# Dealing with Endogeneity: Instrumental Variable

AECN 396/896-002

# Before we start

## Learning objectives

Understand how instrumental variable (IV) estimation works.

## Table of contents

# Endogeneity

**Endogeneity**

$E[u|x_k] \neq 0$ (the error term is not correlated with any of the independent variables)

# Endogeneity

**Endogeneity**

$E[u|x_k] \neq 0$ (the error term is not correlated with any of the independent variables)

**Endogenous independent variable**

If the error term is, for whatever reason, correlated with the independent variable $x_k$, then we say that $x_k$ is an endogenous independent variable.
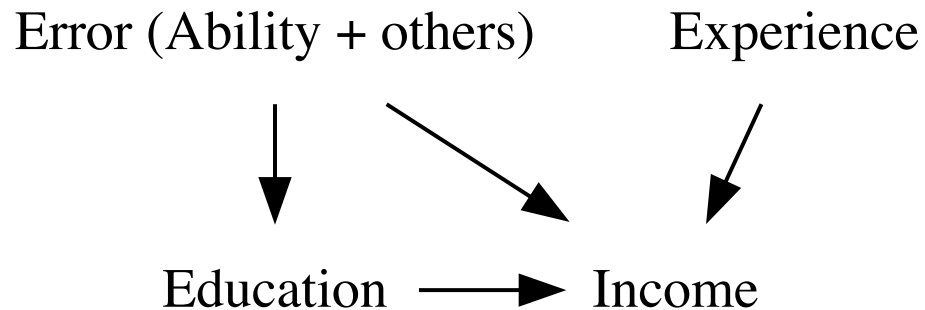
- Omitted variable
- Selection
- Reverse causality
- Measurement error
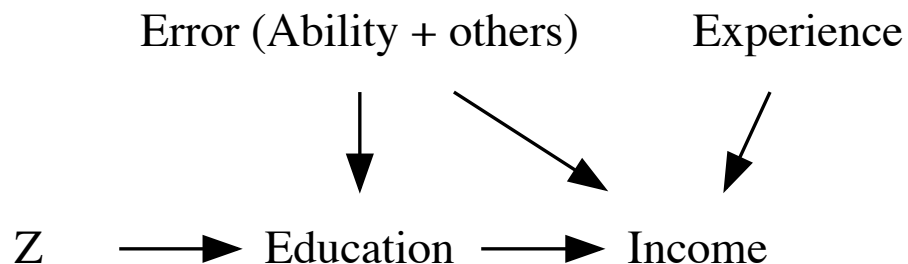
# Instrumental Variable (IV) Approach

# Causal Diagram

You want to estimate the causal impact of education on income.

- Variable of interest: Education
- Dependent variable: Income

Error (Ability + others)     Experience

Education  ⟶  Income

# Rough Idea of IV Approach

Find a variable like $Z$ in the diagram below:

Error (Ability + others)          Experience

$Z$  $\longrightarrow$  Education  $\longrightarrow$  Income

- $Z$ does NOT affect income directly
- $Z$ is correlated with the variable of interest (education)
  - does not matter which causes which (associattion is enough)
- $Z$ is NOT correlated with any of the unobservable variables in the error term (including ability) that is making the vairable of interest (education) endogeneous.
  - $Z$ does not affect ability
  - abiliyt does not affect $Z$

**The Model**

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u$$

- $x_1$ is endogenous: $E[u|x_1] \neq 0$ (or $Cov(u, x_1) \neq 0$)
- $x_2$ is exogenous: $E[u|x_1] = 0$ (or $Cov(u, x_1) = 0$)

**Idea (very loosely put)**

Bring in variable(s) (Instrumental variable(s)) that does NOT belong to the model, but IS related with the endogenous variable,

- Using the instrumental variable(s) (which we denote by $Z$), make the endogenous variable exogenous, which we call instrumented variable(s)

- Use the variation in the instrumented variable instead of the original endogenous variable to estimate the impact of the original variable

# IV estimation procedure

**Step 1**

Using the instrumental variables, make the endogenous variable exogenous, which we call instrumented variable

# IV estimation procedure

Using the instrumental variables, make the endogenous variable exogenous, which we call instrumented variable

- Regress the endogenous variable $(x_1)$ on the instrumental variable(s) $(Z = \{z_1, z_2\}$, two instruments here) and all the other exogenous variables $(x_2$ here)

$$x_1 = \alpha_0 + \sigma_2 x_2 + \alpha_1 z_1 + \alpha_2 z_2 + v$$

# IV estimation procedure

Using the instrumental variables, make the endogenous variable exogenous, which we call instrumented variable

- Regress the endogenous variable $(x_1)$ on the instrumental variable(s) $(Z = \{z_1, z_2\}$, two instruments here) and all the other exogenous variables $(x_2$ here)

$$x_1 = \alpha_0 + \sigma_2 x_2 + \alpha_1 z_1 + \alpha_2 z_2 + v$$

- obtain the predicted value of $x$ from the regression

$$\hat{x}_1 = \hat{\alpha}_0 + \hat{\sigma}_2 x_2 + \hat{\alpha}_1 z_1 + \hat{\alpha}_2 z_2$$

# IV estimation procedure

**Step 2**

use the variation in the instrumented variable instead of the original endogenous variable to estimate the impact of the original variable

# IV estimation procedure

use the variation in the instrumented variable instead of the original endogenous variable to estimate the impact of the original variable

