} \\ &= \frac{\sigma_u^2}{\left(\sum_{i=1}^n \tilde{X}_{ij}^2\right)} \end{aligned}$$

• The \tilde{X}_{ij} is the **residual** obtained from the regression of X_j on all other Xs

$$X_{ij} = \hat{X}_{ij} + \tilde{X}_{ij} = \hat{\delta}_0 + \hat{\delta}_1 X_{i1} + \dots + \hat{\delta}_{j-1} X_{i,j-1} + \hat{\delta}_{j+1} X_{i,j+1} + \dots \hat{\delta}_k X_{ik} + \tilde{X}_{ij}$$

• Then the R-Squared of this partitioned regression is

$$R_j^2 = 1 - \frac{SSR_j}{TSS_j}$$

• The \tilde{X}_{ij} is the **residual** obtained from the regression of X_j on all other Xs

$$X_{ij} = \hat{X}_{ij} + \tilde{X}_{ij} = \hat{\delta}_0 + \hat{\delta}_1 X_{i1} + \dots + \hat{\delta}_{j-1} X_{i,j-1} + \hat{\delta}_{j+1} X_{i,j+1} + \dots \hat{\delta}_k X_{ik} + \tilde{X}_{ij}$$

• Then the **R-Squared** of this partitioned regression is

$$R_j^2 = 1 - \frac{SSR_j}{TSS_j}$$
$$\Rightarrow SSR_j = (1 - R_j^2) \times TSS_j$$

• The \tilde{X}_{ij} is the **residual** obtained from the regression of X_j on all other Xs

$$X_{ij} = \hat{X}_{ij} + \tilde{X}_{ij} = \hat{\delta}_0 + \hat{\delta}_1 X_{i1} + \dots + \hat{\delta}_{j-1} X_{i,j-1} + \hat{\delta}_{j+1} X_{i,j+1} + \dots \hat{\delta}_k X_{ik} + \tilde{X}_{ij}$$

• Then the **R-Squared** of this partitioned regression is

$$R_j^2 = 1 - \frac{SSR_j}{TSS_j}$$

$$\Rightarrow SSR_j = (1 - R_j^2) \times TSS_j$$

$$\Rightarrow \sum_{i=1}^n \tilde{X}_{ij}^2 = (1 - R_j^2) \sum_{i=1}^n (X_{ji} - \bar{X}_j)^2 = (1 - R_j^2)(n - 1)s_j^2$$

• Where SSR_j is the sum of the squared residuals and TSS is the total variances of X_j . And $s_j^2 = \frac{\sum_{i=1}^n (X_{ji} - \bar{X}_j)^2}{n-1}$ is the sample variance of X_j .

The Variance of $\hat{\beta}_j$ under Homoskedasticity

$$Var\left(\hat{\beta}_{j}\right) = \sigma_{\hat{\beta}_{j}}^{2} = \frac{\sigma_{u}^{2}}{(n-1)s_{j}^{2}(1-R_{j}^{2})}$$

The Variance of \hat{eta}_j under Homoskedasticity

$$Var\left(\hat{\beta}_{j}\right) = \sigma_{\hat{\beta}_{j}}^{2} = \frac{\sigma_{u}^{2}}{(n-1)s_{j}^{2}(1-R_{j}^{2})}$$

- How does the variance of \hat{eta}_j change with the following factors?

The Variance of \hat{eta}_j under Homoskedasticity

$$Var\left(\hat{\beta}_{j}\right) = \sigma_{\hat{\beta}_{j}}^{2} = \frac{\sigma_{u}^{2}}{(n-1)s_{j}^{2}(1-R_{j}^{2})}$$

- How does the variance of \hat{eta}_j change with the following factors?

