

# 08 - Causal Inference

ml4econ, HUI 2020

Itamar Caspi

May 18, 2019 (updated: 2020-05-17)

# Replicating this presentation

Use the **pacman** package to install and load packages:

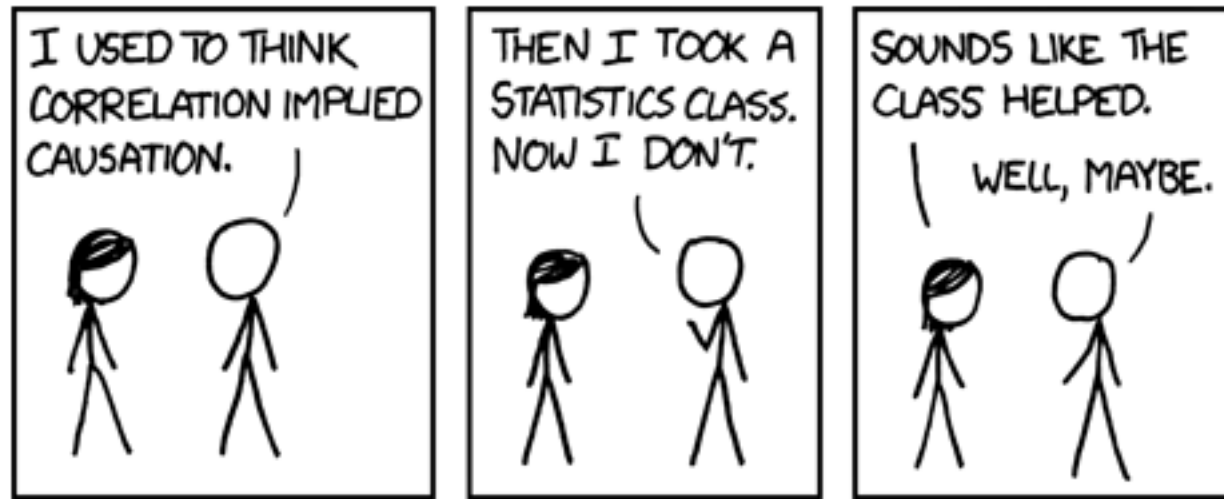
```
if (!require("pacman"))  
  install.packages("pacman")  
  
pacman::p_load(  
  tidyverse,    # for data wrangling and visualization  
  tidymodels,   # for modeling  
  haven,        # for reading dta files  
  here,         # for referencing folders  
  dagitty,      # for generating DAGs  
  ggdag,        # for drawing DAGs  
  knitr         # for printing html tables  
)
```

# Outline

- Causal Inference
- Potential Outcomes
- Directed Acyclic Graphs

# Causal Inference

# Predicting vs. explaining



Source: XKCD

# Looking forward

- Until now, our focus was on prediction.
- However, what we economists mostly care about is causal inference:
  - What is the effect of class size on student performance?
  - What is the effect of education on earnings?
  - What is the effect of government spending on GDP?
  - etc.
- Before we learn how to adjust and apply ML methods to causal inference problems, we need to be explicit about what causal inference is.
- This lecture will review two dominant approaches to causal inference, the statistical/econometric approach and the computer science approach.

# A note on identification

- The primary focus of this lecture is on identification, as opposed to prediction, estimation and inference.
- In short, identification is defined as

*"model parameters or features being uniquely determined from the observable population that generates the data."* - (Lewbel, 2019)

- More specifically, think about identifying the parameter of interest when you have unlimited data (the entire population).

# Potential Outcomes



# Pearl and Rubin



**Source:** The Book of Why (Pearl and Mackenzie)

# Rubin and potential outcomes



**Source:** The Book of Why (Pearl and Mackenzie)

# The road not taken



Source: <https://mru.org/courses/mastering-econometrics/ceteris-paribus>

# Notation

- $Y$  is a random variable
- $X$  is a vector of attributes
- $\mathbf{X}$  is a design matrix

# Treatment and potential outcomes (Rubin, 1974, 1977)

- Treatment

$$D_i = \begin{cases} 1, & \text{if unit } i \text{ received the treatment} \\ 0, & \text{otherwise.} \end{cases}$$

- Treatment and potential outcomes

$Y_{i0}$  is the potential outcome for unit  $i$  with  $D_i = 0$

$Y_{i1}$  is the potential outcome for unit  $i$  with  $D_i = 1$

- Observed outcome: Under the Stable Unit Treatment Value Assumption (SUTVA), The realization of unit  $i$ 's outcome is

$$Y_i = Y_{1i}D_i + Y_{0i}(1 - D_i)$$

**Fundamental problem of causal inference** (Holland, 1986): We cannot observe *both*  $Y_{1i}$  and  $Y_{0i}$ .

# Treatment effect and observed outcomes

- Individual treatment effect: The difference between unit  $i$ 's potential outcomes:

$$\tau_i = Y_{1i} - Y_{0i}$$

- *Average treatment effect* (ATE)

$$\mathbb{E}[\tau_i] = \mathbb{E}[Y_{1i} - Y_{0i}] = \mathbb{E}[Y_{1i}] - \mathbb{E}[Y_{0i}]$$

- *Average treatment effect for the treatment group* (ATT)

$$\mathbb{E}[\tau_i | D_i = 1] = \mathbb{E}[Y_{1i} - Y_{0i} | D_i = 1] = \mathbb{E}[Y_{1i} | D_i = 1] - \mathbb{E}[Y_{0i} | D_i = 1]$$

**NOTE:** The complement of the treatment group is the *control* group.

# Selection bias

A naive estimand for ATE is the difference between average outcomes based on treatment status

However, this might be misleading:

$$\mathbb{E}[Y_i|D_i = 1] - \mathbb{E}[Y_i|D_i = 0] = \underbrace{\mathbb{E}[Y_{1i}|D_i = 1] - \mathbb{E}[Y_{0i}|D_i = 1]}_{\text{ATT}} + \underbrace{\mathbb{E}[Y_{0i}|D_i = 1] - \mathbb{E}[Y_{0i}|D_i = 0]}_{\text{selection bias}}$$

**Causal inference is mostly about eliminating selection-bias**

**EXAMPLE:** Individuals who go to private universities probably have different characteristics than those who go to public universities.

