



ECON526: Quantitative Economics with Data Science Applications

Instrumental Variables

Phil Solimine

philip.solimine@ubc.ca

University of British Columbia

Table of contents

- Overview
- Instrumental Variables
- Example: Returns to Education
- Weak vs Strong Instruments
- Monte-Carlo Tests of 2SLS

Overview

Summary

- Last class we learned about regression, and how we can use it to estimate treatment effects.
- Today we will learn the details of how to estimate the effect of a treatment on an outcome, when the treatment is not randomly assigned.
- Specifically, we will use a method called **instrumental variables**.
- We will also discuss **matching** (the alternative to regression) and **propensity scores**.



Instrumental Variables

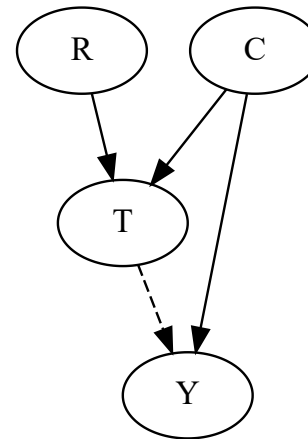
Instrumental Variables

- Suppose we want to estimate the effect of a treatment on an outcome.
- BUT there is a confounder, something that affects both the treatment and the outcome.
 - So we can't distinguish the *pure* effect of the treatment from the effect of the confounder.
 - However, we also have data on a variable that affects the treatment, but does not lie on any dependence path between the treatment and the outcome.

Instrumental Variables

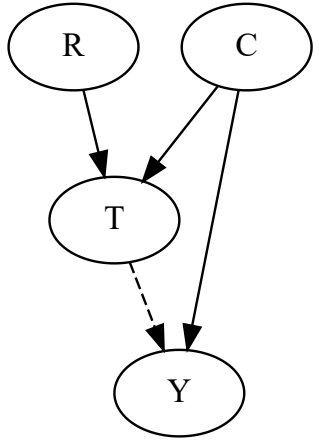
- In the section on DAGs, we saw one example of this, that looked like this:

```
1 import graphviz as gr
2
3 dot = gr.Digraph()
4 dot.edge('T', 'Y', style = 'dashed')
5 dot.edge('C', 'T')
6 dot.edge('C', 'Y')
7 dot.edge('R', 'T')
8 dot
```



- Of course, the regression we *want* to run is $Y = \beta_0 + \beta_1 T + \beta_2 C + \epsilon$.
- But we don't have data on C . We have to settle for
→ $Y = \beta_0 + \kappa T + v$, where $v = \beta C + \epsilon$

Instrumental Variables



- The variable R is an **instrumental variable** for T if
 - R affects T (i.e., R is correlated with T)
 - R does not affect Y except through T (i.e., $R \mid Y \perp T$)
- That means we can *use* R to *mimic* a randomized experiment.

Instrumental Variables

- Since the instrument R is only correlated with the outcome through T , this implies that $\text{cov}(R, \epsilon) = 0$.
 - This is called an **exclusion restriction**.
 - It is the result of our conditional independence assumption.
- We also have that

$$\begin{aligned}\text{cov}(R, Y) &= \text{cov}(R, \beta_0 + \kappa T + v_i) \\ &= \kappa \text{cov}(R, T) + \text{cov}(R, v_i) = \kappa \text{cov}(R, T)\end{aligned}$$

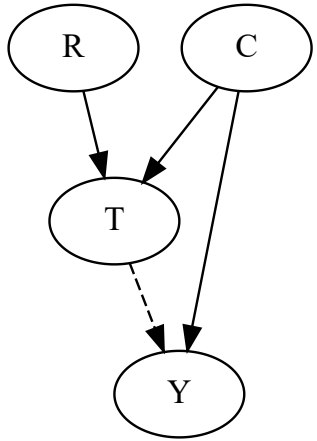
Instrumental Variables

- Divide both sides by $\text{var}(R)$ to get

$$\begin{aligned}\text{cov}(R, Y) &= \kappa \text{cov}(R, T) \\ \text{cov}(R, Y) / \text{var}(R) &= \kappa \text{cov}(R, T) / \text{var}(R) \\ \frac{\text{cov}(R, Y) / \text{var}(R)}{\text{cov}(R, T) / \text{var}(R)} &= \kappa\end{aligned}$$

- Notice that the numerator is regression coefficient of R on Y .
 - The **reduced-form coefficient**
- The denominator is the regression coefficient of T on R .
 - The **first-stage coefficient**
- This amounts to *scaling* the effect of R on Y by the effect of R on T .