**Step 2: Mathematically**

Regress the dependent variable $(y)$ on the instrumented variable $(\hat{x}_1)$,

$$y = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 x_2 + \varepsilon$$

to estimate the coefficient on $x$ in the original model

# Example

$$log(wage) = \beta_0 + \beta_1 educ + \beta_2 exper + (\beta_3 ability + v)$$

- Regress $log(wage)$ on $educ$ and $exper$ ($ability$ not included because you do not observe it)
- $(\beta_3 ability + v)$ is the error term
- $educ$ is considered endogenous (correlated with $ability$)
- $exper$ is considered exogenous (not correlated with $ability$)

# Example

**Model of interest**

$$log(wage) = \beta_0 + \beta_1 educ + \beta_2 exper + (\beta_3 ability + v)$$

- Regress $log(wage)$ on $educ$ and $exper$ ($ability$ not included because you do not observe it)
- $(\beta_3 ability + v)$ is the error term
- $educ$ is considered endogenous (correlated with $ability$)
- $exper$ is considered exogenous (not correlated with $ability$)

**Instruments (Z)**

Suppose you selected the following variables as instruments:

- IQ test score $(IQ)$
- number of siblings $(sibs)$

Regress $educ$ on $exper$, $IQ$, and $sibs$:

$$educ = \alpha_0 + \alpha_1 exper + \alpha_2 IQ + \alpha_3 sibs + u$$

Use the coefficient estimates on $\alpha_0$, $\alpha_1$, $\alpha_2$, and $\alpha_3$ to predict $educ$ as a function of $exper$, $IQ$, and $sibs$.

$$\hat{educ} = \hat{\alpha}_0 + \hat{\alpha}_1 exper + \hat{\alpha}_2 IQ + \hat{\alpha}_3 sibs$$

```r
library(wooldridge)
data("wage2")

#* regress educ on exper, IQ, and sibs
first_reg <- feols(educ ~ exper + IQ + sibs, data = wage2)

#* predict educ as a function of exper, IQ, and sibs
wage2 <- mutate(wage2, educ_hat = first_reg$fitted.values)

#* seed the predicted values
wage2 %>%
  relocate(educ_hat) %>%
  head()
```

```
##    educ_hat wage hours  IQ KWW educ exper tenure age married b
## 1 13.26398  769    40  93  35   12    11      2  31       1
## 2 14.80686  808    50 119  41   18    11     16  37       1
## 3 14.15410  825    40 108  46   14    11      9  33       1
## 4 12.79569  650    40  96  32   12    13      7  32       1
## 5 10.73631  562    40  74  27   11    14      5  34       1
## 6 14.09006 1400    40 116  43   16    14      2  35       1
```

**Step 2:**

Use $\hat{educ}$ in place of $educ$ to estimate the model of interest:

$$log(wage) = \beta_0 + \beta_1 \hat{educ} + \beta_2 exper + u$$

```
#* regression with educ_hat in place of educ
second_reg <- feols(wage ~ educ_hat + exper, data = wage2)

#* see the results
second_reg
```

```
## OLS estimation, Dep. Var.: wage
## Observations: 935
## Standard-errors: IID
##              Estimate Std. Error  t value   Pr(>|t|)
## (Intercept) -1269.7880  214.12821 -5.93004 4.2632e-09 ***
## educ_hat      138.1051   13.10586 10.53766  < 2.2e-16 ***
## exper          31.7955    4.14489  7.67101 4.2899e-14 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## RMSE: 382.0   Adj. R2: 0.104547
```

# When does IV work?

Just like OLS needs to satisy some conditions for it to consistently estimate the coefficients, IV approach needs to satisy some conditions for it to work.

**Estimation Procedure**

- Step 1: $\hat{x}_1 = \hat{\alpha}_0 + \hat{\sigma}_2 x_2 + \hat{\alpha}_1 z_1 + \hat{\alpha}_2 z_2$

- Step 2: $y = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 x_2 + \varepsilon$

**Important question**

What are the conditions under which IV estimation is consistent?

The instruments $(Z)$ need to satisfy two conditions, which we will discuss.

# Condition 1

- Step 1: $\hat{x}_1 = \hat{\alpha}_0 + \hat{\sigma}_2 x_2 + \hat{\alpha}_1 z_1 + \hat{\alpha}_2 z_2$

- Step 2: $y = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 x_2 + \varepsilon$

**Question**

What happens if $Z$ have no power to explain $x_1$ $(\alpha_1 = 0$ and $\alpha_2 = 0)$?

- Step 1: $\hat{x}_1 = \hat{\alpha}_0 + \hat{\sigma}_2 x_2 + \hat{\alpha}_1 z_1 + \hat{\alpha}_2 z_2$

- Step 2: $y = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 x_2 + \varepsilon$

**Question**

What happens if $Z$ have no power to explain $x_1$ $(\alpha_1 = 0$ and $\alpha_2 = 0)$?

**Answer**

- $\hat{x}_1 = \hat{\alpha}_0 + \hat{\sigma}^2 x_2$
- $\hat{\beta}_1$?

**Estimation Procedure**

- Step 1: $\hat{x}_1 = \hat{\alpha}_0 + \hat{\sigma}_2 x_2 + \hat{\alpha}_1 z_1 + \hat{\alpha}_2 z_2$

- Step 2: $y = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 x_2 + \varepsilon$

**Question**

What happens if $Z$ have no power to explain $x_1$ ($\alpha_1 = 0$ and $\alpha_2 = 0$)?

**Answer**

- $\hat{x}_1 = \hat{\alpha}_0 + \hat{\sigma}^2 x_2$
- $\hat{\beta}_1$?