# Randomized control trial (RCT) solves selection bias

In an RCT, the treatments are randomly assigned. This means entails that  $D_i$  is *independent* of potential outcomes, namely

$$\{Y_{1i}, Y_{0i}\} \perp D_i$$

RCTs enables us to estimate ATE using the average difference in outcomes by treatment status:

$$\begin{aligned}\mathbb{E}[Y_i | D_i = 1] - \mathbb{E}[Y_i | D_i = 0] &= \mathbb{E}[Y_{1i} | D_i = 1] - \mathbb{E}[Y_{0i} | D_i = 0] \\ &= \mathbb{E}[Y_{1i} | D_i = 1] - \mathbb{E}[Y_{0i} | D_i = 1] \\ &= \mathbb{E}[Y_{1i} - Y_{0i} | D_i = 1] \\ &= \mathbb{E}[Y_{1i} - Y_{0i}] \\ &= \text{ATE}\end{aligned}$$

**EXAMPLE:** In theory, randomly assigning students to private and public universities would allow us to estimate the ATE going to private school have on future earnings. Clearly, RCT in this case is infeasible.



# Estimands and regression

Assume for now that the treatment effect is constant across all individuals, i.e.,

$$\tau = Y_{1i} - Y_{0i}, \quad \forall i.$$

Accordingly, we can express  $Y_i$  as

$$\begin{aligned} Y_i &= Y_{1i}D_i + Y_{0i}(1 - D_i) \\ &= Y_{0i} + D_i(Y_{1i} - Y_{0i}), \\ &= Y_{0i} + \tau D_i, && \text{since } \tau = Y_{1i} - Y_{0i} \\ &= \mathbb{E}[Y_{0i}] + \tau D_i + Y_{0i} - \mathbb{E}[Y_{0i}], && \text{add and subtract } \mathbb{E}[Y_{0i}] \end{aligned}$$

Or more conveniently

$$Y_i = \alpha + \tau D_i + u_i,$$

where  $\alpha = \mathbb{E}[Y_{0i}]$  and  $u_i = Y_{0i} - \mathbb{E}[Y_{0i}]$  is the random component of  $Y_{0i}$ .

# Unconfoundedness

Typically, in observational studies, treatments are not randomly assigned. (Think of  $D_i = \{\text{private}, \text{public}\}$ .)

In this case, identifying causal effects depended on the *Unconfoundedness* assumption (also known as "selection-on-observable"), which is defined as

$$\{Y_{1i}, Y_{0i}\} \perp D_i | X_i$$

In words: treatment assignment is independent of potential outcomes *conditional* on observed  $X_i$ , i.e., selection bias *disappears* when we control for  $X_I$ .

# Adjusting for confounding factors

The most common approach for controlling for  $X_i$  is by adding them to the regression:

$$Y_i = \alpha + \tau D_i + X_i' \beta + u_i,$$

## COMMENTS:

1. Strictly speaking, the above regression model is valid if we actually *believe* that the "true" model is  $Y_i = \alpha + \tau D_i + X_i' \beta + u_i$ .
2. If  $D_i$  is randomly assigned, adding  $X_i$  to the regression **might** increase the accuracy of ATE.
3. If  $D_i$  is assigned conditional on  $X_i$  (e.g., in observational settings), adding  $X_i$  to the regression eliminates selection bias.

# Illustration: the OHIE data

- The Oregon Health Insurance Experiment (OHIE), is a randomized controlled trial for measuring the treatment effect of Medicaid eligibility.
- Treatment group: Those selected in the Medicaid lottery.
- The outcome, `doc_any_12m`, equals to 1 for patients who saw a primary care physician, and zero otherwise.

# Load the OHIE data

```
descr <-  
  here("08-causal-inference/data",  
        "oregonhie_descriptive_vars.dta") %>%  
  read_dta()  
  
prgm <-  
  here("08-causal-inference/data",  
        "oregonhie_stateprograms_vars.dta") %>%  
  read_dta()  
  
s12 <-  
  here("08-causal-inference/data",  
        "oregonhie_survey12m_vars.dta") %>%  
  read_dta()
```

The entire OHIE data can be found [here](#).

# Preprocessing: Joining datasets

```
ohie_raw <-  
  descr %>%  
  left_join(prgm) %>%  
  left_join(s12) %>%  
  filter(sample_12m_resp == 1) %>%  
  drop_na(  
    weight_12m, doc_any_12m, doc_num_mod_12m,  
    er_any_12m, er_num_mod_12m,  
    hosp_any_12m, hosp_num_mod_12m  
  )
```

# Preprocessing: Refinement

```
ohie <-  
  ohie_raw %>%  
  rename(  
    enrolled = ohp_all_ever_firstn_30sep2009,  
    selected = treatment  
  ) %>%  
  dplyr::select(  
    person_id, household_id,  
    numhh_list, selected,  
    enrolled, doc_any_12m  
  ) %>%  
  mutate(numhh_list = factor(numhh_list, levels = c("1", "2", "3")))
```

# The final dataset

ohie

```
## # A tibble: 23,107 x 6
##   person_id household_id numhh_list      selected      enrolled doc_any_12m
##   <dbl>         <dbl> <fct>         <dbl+lbl>    <dbl+lbl>    <dbl+lbl>
## 1           1           100001 1             1 [Selected]  0 [NOT enro~    0 [No]
## 2           2           100002 1             1 [Selected]  1 [Enrolled]    0 [No]
## 3           5           100005 1             1 [Selected]  0 [NOT enro~    0 [No]
## 4           6           100006 1             1 [Selected]  0 [NOT enro~    1 [Yes]
## 5           8          102094 2             0 [Not selec~  0 [NOT enro~    0 [No]
## 6           9          100009 1             0 [Not selec~  0 [NOT enro~    1 [Yes]
## 7          10          111771 2             0 [Not selec~  0 [NOT enro~    1 [Yes]
## 8          13          100013 1             1 [Selected]  1 [Enrolled]    1 [Yes]
## 9          14          100014 1             1 [Selected]  1 [Enrolled]    0 [No]
## 10         23          115253 2             1 [Selected]  1 [Enrolled]    1 [Yes]
## # ... with 23,097 more rows
```