Instrumental Variables



- In practice we estimate the reduced-form equation, and use R to predict T .

$$\rightarrow \hat{T} \approx \hat{\beta}_0 + \hat{\beta}_1 R$$

- Then we use the predicted values of T in the outcome equation

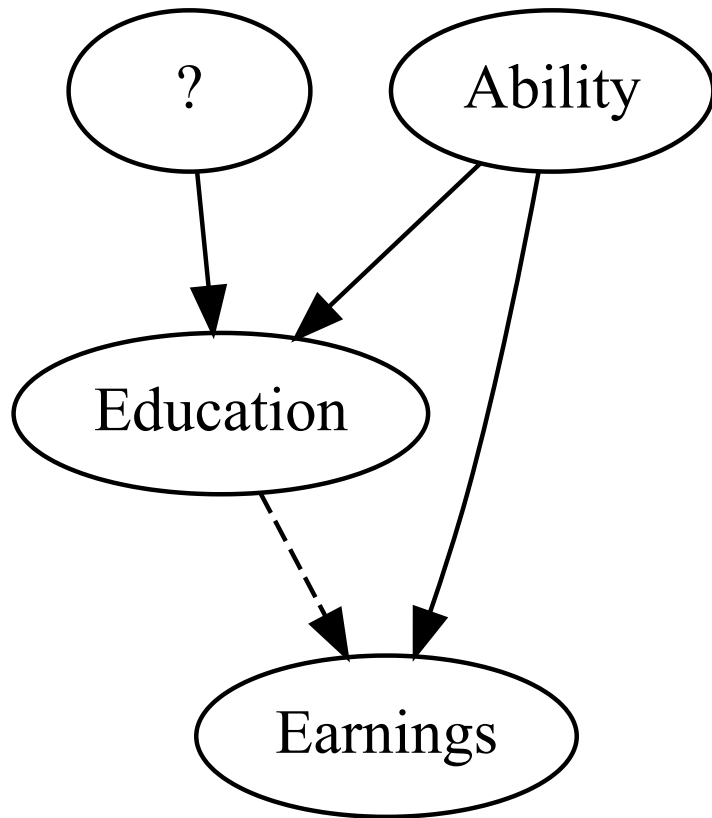
$$\rightarrow Y = \kappa_0 + \kappa \hat{T} + v_i$$

- In other words we are using only the variation in T that is due to R .
- This is called the **two-stage least squares** (2SLS) estimator.

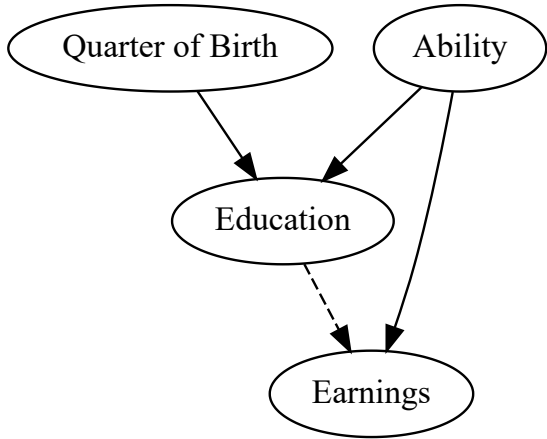
Example: Returns to Education

Example: Returns to Education

- In this example, we'll use data from Angrist and Keueger (QJE, 1991) to estimate the effect of education on earnings.