That is, $\hat{x}_1$ has no information beyond the information $x_2$ possesses.

**Condition 1**

The instrument(s) $Z$ have jointly significant explanatory power on the endogenous variable $x_1$ after you control for all the other exogenous variables (here $x_2$)

# Condition 2

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u$$

- Step 1: $\hat{x}_1 = \hat{\alpha}_0 + \hat{\sigma}_2 x_2 + \hat{\alpha}_1 z_1 + \hat{\alpha}_2 z_2$

- Step 2: $y = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 x_2 + \varepsilon$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u$$

- Step 1: $\hat{x}_1 = \hat{\alpha}_0 + \hat{\sigma}_2 x_2 + \hat{\alpha}_1 z_1 + \hat{\alpha}_2 z_2$

- Step 2: $y = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 x_2 + \varepsilon$

Remember you can break $x_1$ into the predicted part and the residuals.

$$x_1 = \hat{x}_1 + \hat{\varepsilon}$$

where $\hat{\varepsilon}$ is the residual of the first stage estimation.

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u$$

- Step 1: $\hat{x}_1 = \hat{\alpha}_0 + \hat{\sigma}_2 x_2 + \hat{\alpha}_1 z_1 + \hat{\alpha}_2 z_2$

- Step 2: $y = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 x_2 + \varepsilon$

Remember you can break $x_1$ into the predicted part and the residuals.

$$x_1 = \hat{x}_1 + \hat{\varepsilon}$$

where $\hat{\varepsilon}$ is the residual of the first stage estimation.

Plugging in $x_1 = \hat{x}_1 + \hat{\varepsilon}$ into the model of interest,

$$y = \beta_0 + \beta_1(\hat{x}_1 + \hat{\varepsilon}) + \beta_2 x_2 + u$$
$$= \beta_0 + \beta_1 \hat{x}_1 + \beta_2 x_2 + (\beta_1 \hat{\varepsilon} + u)$$

So, if you regress $y$ on $\hat{x}_1$ and $x_2$, then the error term is $(\beta_1 \hat{\varepsilon} + u)$.

**Second stage regression**

$$y = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 x_2 + (\beta_1 \hat{\varepsilon} + u)$$

$$y = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 x_2 + (\beta_1 \hat{\varepsilon} + u)$$

**Question**

What is the condition under which the OLS estimation of $\beta_1$ in the main model is unbiased?

$$y = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 x_2 + (\beta_1 \hat{\varepsilon} + u)$$

**Question**

What is the condition under which the OLS estimation of $\beta_1$ in the main model is unbiased?

**Answer**

$\hat{x}_1$ is not correlated with $(\beta_1 \hat{\varepsilon} + u)$

$$y = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 x_2 + (\beta_1 \hat{\varepsilon} + u)$$

**Question**

What is the condition under which the OLS estimation of $\beta_1$ in the main model is unbiased?

**Answer**

$\hat{x}_1$ is not correlated with $(\beta_1 \hat{\varepsilon} + u)$

This in turn means that $x_2$, $z_1$, and $z_2$ are not correlated with $u$ (the error term of the true model.

($\hat{x}_1$ is always not correlated (orthogonal) with $\varepsilon$)

**Condition 2**

- $z_1$ and $z_2$ do not belong in the main model, meaning they do not have any explanatory power beyond $x_2$ (they should have been included in the model in the first place as independent variables)

- $z_1$ and $z_2$ are not correlated with the error term (there are no unobserved factors in the error term that are correlated with $Z$)

**Question**

Do you think we can test condition 2?

**Question**

Do you think we can test condition 2?

**Answer**

No, because we never observe the error term.

**Question**

Do you think we can test condition 2?

**Answer**

No, because we never observe the error term.

**Important**

- All we can do is to argue that the instruments are not correlated with the error term.

**Question**

Do you think we can test condition 2?

**Answer**

No, because we never observe the error term.

**Important**

- All we can do is to argue that the instruments are not correlated with the error term.

- In journal articles that use IV method, they make careful arguments as to why their choice of instruments are not correlated with the error term.

**Condition 1**

- The instrument(s) $Z$ have jointly significant explanatory power on the endogenous variable $x_1$ after you control for all the other exogenous variables (here $x_2$)}

**Condition 2**

- $z_1$ and $z_2$ do not belong in the main model, meaning they do not have any explanatory power beyond $x_2$ (they should have been included in the model in the first place as independent variables)

- $z_1$ and $z_2$ are not correlated with the error term (there are no unobserved factors in the error term that are correlated with $Z$)

**Condition 1**

- The instrument(s) $Z$ have jointly significant explanatory power on the endogenous variable $x_1$ after you control for all the other exogenous variables (here $x_2$)}

**Condition 2**

- $z_1$ and $z_2$ do not belong in the main model, meaning they do not have any explanatory power beyond $x_2$ (they should have been included in the model in the first place as independent variables)

- $z_1$ and $z_2$ are not correlated with the error term (there are no unobserved factors in the error term that are correlated with $Z$)

**Important**

- Condition 1 is always testable

- Condition 2 is NOT testable (unless you have more instruments than endogenous variables)

**Two-stage Least Square (2SLS)**

IV estimator is also called two-stage least squares estimator (2SLS) because it involves two stages of OLS.

- Step 1: $\hat{x}_1 = \hat{\alpha}_0 + \hat{\sigma}_2 x_2 + \hat{\alpha}_1 z_1 + \hat{\alpha}_2 z_2$
- Step 2: $y = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 x_2 + \varepsilon$

**Two-stage Least Square (2SLS)**

IV estimator is also called two-stage least squares estimator (2SLS) because it involves two stages of OLS.