Factors	symbols	$Var\left(\hat{eta}_{j} ight)$
the variance of u_i	$\sigma_{u}^{2}\uparrow$	1
the sample variance of X_j	$s_{i}^{2}\uparrow$	\downarrow
the R_i^2	$R_i^2 \uparrow$	\uparrow
the sample size	$n\uparrow$	\downarrow

Control Variables
Control Variables: *W*

• The basic regression model is

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

- X_i is the variable of interest and W_i is the control variable which is **NOT** the variable of interest.
- Based on the relationships between W and Y and W and X, we can classify the control variables into several categories:
 - **1**. Whether W has an effect on Y or not:
 - **Relevant Variables**: *W* has a **NONZERO** partial effect on the dependent variable *Y*.
 - Irrelevant Variables: *W* has a ZERO partial effect on the dependent variable *Y*.
 - **2**. Whether W is correlated with X or not:
 - **Uncorrelated Control Variables**: *W* is **NOT** correlated with *X*.
 - **Correlated Control Variables**: *W* is correlated with *X*.
 - Highly-Correlated Control Variables: *W* is highly correlated with *X*.

• We will discuss the potential bias or the precise of the OLS estimators when the control variable is as follows:

W	Irrelevant to Y	Relevant to Y
Uncorrelated	Irrelevant Variables	Non-Omitted Variables
with X		
Correlated with	Irrelevant Variables	Omitted Variables
Х		
Highly-	(Worse)Irrelevant Variables	Omitted Variables and
Correlated with		Multicollinearity
Х		

- Tell me following variables can be classified into which categories?
 - 1. gender composition of the class
 - 2. the weather(temperature) of the exams
 - 3. the share of English learners in the class
 - 4. the size of the classroom

- Assume that we have an **irrelevant** control variable W into the model, thus the model is

$$Y_i = \beta_0 + \beta_1 X_i + \gamma W_i + u_i \tag{7.1}$$

- Since W is ${\bf irrelevant}$ to Y, thus

$$\gamma = 0$$

- Then the model excluding \boldsymbol{W} is

$$Y_i = \tilde{\beta}_0 + \tilde{\beta}_1 X_i + v_i \tag{7.2}$$

• where $v_i = \gamma W_i + u_i$, then $\beta_0 = \tilde{\beta}_0$ and $\beta_1 = \tilde{\beta}_1$.

• Then based on the OVB formula for (7.2),we have

$$plim\hat{\beta}_1 = \tilde{\beta}_1 + \gamma \frac{Cov(X_i, W_i)}{VarX_i} = \tilde{\beta}_1 = \beta_1$$

• The OLS estimator $\hat{\beta}_1$ is still consistent whether you include W or not.

• The variance of \hat{eta}_1 in 7.1 is

$$Var(\hat{\beta}_{1}) = \frac{\sigma_{v}^{2}}{\sum_{i=1}^{n} (X_{i} - \bar{X})^{2} (1 - R_{xw}^{2})}$$

Where R^2_{xw} is the R-Squared of the regression of X on W - The variance of $\hat{\hat{\beta}}$ in 7.2 is

$$Var(\hat{\tilde{\beta}}) = \frac{\sigma_v^2}{\sum_{i=1}^n (X_i - \bar{X})^2}$$
(7.4)

(7.3)

Irrelevant Variables: Variance

• Because
$$v_i = \gamma W_i + u_i$$
 and $\gamma = 0$, thus

$$Var(v_i) = Var(u_i) \Rightarrow \sigma_v^2 = \sigma_u^2$$

- Based on the relationship between W and X in three cases as follows

$$Cov(X_i, W_i) \begin{cases} = 0 & \text{if } X_i \text{ is not correlated with } W_i \\ \neq 0 & \text{if } X_i \text{ is correlated with } X_i \\ \rightarrow 1 & \text{if } X_i \text{ is highly correlated with } X_i : \text{multicollinearity} \end{cases}$$

• Then

.

$$R_{xw}^{2} \begin{cases} = 0 & \text{if } X_{i} \text{ is not correlated with } W_{i} \\ \neq 0 & \text{if } X_{i} \text{ is correlated with } X_{i} \\ \neq 0 \text{ and } \rho \to 1 & \text{if } X_{i} \text{ is highly correlated with } X_{i} : \text{multicollinearity} \end{cases}$$

- Recall: **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*.