Example: Returns to Education



```
1 import pandas as pd
2
3 df = pd.read_csv('data/ak91.csv')
4 df.head()
```

	log_wage	years_of_schooling	year_of_birth	quarter_of_birth	state_of_birth
0	5.790019	12.0	30.0	1.0	45.0
1	5.952494	11.0	30.0	1.0	45.0
2	5.315040	12.0	30.0	1.0	45.0

Example: Returns to Education

- We must assume that
 1. $\text{cov}(QOB, Education) \neq 0$
 2. $Earnings \perp QOB \mid Education$

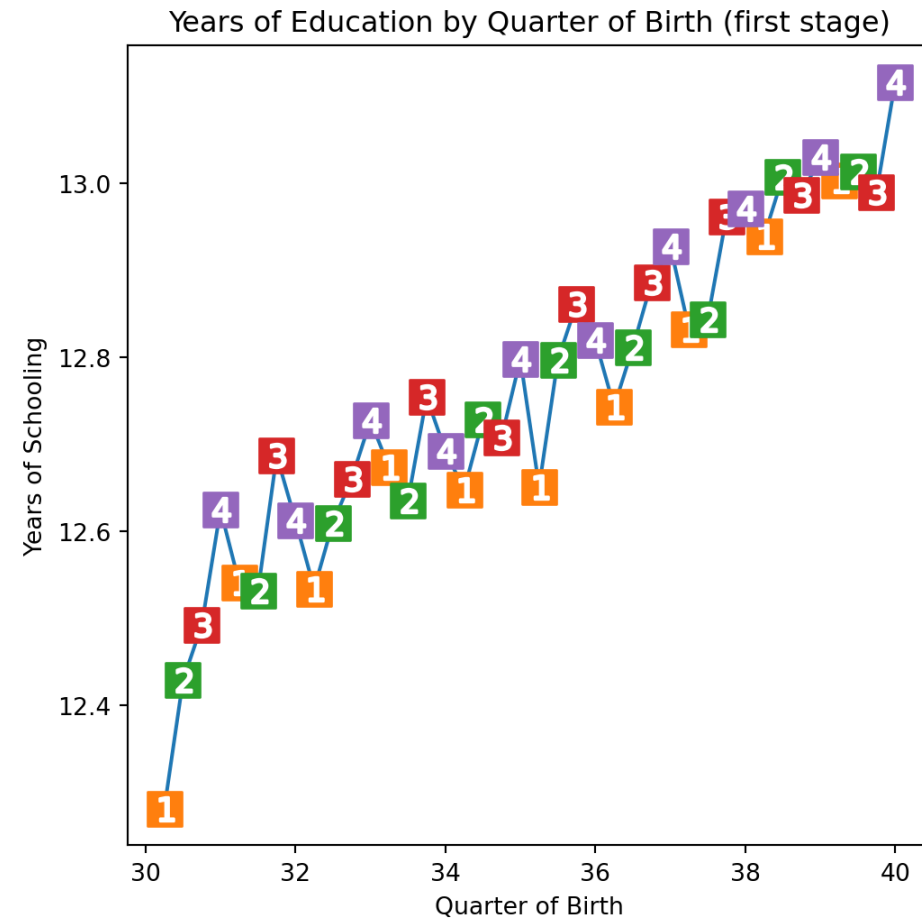
Example: Returns to Education

- The first assumption is easy; we can check directly from the data

```

1 import matplotlib.pyplot as plt
2 group_data = (df
3     .groupby(["year_of_birth", "quarter_of_birth"])
4     [["log_wage", "years_of_schooling"]]
5     .mean()
6     .reset_index()
7     .assign(time_of_birth = lambda d: d["year_of_birth"] + d["quarter_of_birth"] * 4)
8
9 plt.figure(figsize=(6,6))
10 plt.plot(group_data["time_of_birth"], group_data["years_of_schooling"])