- Step 1: $\hat{x}_1 = \hat{\alpha}_0 + \hat{\sigma}_2 x_2 + \hat{\alpha}_1 z_1 + \hat{\alpha}_2 z_2$

- Step 2: $y = \beta_0 + \beta_1 \hat{x}_1 + \beta_2 x_2 + \varepsilon$

- 2SLS framework is a good way to understand conceptually why and how instrumental variable estimation works

- But, IV estimation is done in one-step

# Instrumental variable validity

**The model**

$$log(wage) = \beta_0 + \beta_1 educ + \beta_2 exper + v \ (= \beta_3 ability + u)$$

`educ` is endogenous because of its correlation with `ability`.

**The model**

$$log(wage) = \beta_0 + \beta_1 educ + \beta_2 exper + v \ \ (= \beta_3 ability + u)$$

`educ` is endogenous because of its correlation with `ability`.

**Question**

What conditions would a good instrument $(z)$ satisfy?

**The model**

$$log(wage) = \beta_0 + \beta_1 educ + \beta_2 exper + v \;\; (= \beta_3 ability + u)$$

`educ` is endogenous because of its correlation with `ability`.

**Question**

What conditions would a good instrument $(z)$ satisfy?

**Answer**

- $z$ has explanatory power on $educ$ after you control for the impact of $epxer$ on $educ$

- $z$ is uncorrelated with $v$ $(ability$ and all the other important unobservables)

**The model**

$$log(wage) = \beta_0 + \beta_1 educ + \beta_2 exper + v \ \ (= \beta_3 ability + u)$$

**An example of instruments**

The last digit of an individual's Social Security Number? (this has been actually used in some journal articles)

**Question**

- Is it uncorrelated with $v$ ($ability$ and all the other important unobservables)?

**The model**

$$log(wage) = \beta_0 + \beta_1 educ + \beta_2 exper + v \ \ (= \beta_3 ability + u)$$

**An example of instruments**

The last digit of an individual's Social Security Number? (this has been actually used in some journal articles)

**Question**

- Is it uncorrelated with $v$ ($ability$ and all the other important unobservables)?

- does it have explanatory power on $educ$ after you control for the impact of $epxer$ on $educ$?

**The model**

$$log(wage) = \beta_0 + \beta_1 educ + \beta_2 exper + v \ (= \beta_3 ability + u)$$

**An example of instruments**

IQ test score

**Question**

- Is it uncorrelated with $v$ ($ability$ and all the other important unobservables)?

**The model**

$$log(wage) = \beta_0 + \beta_1 educ + \beta_2 exper + v \;\; (= \beta_3 ability + u)$$

**An example of instruments**

IQ test score

**Question**

- Is it uncorrelated with $v$ ($ability$ and all the other important unobservables)?

- does it have explanatory power on $educ$ after you control for the impact of $epxer$ on $educ$?

**The model**

$$log(wage) = \beta_0 + \beta_1 educ + \beta_2 exper + v \ \ (= \beta_3 ability + u)$$

**An example of instruments**

Mother's education

**Question**

- Is it uncorrelated with $v$ ($ability$ and all the other important unobservables)?

**The model**

$$log(wage) = \beta_0 + \beta_1 educ + \beta_2 exper + v \ (= \beta_3 ability + u)$$

**An example of instruments**

Mother's education

**Question**

- Is it uncorrelated with $v$ ($ability$ and all the other important unobservables)?

- does it have explanatory power on $educ$ after you control for the impact of $epxer$ on $educ$?

**The model**

$$log(wage) = \beta_0 + \beta_1 educ + \beta_2 exper + v \ (= \beta_3 ability + u)$$

**An example of instruments**

Number of siblings

**Question**

- Is it uncorrelated with $v$ ($ability$ and all the other important unobservables)?

**The model**

$$log(wage) = \beta_0 + \beta_1 educ + \beta_2 exper + v \ \ (= \beta_3 ability + u)$$

**An example of instruments**

Number of siblings

**Question**

- Is it uncorrelated with $v$ ($ability$ and all the other important unobservables)?

- does it have explanatory power on $educ$ after you control for the impact of $epxer$ on $educ$?

# Implementation of Instrumental Variable (IV) Estimation in R

**Model**

$$log(wage) = \beta_0 + \beta_1 educ + \beta_2 exper + v \ \ (= \beta_3 ability + u)$$

We believe

- $educ$ is endogenous $(x_1)$

- $exper$ is exogenous $(x_2)$

- we use the number of siblings $(sibs)$ and father's education $(feduc)$ as the instruments ($Z$)

**Model**

$$log(wage) = \beta_0 + \beta_1 educ + \beta_2 exper + v \ \ (= \beta_3 ability + u)$$

We believe

- $educ$ is endogenous $(x_1)$

- $exper$ is exogenous $(x_2)$

- we use the number of siblings $(sibs)$ and father's education $(feduc)$ as the instruments ($Z$)


**Terminology**

- exogenous variable included in the model (here, $exper$) is also called included instruments

- instruments that do not belong to the main model (here, $sibs$ and $feduc$) are also called excluded instruments

- we refer to the collection of included and excluded instruments as instruments

## Dataset

```r
#--- take a look at the data ---#
wage2 %>%
  select(wage, educ, sibs, feduc) %>%
  head()
```

```
##    wage educ sibs feduc
## 1  769   12    1     8
## 2  808   18    1    14
## 3  825   14    1    14
## 4  650   12    4    12
## 5  562   11   10    11
## 6 1400   16    1    NA
```