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(NOT perfect multicollinearity).

- 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).
- Recall: the variance of \hat{eta}_j is

$$Var(\hat{\beta}_{j}) = rac{\sigma_{u}^{2}}{(n-1)s_{j}^{2}(1-R_{j}^{2})}$$

• When R_j^2 is high, the variance of $\hat{\beta}_j$ is high.

Multicollinearity

- To what extent two correlated variables can be seen as "highly correlated"?
- Rule of Thumb:
 - Mild multicollinearity: correlation of 0.5-0.7
 - High multicollinearity: correlation above 0.7-0.8
 - Severe multicollinearity: correlation above 0.9
- Variance Inflation Factor (VIF) is a more comprehensive measure:

$$VIF_j = \frac{1}{1 - R_j^2}$$

- where R_i^2 is the R-Squared of the regression of X_j on all other Xs.
- Rule of Thumb:
 - Generally acceptable: VIF<5
 - Moderate multicollinearity: $5 \le VIF < 10$
 - Severe multicollinearity: $VIF \ge 10$

Venn Diagrams for Multiple Regression Model



- In a simple model (y on X), OLS uses 'Blue' + 'Red' to estimate β.
- When y is regressed on X and W: OLS throws away the red area and just uses blue to estimate β.
- 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





Venn Diagrams for Multicollinearity



• Less information (compare the Blue and Green areas in both figures) is used, the estimation is less precise.

Irrelevant Variables: Warp Up

• Then we have

1. if irrelevant variable W_i is not correlated with X_i , then $R_{xw}^2 = 0$ and

$$Var(\hat{\beta}_1) = Var(\hat{\tilde{\beta}}_1)$$

2. if irrelevant variable W_i is correlated with X_i , then $R^2_{xw} \neq 0$ and

 $Var(\hat{\beta}_1) > Var(\hat{\tilde{\beta}}_1)$

3. if irrelevant variable W_i is highly correlated with X_i , then $R_{xw}^2 \neq 0$ and

$$Var(\hat{\beta}_1) >> Var(\hat{\tilde{\beta}}_1)$$

- What will happen if controlling an irrelevant variable in a regression?*
 - 1. The OLS estimator is still unbiased and consistent.
 - 2. It increase the variance of estimator, in other words, make the estimate less precise.
 - 3. If multicollinearity exists, it will make the estimate very imprecise.
- Conclusion: we should avoid to put irrelevant variables into our regression.

Relevant Variables and Non-Omitted Variables: Estimate

• Our regression models is still (7.1) and (7.2), but *W* now is not an irrelevant variable but a **Non-omitted variable**, thus

 $Cov(X_i, W_i) = 0 \& \gamma \neq 0$

• Then based on the OVB formula for (7.2),we still have

$$plim\hat{\hat{\beta}}_1 = \tilde{\beta}_1 + \gamma \frac{Cov(X_i, W_i)}{VarX_i} = \tilde{\beta}_1$$

- The OLS estimator $\hat{\beta}_1$ is still consistent whether you include W or not.

Relevant Variables and Non-Omitted Variables: Variance

• Because $\gamma \neq 0$ and $v_i = \gamma W_i + u_i$,thus

$$Var(u_i) \le Var(v_i) \Rightarrow \sigma_u^2 \le \sigma_v^2$$

• Since we have $Cov(X_i, W_i) = 0$, thus $R_{xw}^2 = 0$, then

$$Var(\hat{\beta}_1) \le Var(\hat{\tilde{\beta}}_1)$$

- It **decrease** the variance of estimator, in other words, it will make the estimate **more precise**.
- **Conclusion**: we should always put Relevant but Non-Omitted Variables into our regression.