# Distribution of treated-control

```
ohie %>%  
  count(selected) %>%  
  kable(format = "html")
```

selected	n
0	11629
1	11478

# Estimating ATE

The estimated model

$$doc\_any\_12m_i = \alpha + \tau \times selected_i + \varepsilon_i$$

In R:

```
fit <- lm(doc_any_12m ~ selected, data = ohie)
```

# Results

```
fit %>%  
  tidy(conf.int = TRUE) %>%  
  filter(term != "(Intercept)") %>%  
  dplyr::select(term, estimate, starts_with("conf.")) %>%  
  kable(digits = 4, format = "html")
```

term	estimate	conf.low	conf.high
selected	0.0575	0.0449	0.0701

**Interpretation:** being selected in the lottery increases the probability that you visit primary care physician in the following year by 5.75 [4.49, 7.01] percentage points.

# Adjustments

One issue with OHIE is that people are able to apply for Medicaid for their entire household.

This fact undermines the critical random assignment assumption since belonging to larger households increases the chances of being selected to Medicaide.

```
ohie %>%  
  count(selected, numhh_list) %>%  
  kable(format = "html")
```

selected	numhh_list	n
0	1	8684
0	2	2939
0	3	6
1	1	7525
1	2	3902
1	3	51

# ATE under adjustment

```
lm(doc_any_12m ~ selected + numhh_list, data = ohie) %>%  
  tidy(conf.int = TRUE) %>%  
  dplyr::select(term, estimate, starts_with("conf.")) %>%  
  kable(digits = 4, format = "html")
```

term	estimate	conf.low	conf.high
(Intercept)	0.5902	0.5807	0.5997
selected	0.0639	0.0512	0.0765
numhh_list2	-0.0657	-0.0796	-0.0519
numhh_list3	-0.1737	-0.3006	-0.0467

After adjusting for numhh\_list, ATE has increased from 0.057 to 0.064. Can you guess why?

# Directed Acyclic Graphs

# Pearl and DAGs



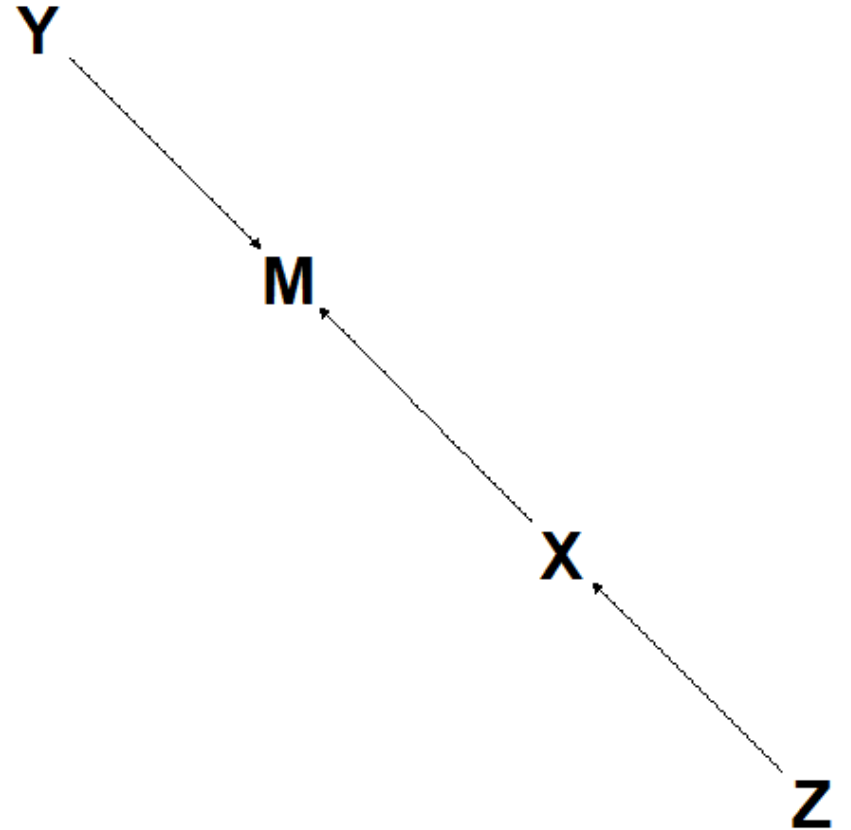
**Source:** The Book of Why (Pearl and Mackenzie)

---

# DAGs (Pearl)

A DAG is a way to model a system of causal interactions using graphs.

- **Nodes** represents random variables, e.g.,  $X$ ,  $Y$ , etc.
- **Arrows** (or directed edges) represent "from  $\rightarrow$  to" causal effects. For example,  $Z \rightarrow X$  reads " $Z$  causes  $X$ ".
- A **path** is a sequence of edges connecting two nodes. For example,  $Z \rightarrow X \rightarrow M \leftarrow Y$  describes a path from  $Z$  to  $Y$ .
  - In a **direct path** arrows point to the same direction:  $Z \rightarrow X \rightarrow M$