11 for q in range(1, 5):
12     x = group_data.query(f"quarter_of_birth=={q}")["year_of_birth"]
13     y = group_data.query(f"quarter_of_birth=={q}")["years_of_schooling"]
14     plt.scatter(x, y, marker="s", s=200, c=f"C{q}")
15     plt.scatter(x, y, marker=f"${q}$", s=100, c=f"wh")
16
17 plt.title("Years of Education by Quarter of Birth (first stage)")

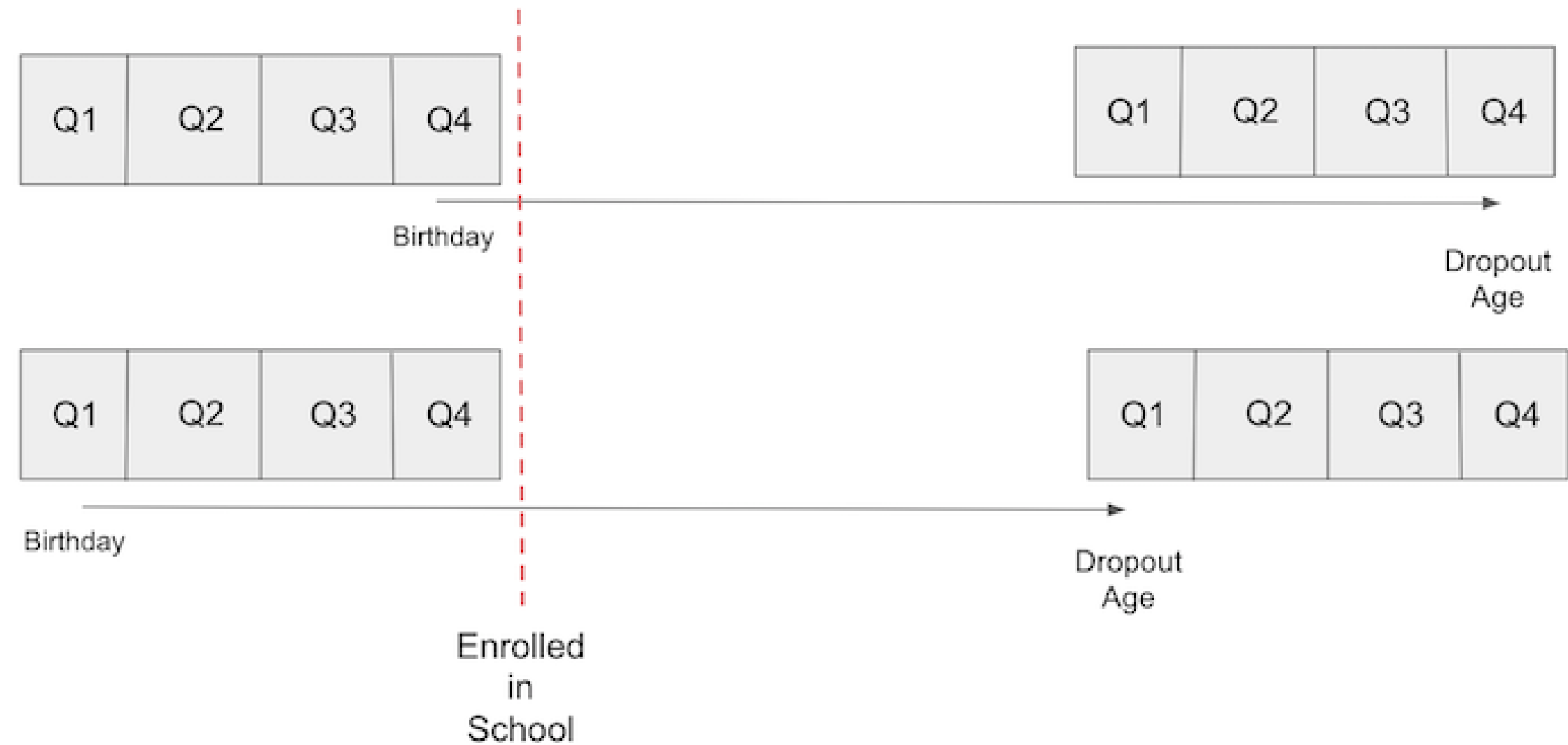
```



Example: Returns to Education

- Angrist and Krueger (1991) argue that the first assumption is reasonable because
 - The quarter of birth is randomly assigned
 - The quarter of birth affects education through compulsory schooling law
- The compulsory schooling law is a law that requires students to stay in school until a certain age.
 - In the US, this age is 16.
 - In Canada, this age is 18.
- The law is enforced by the government, so it is not affected by the ability of the student.