Relevant Variables: Omitted Variables

• What about a **Relevant Variables** and correlated with X in our regression model?

```
Cov(X_i, W_i) \neq 0 \& \gamma \neq 0
```

- thus the standard definition of **Omitted Variable Bias** if we left it out in our regression model.
- Then based on the OVB formula for (7.2),we still have

$$plim\hat{\hat{\beta}}_1 = \tilde{\beta}_1 + \gamma \frac{Cov(X_i, W_i)}{VarX_i} \neq \tilde{\beta}_1$$

- The OLS estimator $\hat{\beta}_1$ is inconsistent if you do not include W in the regression model.
- **Conclusion**: We should put Omitted Variables into our regression.

Good Controls v.s Bad Controls

Bad Controls v.s Omitted Variable Bias

- It seems that controlling for more covariates always increases the likelihood that regression estimates have a causal interpretation.
 - often true, but not always.
- eg. Some researchers regressing earnings(Y_i) on schooling(S_i) (and experience) include controls for occupation(O_i). Thus our regression model is

$$Y_i = \beta_0 + \beta_1 S_i + \gamma O_i + u_i$$

where β_1 is the most of interest coefficient.

- Clearly we can also think of schooling(S_i) affecting the access to higher level occupations(O_i),
 - e.g. you need a Ph.D. to become a university professor. thus

$$O_i = \lambda_0 + \lambda_1 S_i + e_i$$

Bad Controls v.s Omitted Variable Bias

• Assume that the true relation is a two equation system: a simultaneous equations system

$$Y_i = \beta_0 + \beta_1 S_i + \gamma O_i + e_i$$
$$O_i = \lambda_0 + \lambda_1 S_i + u_i$$

- In the case, Occupation O_i is an *endogenous variable*.
- As a result, you could not necessarily estimate the first equation by OLS, which means that the estimation of β_1 is not *unbiased* and *consistent*, because of controlling Occupation(O_i).

- Let us come back to the wage premium of college graduation: the conditional expectation.But now we have additional control variable-*occupations*: *white-color* and *blue-olor*
- Two reasonable assumptions:
 - 1. white-collar jobs, on average, pay more than blue-collar jobs.
 - 2. graduating college increases the likelihood of a white-collar job.
- **Question 1**: Is occupation an omitted variable in the regression of college degree on wage?
- **Question 2**: Should we control for occupations when considering the effect of college graduation on wages?

- Assume that college degrees are randomly assigned, then we just need to compare the wage difference between workers with college degrees and those without degrees.
- Now we **control** the occupation, which means when we do as follows conditional on occupation:
 - compare degree-earners who chose blue-collar jobs to non-degree-earners who chose blue-collar jobs.
 - or compare degree-earners who chose white-collar jobs to non-degree-earners who chose white-collar jobs.
- Note: the assumption of random degrees says nothing about random job selection.

More formally,

- \mathbf{Y}_i denotes *i*'s earnings
- W_i is also a dummy for whether individual i has a white-collar job
- D_i a dummy variable, refers to *i*'s college-graduation status which is randomly assigned, which indicates

$$\{Y_1, Y_0 \perp D\}$$
 and $\{W_1, W_0 \perp D\}$

Then

$$\mathbf{Y}_i = \mathbf{D}_i \mathbf{Y}_{1i} + (1 - \mathbf{D}_i) \mathbf{Y}_{0i}$$
$$\mathbf{W}_i = \mathbf{D}_i \mathbf{W}_{1i} + (1 - \mathbf{D}_i) \mathbf{W}_{0i}$$