# DAGs and SEM

- Another way to think about DAGs is as non-parametric **structural equation models** (SEM)
- For example, the single-confounder DAG,  $D \leftarrow X \rightarrow Y$ , can be represented by a set of three equations:

$$X \leftarrow f_X(U_X)$$

$$D \leftarrow f_D(X, U_D)$$

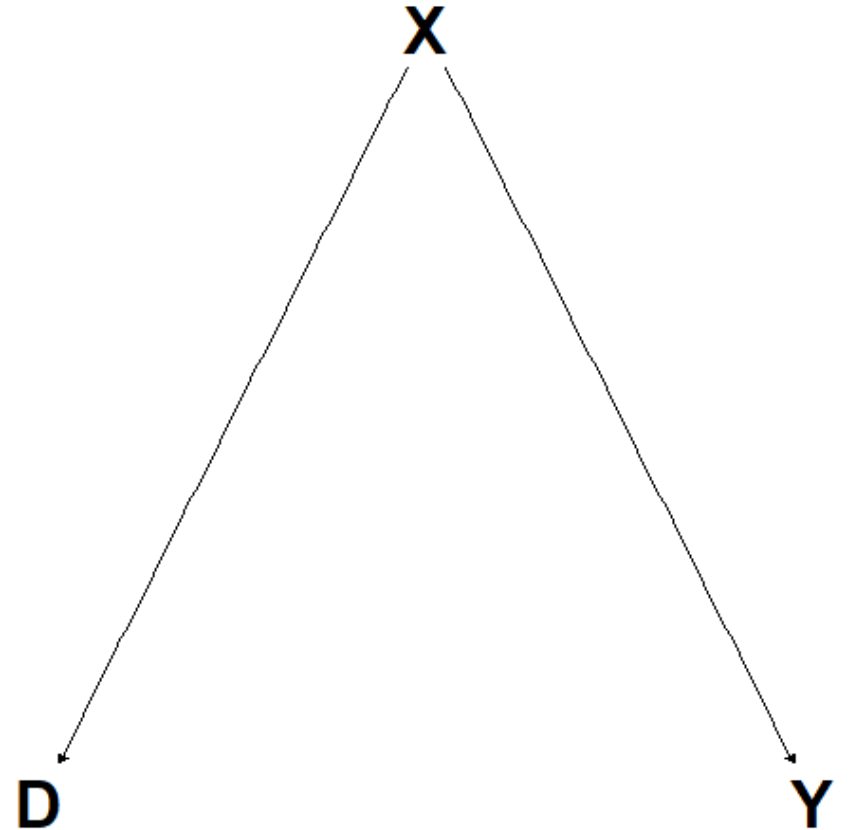
$$Y \leftarrow f_Y(D, X, U_Y)$$

where

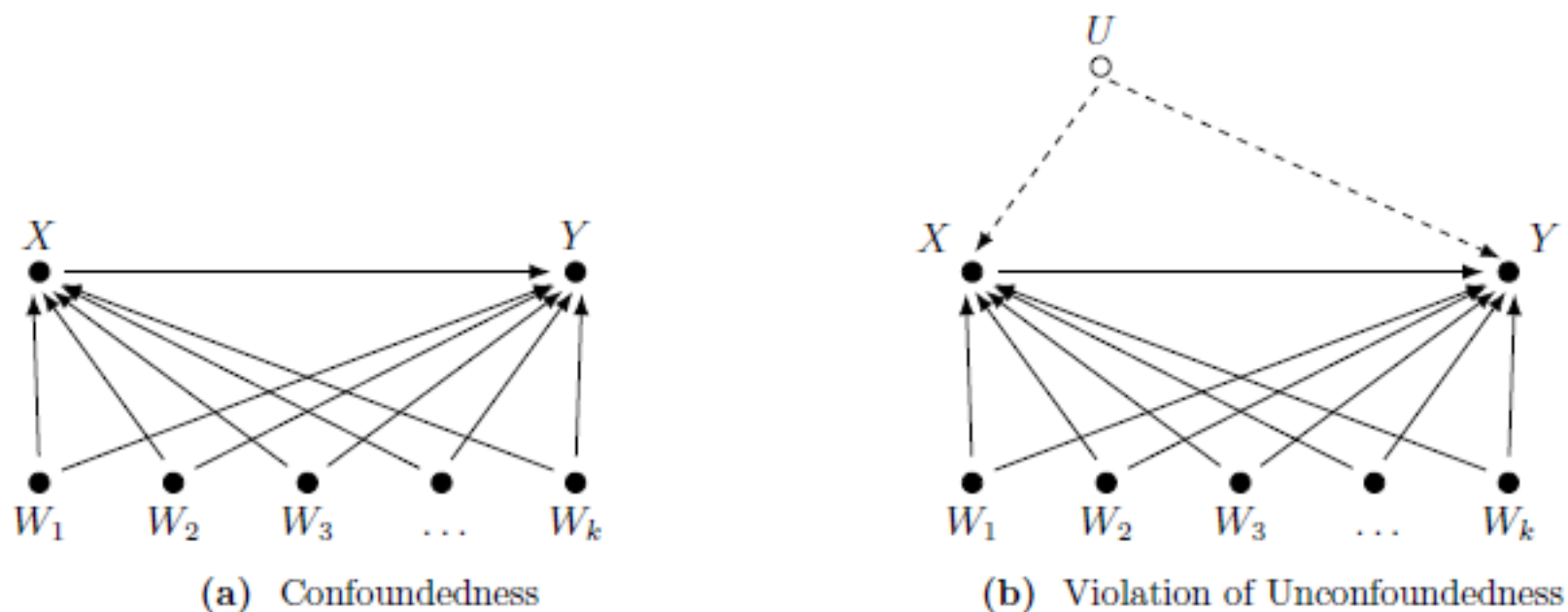
- The  $f_i$ 's denote the causal mechanisms in the model. Are not restricted to be linear.
- $U_X, U_D, U_Y$  denote independent background factors that the we chooses not to include in the analysis.
- Assignment operator ( $\leftarrow$ ) captures asymmetry of causal relationships.

# Confounder

- $X$  is a common cause of  $D$  and  $Y$ .
- conditioning on  $X$  removes dependency between  $D$  and  $Y$
- In DAG terms, controlling for  $X$  "closes the backdoor path" between  $D$  and  $Y$ , and leaves open the direct path.
- The notion of closing the backdoor path is similar to dealing with the omitted variable bias.



# Unconfoundedness in DAGs



**Figure 8:** Unconfoundedness with Multiple Observed Confounders

Source: Imbens (forthcoming).