Example: Returns to Education



Example: Returns to Education

- Unfortunately, there is no way to test the second assumption from the data.
 - We have to assume that it is true.
 - In this case, it is hard to come up with another plausible explanation for why quarter of birth would affect earnings.
- In the paper, they also add that there is no correlation of QOB with college graduation rates, indicating that compulsory schooling laws are really driving the effect, not ability.

Example: Returns to Education

- First, we run the first stage. To aid interpretation, we'll use just **q4** as our instrument.

```
1 # Convert the quarter of birth to dummy variables
2 factor_data = df.assign(**{f"q{int(q)}": (df["quarter_of_birth" == q])
3                             for q in df["quarter_of_birth"].unique()})
4
5 # Run the first stage regression
6 import statsmodels.formula.api as smf
7
8 first_stage = smf.ols("years_of_schooling ~ C(year_of_birth) + C(quarter_of_birth)",
9                       data=df).fit()
10 print("q4 parameter estimate:", ", first_stage.params["q4"]
11 print("q4 p-value:", ", first_stage.pvalues["q4"])
```

```
q4 parameter estimate:, 0.10085809272790978
q4 p-value:, 5.464829416479072e-15
```

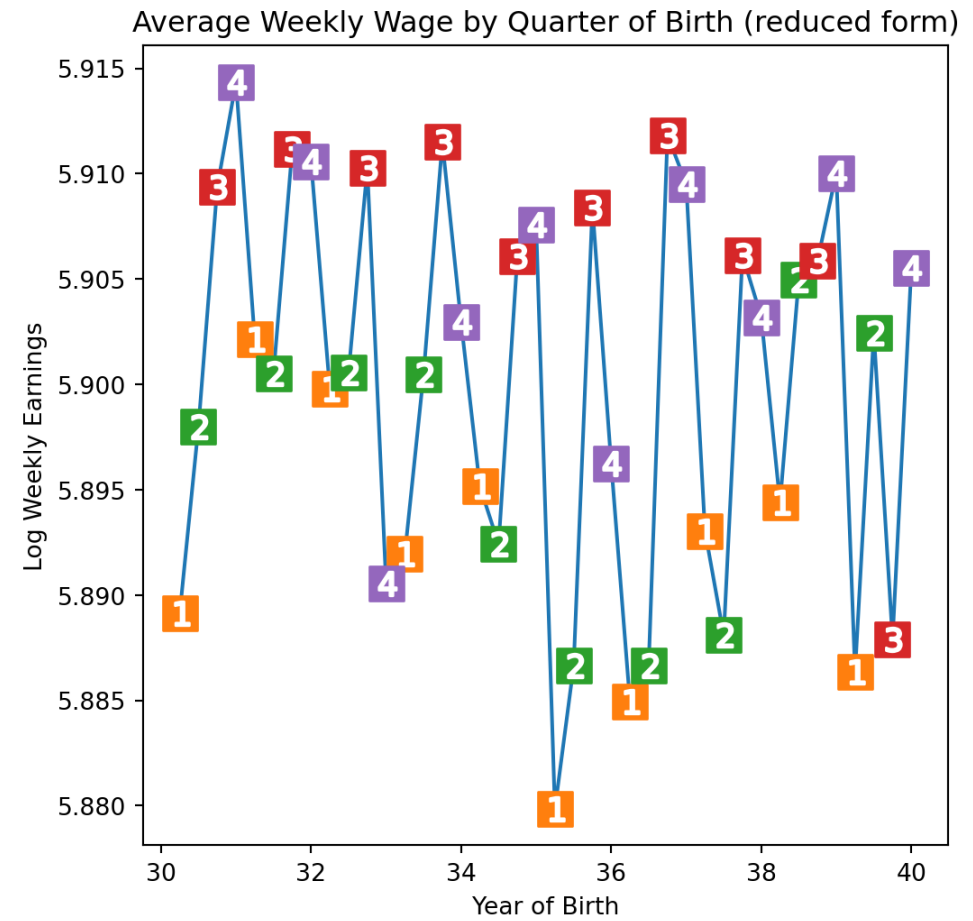
Example: Returns to Education

- Next, we start looking at the second stage. Do we think this will really work?

```

1 # Plot the reduced form
2 plt.figure(figsize=(6,6))
3 plt.plot(group_data["time_of_birth"], group_data["lo
4 for q in range(1, 5):
5     x = group_data.query(f"quarter_of_birth=={q}")["
6     y = group_data.query(f"quarter_of_birth=={q}")["
7     plt.scatter(x, y, marker="s", s=200, c=f"C{q}")
8     plt.scatter(x, y, marker=f"${q}$", s=100, c=f"wh
9
10 plt.title("Average Weekly Wage by Quarter of Birth (
11 plt.xlabel("Year of Birth")
12 plt.ylabel("Log Weekly Earnings");

```



Example: Returns to Education

- The reduced form looks pretty good, but we can't be sure until we run the second stage.

```
1 # Run the reduced form
2 reduced_form = smf.ols("log_wage ~ C(year_of_birth)