• Because we've assumed D_i is randomly assigned, differences in means yield causal estimates, *i.e.*

$$E[\mathbf{Y}_i \mid \mathbf{D}_i = 1] - E[\mathbf{Y}_i \mid \mathbf{D}_i = 0] = E[\mathbf{Y}_{1i} - \mathbf{Y}_{0i}]$$
$$E[\mathbf{W}_i \mid \mathbf{D}_i = 1] - E[\mathbf{W}_i \mid \mathbf{D}_i = 0] = E[\mathbf{W}_{1i} - \mathbf{W}_{0i}]$$

$$E[\mathbf{Y}_{i} | \mathbf{W}_{i} = 1, \mathbf{D}_{i} = 1] - E[\mathbf{Y}_{i} | \mathbf{W}_{i} = 1, \mathbf{D}_{i} = 0]$$

$$E[\mathbf{Y}_{i} | \mathbf{W}_{i} = 1, \mathbf{D}_{i} = 1] - E[\mathbf{Y}_{i} | \mathbf{W}_{i} = 1, \mathbf{D}_{i} = 0]$$

= $E[\mathbf{Y}_{1i} | \mathbf{W}_{1i} = 1, \mathbf{D}_{i} = 1] - E[\mathbf{Y}_{0i} | \mathbf{W}_{0i} = 1, \mathbf{D}_{i} = 0]$

$$E[\mathbf{Y}_{i} | \mathbf{W}_{i} = 1, \mathbf{D}_{i} = 1] - E[\mathbf{Y}_{i} | \mathbf{W}_{i} = 1, \mathbf{D}_{i} = 0]$$

= $E[\mathbf{Y}_{1i} | \mathbf{W}_{1i} = 1, \mathbf{D}_{i} = 1] - E[\mathbf{Y}_{0i} | \mathbf{W}_{0i} = 1, \mathbf{D}_{i} = 0]$
= $E[\mathbf{Y}_{1i} | \mathbf{W}_{1i} = 1] - E[\mathbf{Y}_{0i} | \mathbf{W}_{0i} = 1]$

$$E [\mathbf{Y}_{i} | \mathbf{W}_{i} = 1, \mathbf{D}_{i} = 1] - E [\mathbf{Y}_{i} | \mathbf{W}_{i} = 1, \mathbf{D}_{i} = 0]$$

= $E [\mathbf{Y}_{1i} | \mathbf{W}_{1i} = 1, \mathbf{D}_{i} = 1] - E [\mathbf{Y}_{0i} | \mathbf{W}_{0i} = 1, \mathbf{D}_{i} = 0]$
= $E [\mathbf{Y}_{1i} | \mathbf{W}_{1i} = 1] - E [\mathbf{Y}_{0i} | \mathbf{W}_{0i} = 1]$
= $E [\mathbf{Y}_{1i} | \mathbf{W}_{1i} = 1] - E [\mathbf{Y}_{0i} | \mathbf{W}_{1i} = 1] + E [\mathbf{Y}_{0i} | \mathbf{W}_{1i} = 1] - E [\mathbf{Y}_{0i} | \mathbf{W}_{0i} = 1]$

• What happens when we estimate the wage-effect of college graduation for white-collar jobs by controlling occupations?

$$\begin{split} E\left[\mathbf{Y}_{i} \mid \mathbf{W}_{i} = 1, \ \mathbf{D}_{i} = 1\right] - E\left[\mathbf{Y}_{i} \mid \mathbf{W}_{i} = 1, \ \mathbf{D}_{i} = 0\right] \\ = E\left[\mathbf{Y}_{1i} \mid \mathbf{W}_{1i} = 1, \ \mathbf{D}_{i} = 1\right] - E\left[\mathbf{Y}_{0i} \mid \mathbf{W}_{0i} = 1, \ \mathbf{D}_{i} = 0\right] \\ = E\left[\mathbf{Y}_{1i} \mid \mathbf{W}_{1i} = 1\right] - E\left[\mathbf{Y}_{0i} \mid \mathbf{W}_{0i} = 1\right] \\ = E\left[\mathbf{Y}_{1i} \mid \mathbf{W}_{1i} = 1\right] - E\left[\mathbf{Y}_{0i} \mid \mathbf{W}_{1i} = 1\right] + E\left[\mathbf{Y}_{0i} \mid \mathbf{W}_{1i} = 1\right] - E\left[\mathbf{Y}_{0i} \mid \mathbf{W}_{0i} = 1\right] \\ = \underbrace{E\left[\mathbf{Y}_{1i} - \mathbf{Y}_{0i} \mid \mathbf{W}_{1i} = 1\right]}_{\mathbf{ATT on white-collar workers}} + \underbrace{E\left[\mathbf{Y}_{0i} \mid \mathbf{W}_{1i} = 1\right] - E\left[\mathbf{Y}_{0i} \mid \mathbf{W}_{0i} = 1\right]}_{\mathbf{Selection bias}} \end{split}$$