We can continue to use the `fixest` package to run IV estimation method.

```
library(fixest)
```

```
felm(dep var ~ included instruments|first stage formula, data = dataset)
```

- `included instruments`: exogenous included variables (do not include endogenous variables here)

We can continue to use the `fixest` package to run IV estimation method.

```
library(fixest)
```

**Syntax**

```
felm(dep var ~ included instruments|first stage formula, data = dataset)
```

- `included instruments`: exogenous included variables (do not include endogenous variables here)

**first stage formula**

```
(endogenous vars ~ excluded instruments)
```

We can continue to use the `fixest` package to run IV estimation method.

```
library(fixest)
```

**Syntax**

```
felm(dep var ~ included instruments|first stage formula, data = dataset)
```

- `included instruments`: exogenous included variables (do not include endogenous variables here)

**first stage formula**

```
(endogenous vars ~ excluded instruments)
```

**Example**

```
iv_res <- feols(log(wage) ~ exper | educ ~ sibs + feduc, data = wage2)
```

- `included variable`:
  - exogenous included variables: `exper`
  - endogenous included variables: `educ`
- `instruments`:
  - included instruments: `exper`
  - excluded instruments: `sibs` and `feduc`

## IV regression results

```
iv_res
```

```
## TSLS estimation, Dep. Var.: log(wage), Endo.: educ, Instr.: sibs, feduc
## Second stage: Dep. Var.: log(wage)
## Observations: 741
## Standard-errors: IID
##              Estimate Std. Error  t value   Pr(>|t|)
## (Intercept) 4.507316   0.315735 14.27564  < 2.2e-16 ***
## fit_educ    0.137405   0.019215  7.15104 2.0766e-12 ***
## exper       0.037029   0.005694  6.50306 1.4502e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## RMSE: 0.406208   Adj. R2: 0.049979
## F-test (1st stage), educ: stat = 65.6      , p < 2.2e-16 , on 2 and 737 DoF.
##                Wu-Hausman: stat = 13.2      , p = 3.051e-4, on 1 and 737 DoF.
##                    Sargan: stat =  0.230925, p = 0.630838, on 1 DoF.
```

## Note

- When variable `x` is the endogenous variable, `fixest` changes the name of `x` to `x(fit)`.

- Here, `educ` has become `educ(fit)`.

**Comparison of OLS and IV Estimation Results**

|              | Model 1   | Model 2   |
| ------------ | --------- | --------- |
| (Intercept)  | 5.503***  | 4.507***  |
|              | (0.112)   | (0.316)   |
| educ         | 0.078***  |           |
|              | (0.007)   |           |
| exper        | 0.020***  | 0.037***  |
|              | (0.003)   | (0.006)   |
| fit_educ     |           | 0.137***  |
|              |           | (0.019)   |
| Num.Obs.     | 935       | 741       |
| R2           | 0.131     | 0.053     |
| Std. errors  | IID       | IID       |

$* p < 0.1$, $** p < 0.05$, $*** p < 0.01$

|               | Model 1   | Model 2   |
| ------------- | --------- | --------- |
| (Intercept)   | 5.503***  | 4.507***  |
|               | (0.112)   | (0.316)   |
| educ          | 0.078***  |           |
|               | (0.007)   |           |
| exper         | 0.020***  | 0.037***  |
|               | (0.003)   | (0.006)   |
| fit_educ      |           | 0.137***  |
|               |           | (0.019)   |
| Num.Obs.      | 935       | 741       |
| R2            | 0.131     | 0.053     |
| Std. errors   | IID       | IID       |

* p < 0.1, ** p < 0.05, *** p < 0.01

**Question**

Do you think $sibs$ and $feduc$ are good instruments?

- Condition 1: weak instruments?
- Condition 2: uncorrelated with the error term?

**Weak Instrument Test**

We can always test if the excluded instruments are weak or not (test of condition 1).

**Weak Instrument Test**

We can always test if the excluded instruments are weak or not (test of condition 1).

**How**

- Run the 1st stage regression

$$educ = \alpha_0 + \alpha_1 exper + \alpha_2 sibs + \alpha_3 feduc + v$$

**Weak Instrument Test**

We can always test if the excluded instruments are weak or not (test of condition 1).

**How**

- Run the 1st stage regression

$$educ = \alpha_0 + \alpha_1 exper + \alpha_2 sibs + \alpha_3 feduc + v$$

- test the joint significance of $\alpha_2$ and $\alpha_3$ ($F$-test)

If excluded instruments ($sibs$ and $feduc$) are jointly significant, then it would mean that $sibs$ and $feduc$ are not weak instruments, satisfying condition 1.

When we ran the IV estimation using `fixest::feols()` earlier, it automatically calculated the F-statistic for the weak instrument test.