3
4 print("q4 parameter estimate:, ", reduced_form.param
5 print("q4 p-value:, ", reduced_form.pvalues["q4"])
```

```
q4 parameter estimate:, 0.008603484260163176
q4 p-value:, 0.0014949127183224312
```

Example: Returns to Education

- The second stage looks good, so we can estimate the effect of education on earnings.
- We can do this “by hand” using the formula for κ .

```
1 ate_iv = reduced_form.params["q4"] / first_stage.par  
2 print("ATE (IV):", ate_iv)
```

ATE (IV): 0.08530286492104555

- This means that we expect earnings to increase by 8% for every additional year of school.

Example: Returns to Education

- Of course, we wouldn't usually. We can use the `IV2SLS` function from `linearmodels`, which will also give us confidence intervals
→ There is also a `IVGMM` function available in `statsmodels`.

```
1 from linearmodels.iv import IV2SLS
2
3 formula = 'log_wage ~ 1 + C(year_of_birth) + C(state_of_birth) + [years_of_schooling ~ q4]'
4 iv = IV2SLS.from_formula(formula, data=factor_data).fit()
5 print(iv.summary)
```

IV-2SLS Estimation Summary			
=====			
Dep. Variable:	log_wage	R-squared:	0.1217
Estimator:	IV-2SLS	Adj. R-squared:	0.1215
No. Observations:	329509	F-statistic:	1.028e+04
Date:	Wed, Nov 01 2023	P-value (F-stat)	0.0000
Time:	00:17:36	Distribution:	chi2(60)
Cov. Estimator:	robust		

Parameter Estimates						
=====						
	Parameter	Std. Err.	T-stat	P-value	Lower CI	Upper CI

Intercept	4.7468	0.2904	16.348	0.0000	4.1777	5.3158

Weak vs Strong Instruments

Weak vs Strong Instruments

- Standard errors for instrumental variables are larger than for OLS.
- The difference comes from the fact that our prediction of T from just R is imperfect.
- The more *correlated* R is with T , the better our prediction will be.
 - If $\text{cov}(R, T)$ is large, we call R a **strong instrument**.
 - If $\text{cov}(R, T)$ is small, we call R a **weak instrument**.

Monte-Carlo Tests of 2SLS

Weak vs Strong Instruments

- To explore the 2SLS estimator, we could use Monte-Carlo simulation.
- We'll simulate data from the following model:

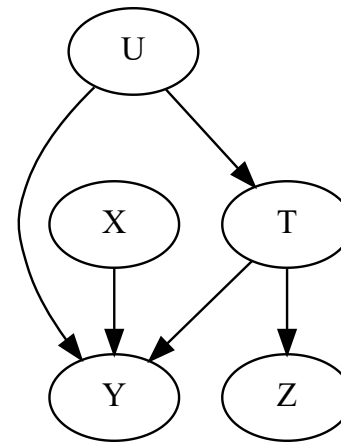
$$X \sim N(0, 2^2)$$

$$U \sim N(0, 2^2)$$

$$T \sim N(1 + 0.5U, 5^2)$$

$$Y \sim N(2 + X - 0.5U + 2T, 5^2)$$

$$Z \sim N(T, \sigma^2) \text{ for } \sigma^2 \text{ in } 0.1 \text{ to } 100$$



Weak vs Strong Instruments

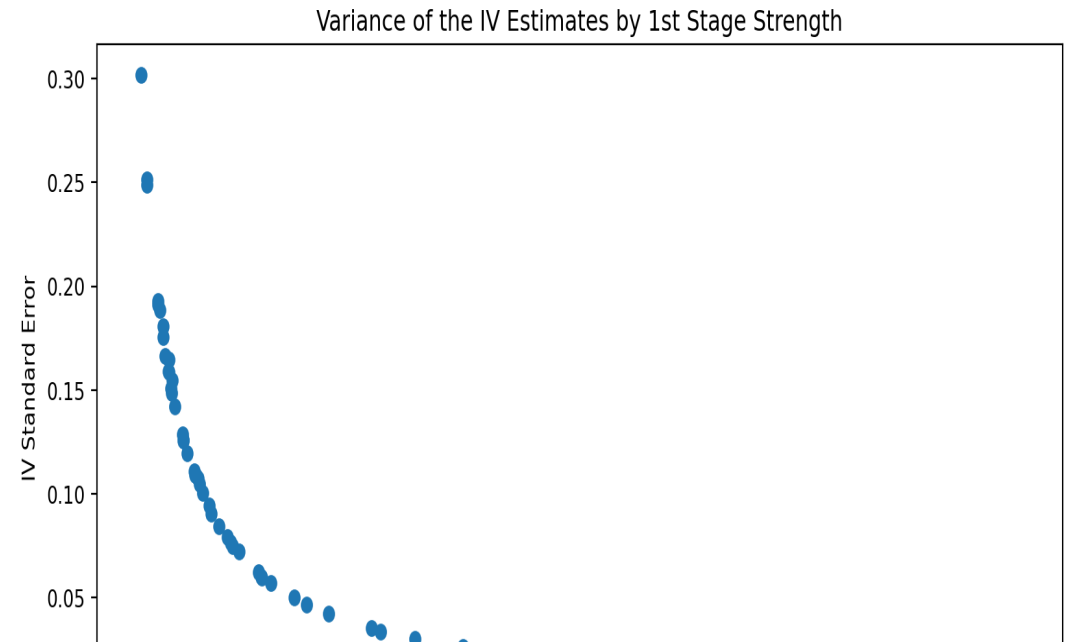
```
1 import numpy as np
2 np.random.seed(12)
3 n = 10000
4 X = np.random.normal(0, 2, n) # observable variable
5 U = np.random.normal(0, 2, n) # unobservable (omitted) variable
6 T = np.random.normal(1 + 0.5*U, 5, n) # treatment
7 Y = np.random.normal(2 + X - 0.5*U + 2*T, 5, n) # outcome
8
9 stddevs = np.linspace(0.1, 100, 50)
10 Zs = {f"Z_{z}": np.random.normal(T, s, n) for z, s in enumerate(stddevs)} # instruments with decreasing Cov(Z, T)
11
12 sim_data = pd.DataFrame(dict(U=U, T=T, Y=Y)).assign(**Zs)
13
14 sim_data.head()
```

	U	T	Y	Z_0	Z_1	Z_2
0	2.696148	8.056988	18.388910	8.233315	9.028779	16.430365
1	2.570240	0.245067	2.015052	0.455988	-0.901285	-6.442245
2	0.664741	5.597510	11.939170	5.528384	6.148148	10.141348

Weak vs Strong Instruments

- Now we can run the 2SLS estimator for each value of σ , as the covariance between the instrument and treatment changes.
 - We can see that the standard error of the 2SLS estimator increases as the covariance decreases.