• By introducing a *bad control*, we introduced **selection bias** into a setting that did not have selection bias without controls.

• Specifically,

$$\underbrace{E\left[\mathbf{Y}_{1i} - \mathbf{Y}_{0i} \mid \mathbf{W}_{1i} = 1\right]}_{\text{ATT on white-collar workers}} + \underbrace{E\left[\mathbf{Y}_{0i} \mid \mathbf{W}_{1i} = 1\right] - E\left[\mathbf{Y}_{0i} \mid \mathbf{W}_{0i} = 1\right]}_{\text{Selection bias}}$$

- The First term: Expected potential non-college earnings, given that potential white collar status associated with college education is equal to 1.
- If the occupational choice between white-collor and blue-collor is randomly assigned, then

$$E[\mathbf{Y}_{0i} | \mathbf{W}_{1i} = 1] = E[\mathbf{Y}_{0i} | \mathbf{W}_{0i} = 1]$$

- It describes how college graduation changes the composition of the pool of white-collar workers, which in turn change the wage premium between college and high school graduates.
- Even if the true wage causal effect is zero, this selection bias need not be zero.

Bad Controls v.s Omitted Variable Bias

- Bad Controls: Including too many control variables can be counterproductive. They may introduce post-treatment bias when these variables are themselves outcomes of the treatment variable X (essentially functioning as another dependent variable).
- OVB: but if you don't control more variables,you may suffer OVB, which also lead a unbiased and inconsistent estimate.

Bad Controls v.s Omitted Variable Bias

- Bad Controls: Including too many control variables can be counterproductive. They may introduce post-treatment bias when these variables are themselves outcomes of the treatment variable X (essentially functioning as another dependent variable).
- OVB: but if you don't control more variables,you may suffer OVB, which also lead a unbiased and inconsistent estimate.



- "A Hard and Unsolved Problem in Social Sciences" by Gary King(2010), Statistician at Harvard University.
- It is the most "artistic" part of the econometrics in my opinion.

Wrap up

- Which variables should be included on the right-hand side of a regression equation?
 - 1. **Relevant and Omitted Variables**: These are variables that correlate with both the treatment and the outcome.
 - However, be careful of **bad controls**, as they can introduce more bias.
 - 2. **Relevant but Non-Omitted Variables**: These are variables that don't correlate with the treatment but do correlate with the outcome.
 - Including these variables may help reduce standard errors.
- Which variables should NOT be included on the right-hand side of the equation?
 - Variables that are irrelevant.
 - Variables that are **highly correlated with treatment variable** as they may introduce **multicollinearity**.
 - Variables that are **outcome of the treatment variable** as they may introduce **post-treatment bias**.

DAGs and Control Variables

- **Directed Acyclic Graphs (DAGs)** graphically illustrates the causal relationships and non-causal associations within a network of random variables.
 - Like Mind Map, but more complex for causal inference.
- It can be seen as the other framework to think about the **causality** between variables, besides the **potential outcome framework** we have learned in the second lecture.
 - In my personal opinion, it is a more intuitive and easier way to understand the **causality** between variables.
 - especially, it can help us identify **bad controls** and **omitted variable bias**.
Judea Pearl and Causal DAGs



Computer Scientist, Turing Award
Winner in 2011.