# Example: Identifying the Returns to Education

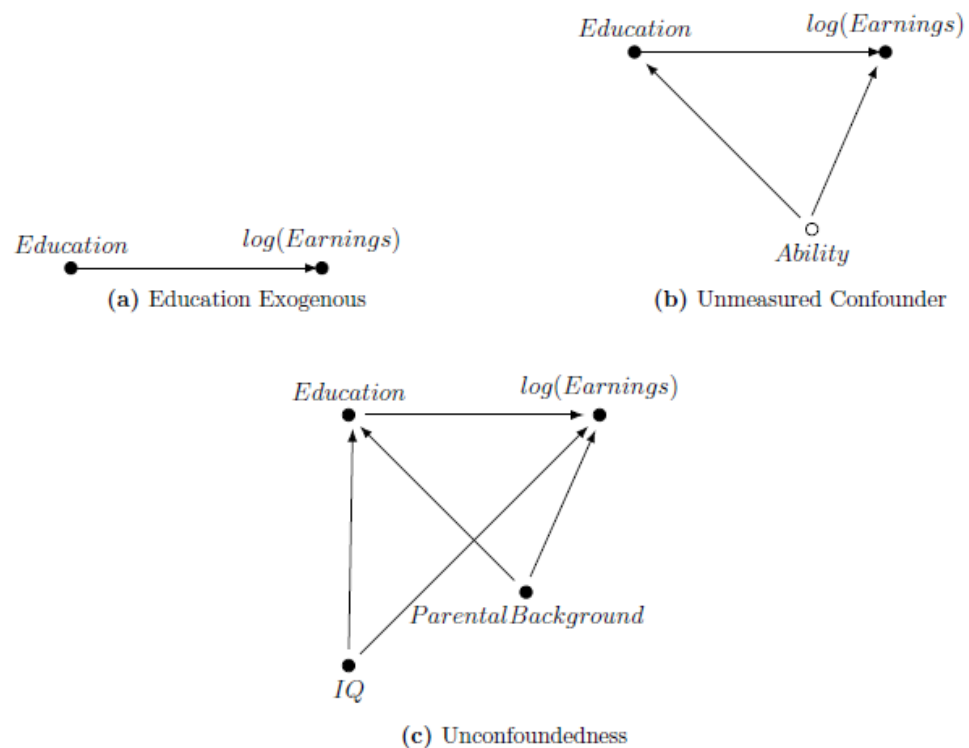
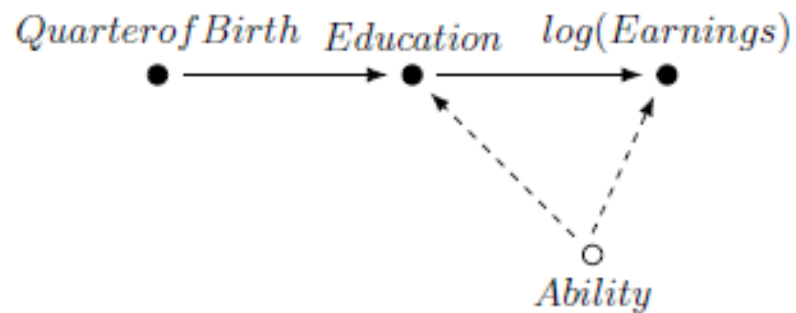
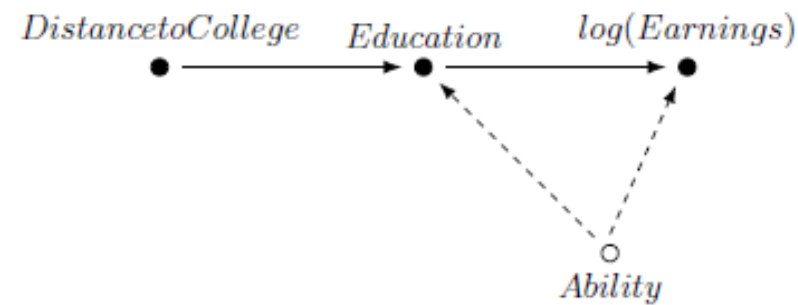


Figure 15: DAGs for the Returns to Education (I)

# Instrumental variables in DAGs



(a) Instrumental Variables: Quarter of Birth



(b) Instrumental Variables: Distance to College

Source: Imbens (forthcoming).

# Mediator

- $D$  causes  $M$  causes  $Y$ .
- $M$  mediates the causal effect of  $D$  on  $Y$
- conditioning on  $M$  removes dependency between  $D$  and  $Y$
- We've essentially closed a direct path (the only direct path between  $D$  and  $Y$ )  
.



# Collider

- $D$  and  $Y$  are independent.
- $D$  and  $Y$  jointly cause  $C$ .
- conditioning on  $C$  creates dependency between  $D$  and  $Y$