When we ran the IV estimation using `fixest::feols()` earlier, it automatically calculated the F-statistic for the weak instrument test.

```
iv_res
```

```
## TSLS estimation, Dep. Var.: log(wage), Endo.: educ, Instr.: sibs, feduc
## Second stage: Dep. Var.: log(wage)
## Observations: 741
## Standard-errors: IID
##              Estimate Std. Error  t value   Pr(>|t|)
## (Intercept) 4.507316   0.315735 14.27564  < 2.2e-16 ***
## fit_educ    0.137405   0.019215  7.15104 2.0766e-12 ***
## exper       0.037029   0.005694  6.50306 1.4502e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## RMSE: 0.406208   Adj. R2: 0.049979
## F-test (1st stage), educ: stat = 65.6     , p < 2.2e-16 , on 2 and 737 DoF.
##              Wu-Hausman: stat = 13.2     , p = 3.051e-4, on 1 and 737 DoF.
##                  Sargan: stat =  0.230925, p = 0.630838, on 1 DoF.
```

Here, F-test for the null hypothesis of the excluded instruments (`sibs` and `feduc`) do not have any explanatory power on the endogenous variable (`educ`) beyond the included instrument (`exper`) is rejected.

Alternatively, you can access the `iv_first_stage` component of the regression results.

```
iv_res$iv_first_stage
```

```
## $educ
## TSLS estimation, Dep. Var.: educ, Endo.: educ, Instr.: sibs, feduc
## First stage: Dep. Var.: educ
## Observations: 741
## Standard-errors: IID
##              Estimate Std. Error   t value    Pr(>|t|)
## (Intercept) 14.075273   0.358595  39.25116  < 2.2e-16 ***
## sibs        -0.131009   0.030800  -4.25357 2.3749e-05 ***
## feduc        0.205169   0.021909   9.36459  < 2.2e-16 ***
## exper       -0.191535   0.016373 -11.69819  < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## RMSE: 1.84505   Adj. R2: 0.319802
## F-test (1st stage): stat = 65.6, p < 2.2e-16, on 2 and 737 DoF.
```

**Notes**

- It is generally recommended that you have $F$-stat of over $10$ (this is not a clear-cut criteria that applied to all the empirical cases)

- Even if you reject the null if $F$-stat is small, you may have a problem

- You know nothing about if your excluded instruments satisfy Condition 2.

- If you cannot reject the null, it is a strong indication that your instruments are weak. Look for other instruments.

- Always, always report this test. There is no reason not to.

# Consequences of weak instruments

**Data generation**

```r
set.seed(73289)
N <- 500 # number of observations

u_common <- runif(N) # the term shared by the endogenous variable and the error term
z_common <- runif(N) # the term shared by the endogenous variable and instruments
x_end <- u_common + z_common + runif(N) # the endogenous variable
z_strong <- z_common + runif(N) # strong instrument
z_weak <- 0.01 * z_common + 0.99995 * runif(N) # weak instrument
u <- u_common + runif(N) # error term
y <- x_end + u # dependent variable

data <- data.frame(y, x_end, z_strong, z_weak)
```

**Correlation**

```
cor(data)
```

```
##                    y         x_end    z_strong        z_weak
## y         1.0000000   0.86492868  0.298704509  -0.108007146
## x_end     0.8649287   1.00000000  0.419011491  -0.074224622
## z_strong  0.2987045   0.41901149  1.000000000   0.003839565
## z_weak   -0.1080071  -0.07422462  0.003839565   1.000000000
```

**Estimation with the strong instrumental variable**

```r
#--- IV estimation (strong) ---#
iv_strong <- feols(y ~ 1 | x_end ~ z_strong, data = data)
```

**Estimation with the weak instrumental variable**

```r
#--- IV estimation (weak) ---#
iv_weak <- feols(y ~ 1 | x_end ~ z_weak, data = data)
```

```
#--- coefs (strong) ---#
tidy(iv_strong)
```

```
## # A tibble: 2 × 5
##   term          estimate std.error statistic  p.value
##   <chr>            <dbl>     <dbl>     <dbl>     <dbl>
## 1 (Intercept)      0.883     0.133      6.64 8.20e-11
## 2 fit_x_end        1.09      0.0856    12.7  2.96e-32
```

```
#--- coefs (weak) ---#
tidy(iv_weak)
```

```
## # A tibble: 2 × 5
##   term          estimate std.error statistic p.value
##   <chr>            <dbl>     <dbl>     <dbl>    <dbl>
## 1 (Intercept)     -0.862     1.10     -0.784 0.434
## 2 fit_x_end        2.22      0.714     3.11  0.00197
```

**Question**

Any notable differences?

```
#--- coefs (strong) ---#
tidy(iv_strong)
```

```
## # A tibble: 2 × 5
##   term        estimate std.error statistic  p.value
##   <chr>          <dbl>     <dbl>     <dbl>    <dbl>
## 1 (Intercept)    0.883    0.133      6.64 8.20e-11
## 2 fit_x_end      1.09     0.0856    12.7  2.96e-32
```

```
#--- coefs (weak) ---#
tidy(iv_weak)
```

```
## # A tibble: 2 × 5
##   term        estimate std.error statistic p.value
##   <chr>          <dbl>     <dbl>     <dbl>   <dbl>
## 1 (Intercept)   -0.862    1.10     -0.784 0.434
## 2 fit_x_end      2.22     0.714     3.11  0.00197
```

**Question**

Any notable differences?

The coefficient estimate on $x\_end$ is far away from the true value in the weak instrument case.

## Comparison of the weak instrument tests

```
#--- diagnostics (strong) ---#
iv_strong$iv_first_stage
```

```
## $x_end
## TSLS estimation, Dep. Var.: x_end, Endo.: x_end, Instr.: z_strong
## First stage: Dep. Var.: x_end
## Observations: 500
## Standard errors: IID
```

```
#--- diagnostics (weak) ---#
iv_weak$iv_first_stage
```

```
## $x_end
## TSLS estimation, Dep. Var.: x_end, Endo.: x_end, Instr.: z_weak
## First stage: Dep. Var.: x_end
## Observations: 500
## Standard errors: IID
```

## Question

Any notable differences?