- IUDEA PEARL WINNER OF THE TURING AND DANA MACKENZIE THE BOOK OF WHY THE NEW SCIENCE OF CAUSE AND EFFECT
- "for fundamental contributions to artificial intelligence through the development of a calculus for probabilistic and causal reasoning".



• In graph theory, a graph is a collection of **nodes** connected by **edges**



- Nodes(结点) connected by an edge(边或者连线) are called adjacent(邻近).
- Paths run along adjacent nodes: A B C
- The graph above is **undirected**, since the edges don't have direction.

Directed Graphs

- Directed graphs have edges with direction: $A \rightarrow B \rightarrow C$



- Ancestors are nodes that precede a given node in a directed path.
- Descendants come after the ancestor node. eg. D is a descendant of C.

Cycles in Graphs

• If a node is its own descendant, then the graph has a cycle.



• **Directed Acyclic Graph(DAG)** is the directed graph which does not have any cycles.

- Focusing the relationship between variables(nodes)
 - Dependent or Independent
- 1. Two unconnected nodes



• A and B are independent—no link between the nodes.

2. Two connected nodes



• There is clear (causal) dependence: A is a cause of B.

3. Chains



- A and B are dependent.
- B and C are dependent.
- Then are A and C dependent?
 - Mostly yes.through the chain A-B-C.

3. Chains with conditions



- Quesiton: how does conditioning on B affect the association between A and C?
- Answer: It breaks the chain between A and C, thus A and C are independent.

4. Forks or Confounds



- A and C are usually associated in forks.
 - B induces changes in A
 - B also induces changes in C
- A and C are associated due to their common cause, typically OVB.

4. Blocked forks



- Question: What happens when we condition on B?
- Answer: Conditioning on B makes A and C independent.
 - A and C are only associated due to their common cause B.
 - When we shutdown (hold constant) this common cause (B), there is no way for A and C to associate.

5. Immoralities or Colliders



• An **immorality** occurs when two nodes share a child without being otherwise connected

$$\mathbf{A}
ightarrow \mathbf{B} \leftarrow \mathbf{C}$$

• The child (here:B) at the center of this immorality is called a collider.

5. Immoralities or Colliders



- Question: Are A and C independent?
- Answer: Yes. A \perp C.
 - Causal effects flow from A and C and stop there.
 - Neither A nor C is a descendant of the other.
 - A and C do not share any common causes.

5. Immoralities with conditions



- Question: What happens when we condition on B?
- Answer: We unblock or open the previously blocked (closed) path.
- While A and C are independent, they are conditionally dependent.
- When you condition on a collider, you open up the path.

- A path between X and Y is **blocked** by conditioning on a set of variables Z if either of the following statements is true:
 - 1. On the path, there is a chain $(\dots \to W \to \dots)$ or a fork $(\dots \leftarrow W \to \dots)$, and we condition on $W (W \in Z)$.
 - 2. On the path, there is a collider $(\cdots \rightarrow W \leftarrow \dots)$, and we do NOT condition on W $(W \not\in Z)$ or any of its descendants $(de(W) \not\subseteq Z)$.



- Nodes are random variables.
 - *Y* is the dependent variable
 - *D* is the treatment variable
 - W is the Omitted variable



- Edges depict causal links.
- Causality flows in the direction of the arrows.
- Both connections and non-connections matter!



- There are two pathways from D to Y:
 - 1. The path from D to Y is our casual relationship of interest.
 - 2. The path $Y \leftarrow W \rightarrow D$ creates non-causal association between D and Y.



• To shut down this pathway creating a non-causal association, we have to **control/balance/adjust** on W.

Mediation in a DAG



- Question: How to control variables to isolate the causal effect of D on Y?
- Answer: Control on W and Don't control on M

Partial Mediation in a DAG



- Question: How to control variables to isolate the causal effect of D on Y?
- Answer: Control on W and Don't control on M, either.
 - Why? Because our causal effect of interest is an aggregate effect, not a partial effect.
 - Or you could underestimate the effect of D on Y.