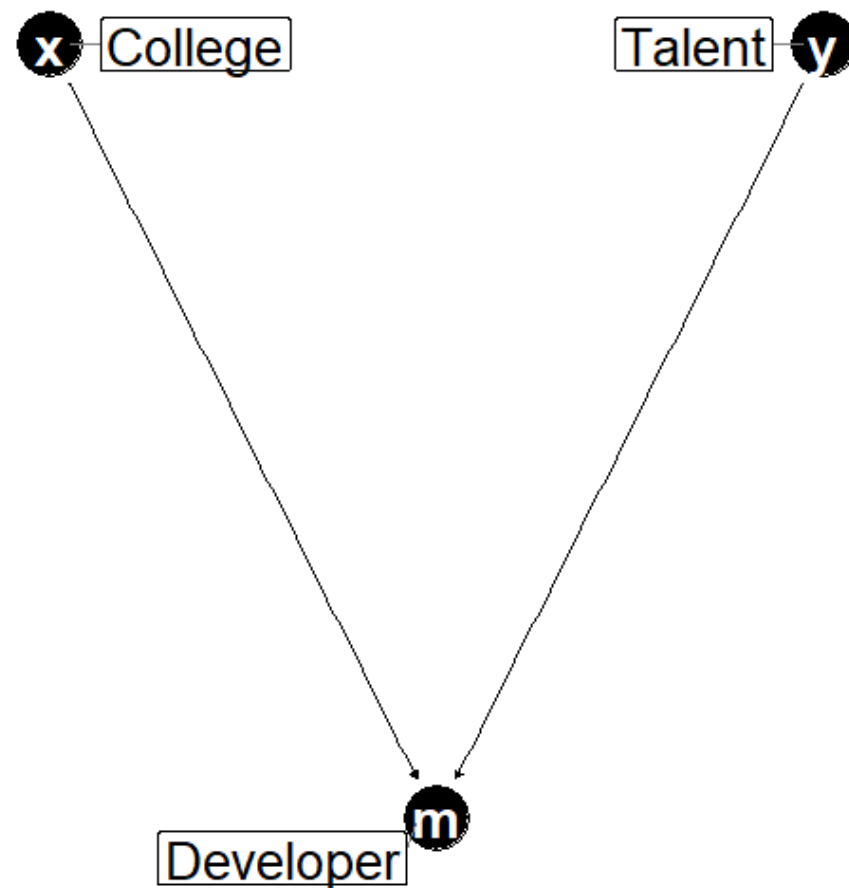


# Example: "Bad controls"

- "Bad controls" are variables that are themselves outcome variables.
- This distinction becomes important when dealing with high-dimensional data

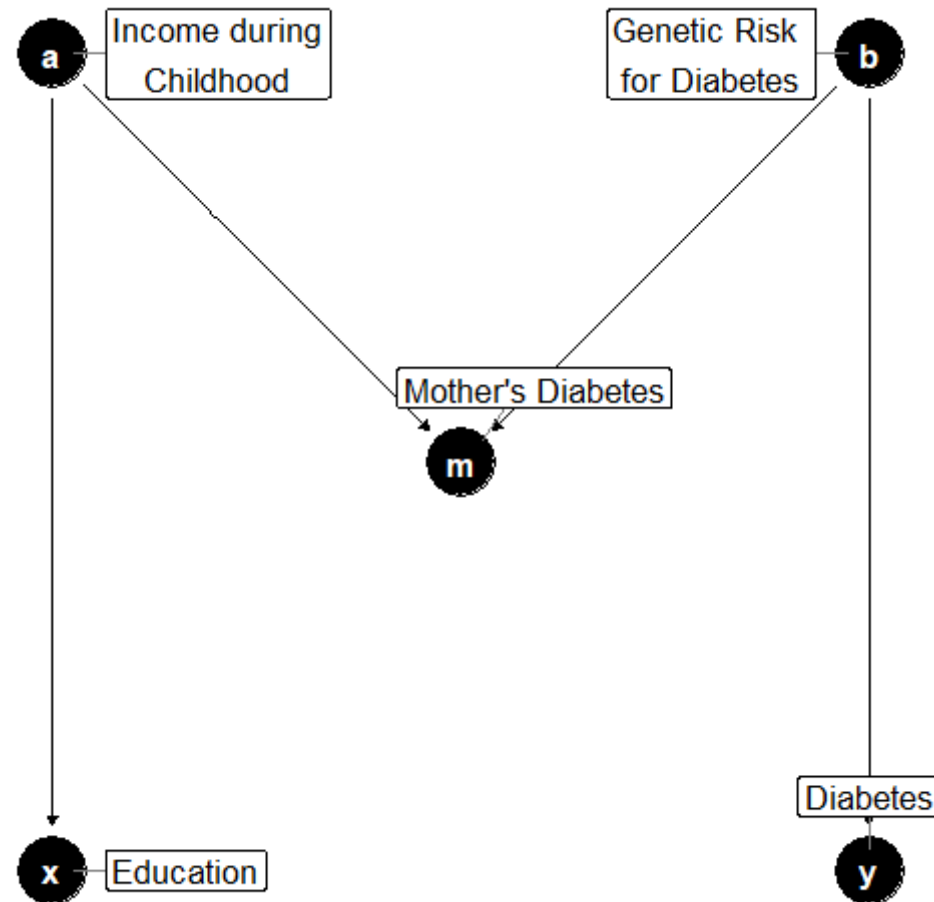
**EXAMPLE:** Occupation as control in a return to years of schooling regression.

Discovering that a person works as a developer in a high-tech firm changes things; knowing that the person does not have a college degree tells us immediately that he is likely to be highly capable.





# Collider: M-bias



# Simulation I: De-confounding

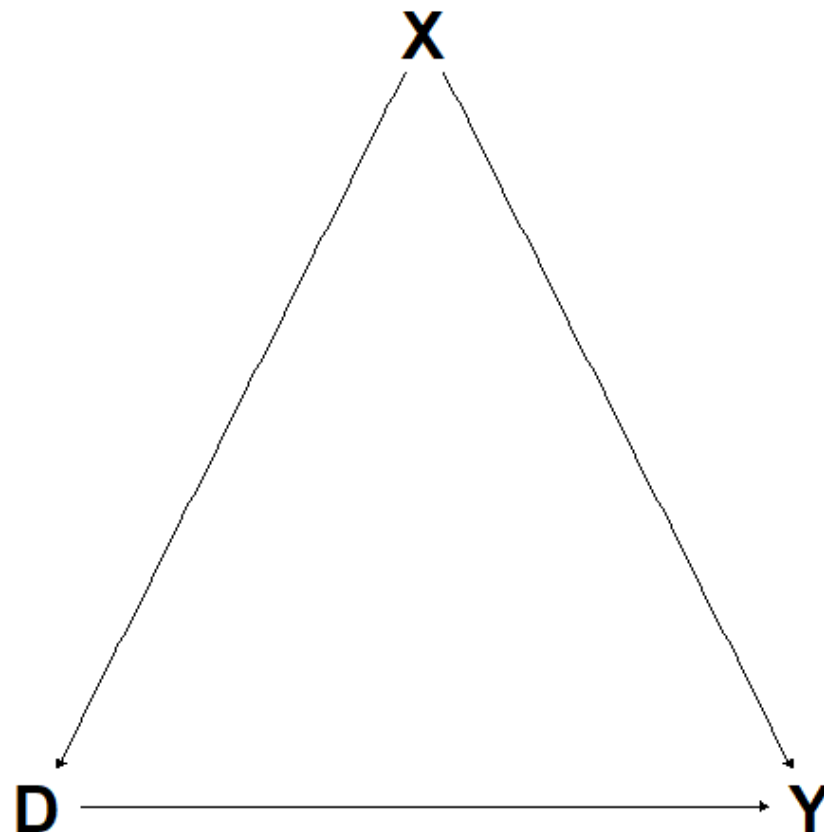
```
n <- 1000
p <- 3

u <- matrix(rnorm(n * p), n, p)

x <- u[,2]
d <- 0.8 * x + 0.6 * u[,1]
y <- 0 * d + 0.2 * x + u[,3]
```

```
cor(cbind(y,x,d)) %>%
  kable(digits = 1, format = "html")
```