**Comparison of the weak instrument tests**

```
#--- diagnostics (strong) ---#
iv_strong$iv_first_stage
```

```
## $x_end
## TSLS estimation, Dep. Var.: x_end, Endo.: x_end, Instr.: z_strong
## First stage: Dep. Var.: x_end
## Observations: 500
## Standard errors: IID
```

```
#--- diagnostics (weak) ---#
iv_weak$iv_first_stage
```

```
## $x_end
## TSLS estimation, Dep. Var.: x_end, Endo.: x_end, Instr.: z_weak
## First stage: Dep. Var.: x_end
## Observations: 500
## Standard errors: IID
```

**Question**

Any notable differences?

You cannot reject the null hypothesis of weak instrument in the weak instrument case.

## MC simulation

```r
B <- 1000 # the number of experiments
beta_hat_store <- matrix(0, B, 2) # storage of beta hat

for (i in 1:B) {

  #--- data generation ---#
  u_common <- runif(N)
  z_common <- runif(N)
  x_end <- u_common + z_common + runif(N)
  z_strong <- z_common + runif(N)
  z_weak <- 0.01 * z_common + 0.99995 * runif(N)
  u <- u_common + runif(N)
  y <- x_end + u
  data <- data.table(y, x_end, z_strong, z_weak)

  #--- IV estimation with a strong instrument ---#
  iv_strong <- feols(y ~ 1 | x_end ~ z_strong, data = data)
  beta_hat_store[i, 1] <- iv_strong$coefficients[2]

  #--- IV estimation with a weak instrument ---#
  iv_weak <- feols(y ~ 1 | x_end ~ z_weak, data = data)
  beta_hat_store[i, 2] <- iv_weak$coefficients[2]
}
```
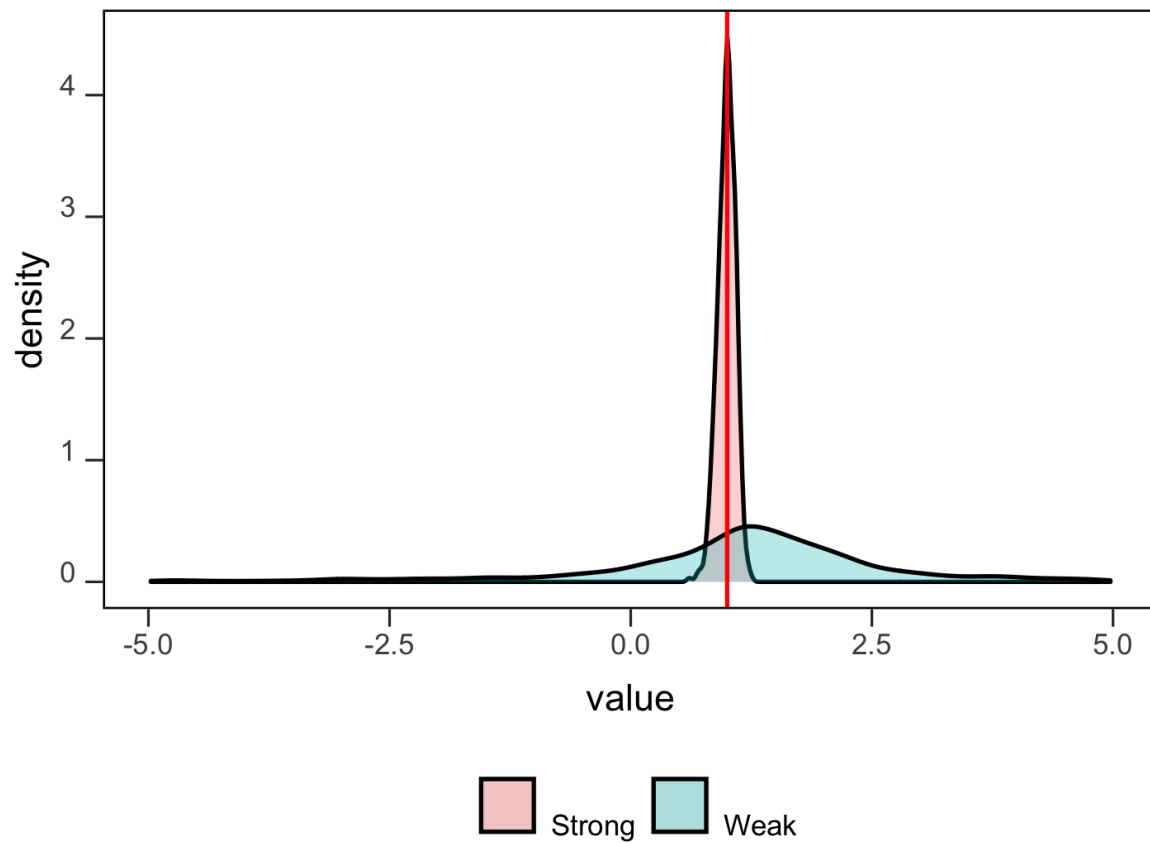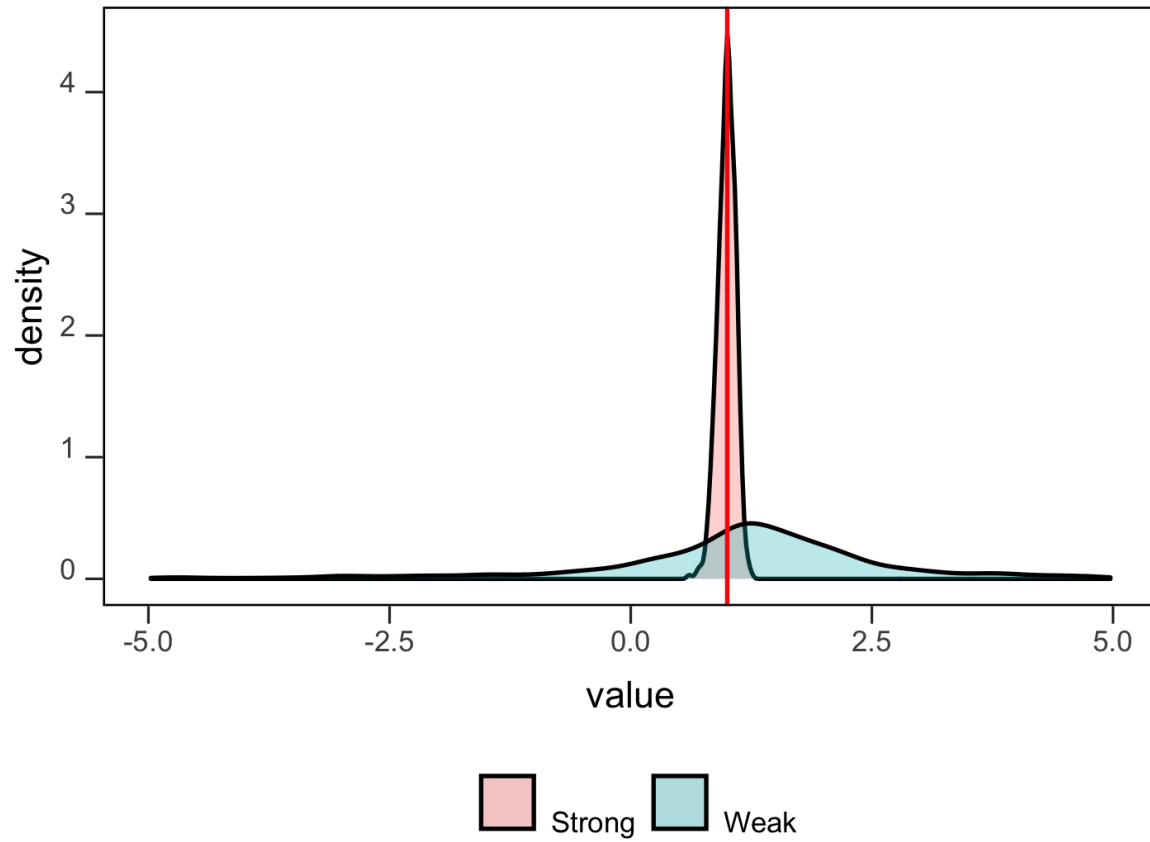
# Visualization of the MC Results