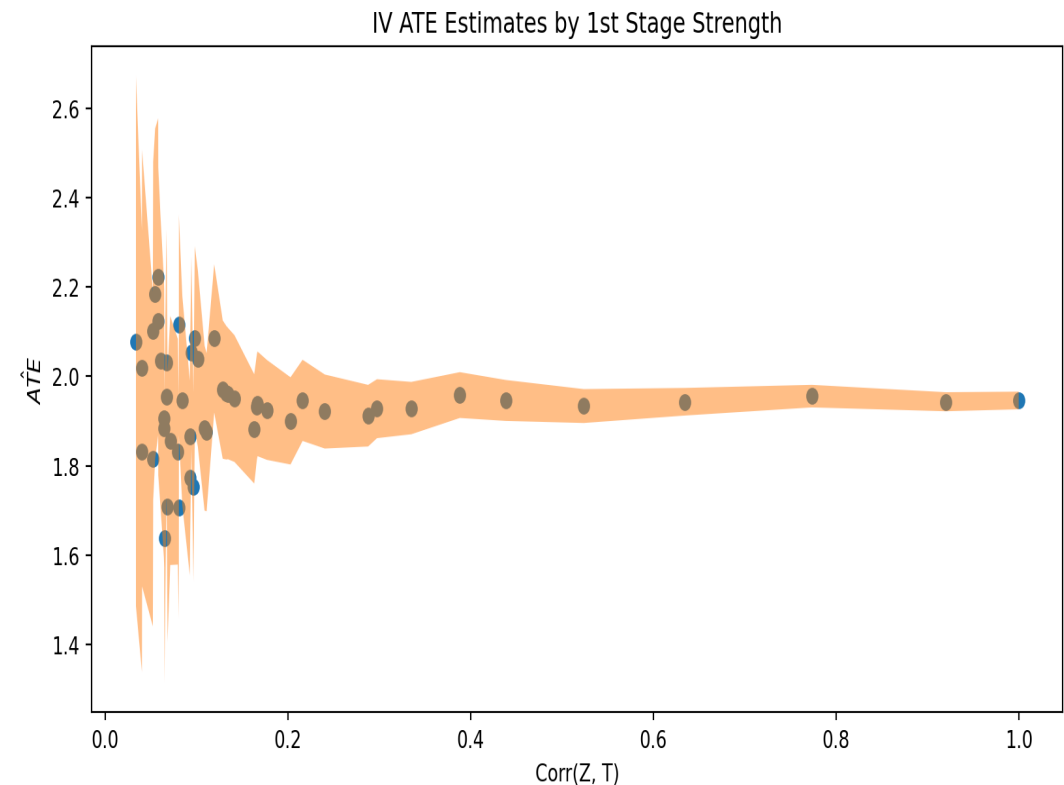
```
1 corr = (sim_data.corr()["T"]
2     [lambda d: d.index.str.startswith("Z")])
3
4 se = []
5 ate = []
6 for z in range(len(Zs)):
7     formula = f'Y ~ 1 + X + [T ~ Z_{z}]'
8     iv = IV2SLS.from_formula(formula, sim_data).fit()
9     se.append(iv.std_errors["T"])
10    ate.append(iv.params["T"])
11
12 plot_data = pd.DataFrame(dict(se=se, ate=ate, corr=c
13
14 plt.scatter(plot_data["corr"], plot_data["se"])
```



Bias of OLS

- Now, let's look at the point estimates themselves

```
1 plt.scatter(plot_data["corr"], plot_data["ate"])
2 plt.fill_between(plot_data["corr"],
3     plot_data["ate"]+1.96*plot_data["se"],
4     plot_data["ate"]-1.96*plot_data["se"], alpha=.5)
5 plt.xlabel("Corr(Z, T)")
6 plt.ylabel("$\hat{ATE}$");
7 plt.title("IV ATE Estimates by 1st Stage Strength");
```



Bias of 2SLS

- We can see that the estimates from 2SLS are still biased.
- This is because the instrument is not able to “mimic” the randomized experiment as well, there will always be some variation in T that is not explained by Z .
- 2SLS is biased in the same direction as OLS would be.
 - But it's consistent!
 - So we can get rid of the bias by increasing the sample size.

Credits

This lecture draws heavily from [Causal Inference for the Brave and True: Chapter 08 - Instrumental Variables](#) by Matheus Facure.

As well as [The Effect: Chapter 19 - Instrumental Variables](#) by Nick Huntington-Klein.