	<b>y</b>	<b>x</b>	<b>d</b>
y	1.0	0.2	0.1
x	0.2	1.0	0.8
d	0.1	0.8	1.0



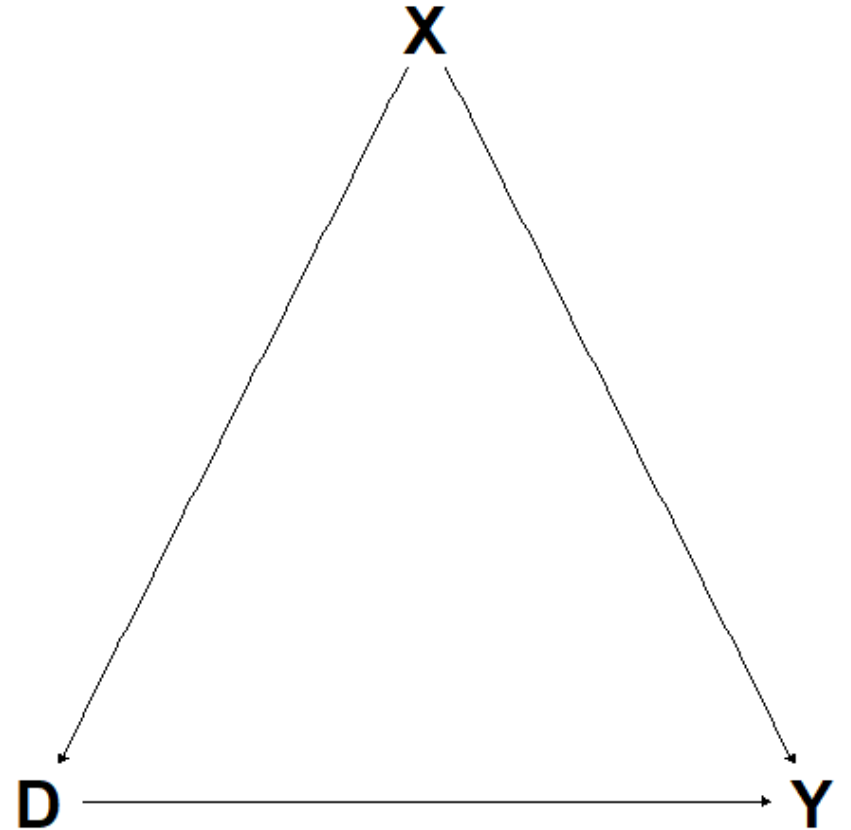
# Simulation I: De-confounding

$y \sim d + x$

term	estimate	p.value
(Intercept)	0.03	0.27
d	-0.05	0.29
x	0.22	0.00

$y \sim d$

term	estimate	p.value
(Intercept)	0.03	0.28
d	0.11	0.00



# Simulation II: M-bias

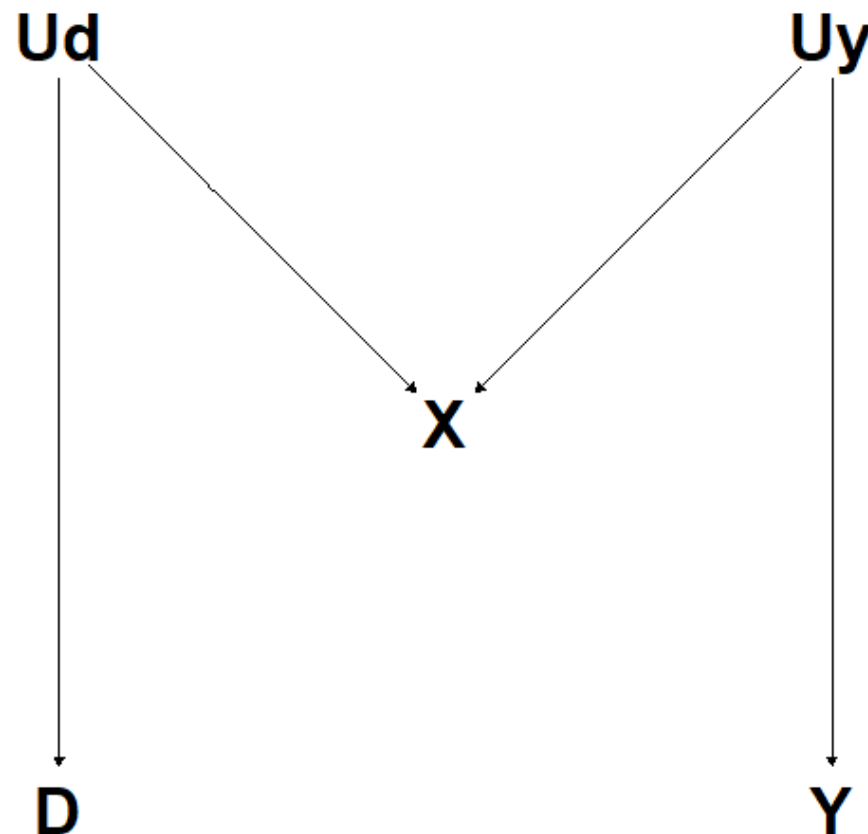
```
n <- 1000
p <- 3

u <- matrix(rnorm(n * p), n, p)

d <- u[,1]
x <- 0.8 * u[,1] + 0.2 * u[,2] + 0.6 * u[,3]
y <- 0 * d + u[,2]
```

```
cor(cbind(y,x,d)) %>%
  kable(digits = 1, format = "html")
```