DAGs: Example 4 Non-mediator descendants



- Question: How to control variables to isolate the causal effect of D on Y?
- Answer: Control on W and Don't control on Z.
 - here Z is a collider variable, controlling on it will open the backdoor path.

DAGs: Example 5 M-Bias



- Question: How to control variables to isolate the causal effect of D on Y?
- Answer: Control on W and Don't control on C.
 - here C is a collider variable, controlling on it will open the backdoor path.
- Question: How about controlling on A and B?
- Answer:
 - Control on A which is more like a relevant variable.
 - Don't control on B which is more like a irrelevant and multicollinearity variable.

DAG Applications: Going to College on Earnings

• Our Simple OLS Regression:

$$\log(\mathbf{Y}) = \beta_0 + \beta_1 \mathbf{D} + \varepsilon$$

- where Y is earnings, D is a dummy variable for going to college, ε is the error term. The coefficient β_1 represents the causal effect of going to college on earnings or *the return to college*.
- However, the OLS estimator is biased and inconsistent because of the OVB and confounding factors.
- There are some other factors that affect the wage and schooling as well.
 - Some are observed: Parental education(P), Occupation(O), Social Network(N), etc.
 - Some are **unobserved**: Ability from genetics(A),etc.

DAG Applications: Going to College on Earnings

- Suppose
 - Children's Schooling(D) can affect Occupation(O), Earnings(Y) and Health(H).
 - Parents' SES(P) can affect Children's Schooling(D) and Occupation(O).
 - Children's Earnings(Y) can affect their Health(H).
 - Children's Occupation(O) can affect their Earnings(Y).
 - Ability from genetics(A) of family can affect Children's schooling(D) and Parents' SES(P).

DAG Applications: Going to College on Earnings

• Then the DAG is:



- The path from D to Y is our causal relationship of interest. Then list all the paths from D to Y.
 - * $D \rightarrow Y$ (causal relationship of interest)
 - $D \leftarrow P \rightarrow O \rightarrow Y$ (Partial Mediator)
 - $D \leftarrow P \rightarrow O \rightarrow Y$ (confounder)
 - $D \leftarrow A \rightarrow P \rightarrow O \rightarrow Y$ (confounder)
 - $D \rightarrow H \leftarrow Y$ (collider)
- Question: Which variable should we control on?

- DAGs are typically depicted without "noise variables" or disturbances.
- However, these disturbances still exist—they're just "outside of the model."
- Simultaneity defines causality as unidirectional and disallows cycles.
- Dynamics allow a variable to somewhat influence itself, especially in time series.
- Essentially, this is a form of logical deduction that proves to be very useful.

- DAGs are a powerful tool for causal inference and can help us identify **bad controls** and **omitted variable bias**.
- However, they heavily depends on the **knowledge of the causal relationships** and **assumptions** we made.
 - Knowledge of the causal relationships: from theory, model, observation, experience, priors, etc.
 - The best part of DAGs is that it makes our assumptions and mechanisms in our regression model **explicit**.
- DAGs can not be a replacement for statistical models, but rather a complement to them.

Some practical tips about Control Variables

- Use theory and prior research to guide control selection:
 - Tell which are relevant and which are not and which can be OVB
 - Control for variables that are likely to be confounders.
 - Avoid controlling for variables that might be outcomes of the treatment (post-treatment variables or collider variables) and multicollinearity variables.

- Sensitivity analysis:
 - Start with a simple model which only includes the treatment variable and essential controls.
 - Add controls one by one(not too many, can be grouped) to see how they affect your estimates
 - If adding controls dramatically changes your estimate, investigate the reason why.
 - If adding controls does not change your estimate, you can be more confident about your estimate.
 - Be skeptical about "kitchen sink" regressions at first, which include every available control.

• Remember the fundamental trade-off:

- More controls can reduce omitted variable bias.
- But they can also increase variance and potentially introduce new biases.
- The goal is to find the "sweet spot" that balances these concerns.