**Flow of IV Estimation in Practice**

- Identify endogenous variable(s) and included instrument(s)

- Identify potential excluded instrument(s)

- Argue why the excluded instrument(s) you pick is uncorrelated with the error term ( **condition 2** )

- Once you decide what variable(s) to use as excluded instruments, test whether the excluded instrument(s) is weak or not ( **condition 1** )

- Implement IV estimation and report the results

You can include fixed effects in your IV estimation.

**Syntax**

```
feols(dep var ~ included instruments | FE | 1st stage formula, data = dataset)
```

**Example**

Include `married` and `south` as fixed effects.

```
feols(log(wage) ~ exper | married + south | educ ~ feduc + sibs, data = wage2)
```

```
## TSLS estimation, Dep. Var.: log(wage), Endo.: educ, Instr.: feduc, sibs
## Second stage: Dep. Var.: log(wage)
## Observations: 741
## Fixed-effects: married: 2,  south: 2
## Standard-errors: Clustered (married)
##           Estimate Std. Error t value Pr(>|t|)
## fit_educ 0.124355   0.003627 34.2906 0.018560 *
## exper    0.032128   0.002260 14.2144 0.044713 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## RMSE: 0.391178     Adj. R2: 0.116588
##                  Within R2: 0.069595
## F-test (1st stage), educ: stat = 61.1      , p < 2.2e-16 , on 2 and 736 DoF.
##              Wu-Hausman: stat =  8.98498 , p = 0.002814, on 1 and 735 DoF.
##                  Sargan: stat =  0.169226, p = 0.6808   , on 1 DoF.
```

Clustered SE? You can just add `cluster =` option just like we previously did.

```
feols(log(wage) ~ exper | married + south | educ ~ feduc + sibs, cluster = ~black, data = wage2)
```

```
## TSLS estimation, Dep. Var.: log(wage), Endo.: educ, Instr.: feduc, sibs
## Second stage: Dep. Var.: log(wage)
## Observations: 741
## Fixed-effects: married: 2,  south: 2
## Standard-errors: Clustered (black)
##           Estimate Std. Error t value Pr(>|t|)
## fit_educ 0.124355   0.005258 23.6526 0.026899 *
## exper    0.032128   0.002798 11.4842 0.055295 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## RMSE: 0.391178     Adj. R2: 0.116588
##                  Within R2: 0.069595
## F-test (1st stage), educ: stat = 61.9      , p < 2.2e-16 , on 2 and 735 DoF.
##               Wu-Hausman: stat =  8.98498 , p = 0.002814, on 1 and 735 DoF.
##                   Sargan: stat =  0.169226, p = 0.6808   , on 1 DoF.
```