	<b>y</b>	<b>x</b>	<b>d</b>
y	1.0	0.2	0.0
x	0.2	1.0	0.8
d	0.0	0.8	1.0



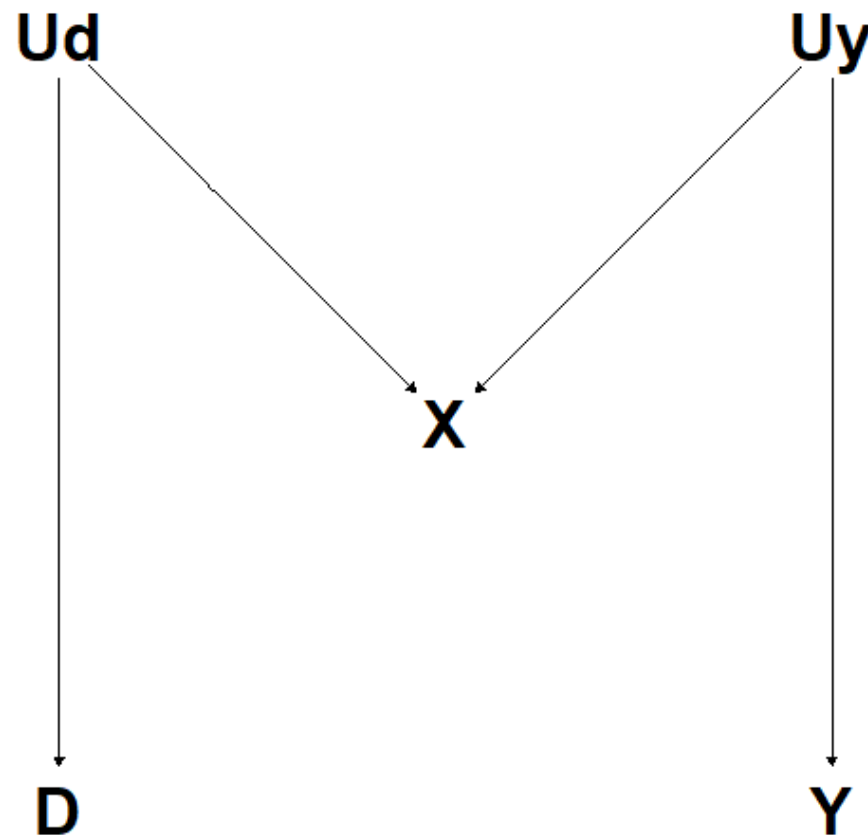
# Simulation II: M-bias

$y \sim d + x$

term	estimate	p.value
(Intercept)	-0.01	0.79
d	-0.40	0.00
x	0.46	0.00

$y \sim d$

term	estimate	p.value
(Intercept)	-0.01	0.81
d	-0.02	0.58



# Simulation III: Mediator

```
n <- 1000
p <- 3

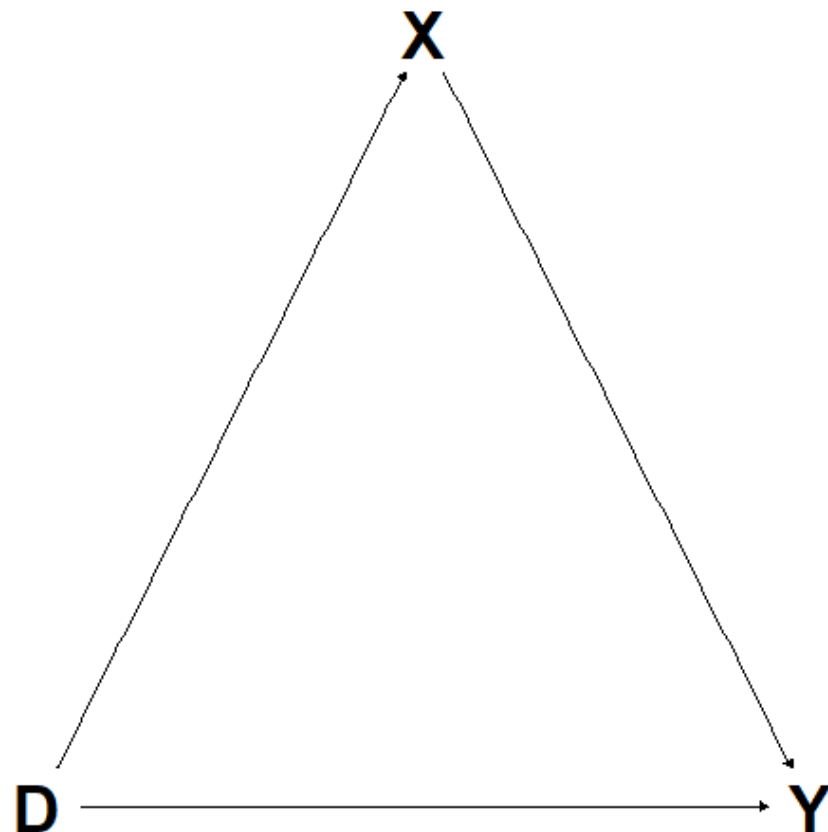
u <- matrix(rnorm(n * p), n, p)

d <- u[,1]
x <- 1.3 * d + u[,2]
y <- 0.1 * d + 0.07 * x + u[,3]
```

```
cor(cbind(y,x,d)) %>%
  kable(digits = 1, format = )
```

	y	x	d
--	---	---	---

y	1.0	0.1	0.2
x	0.1	1.0	0.8
d	0.2	0.8	1.0



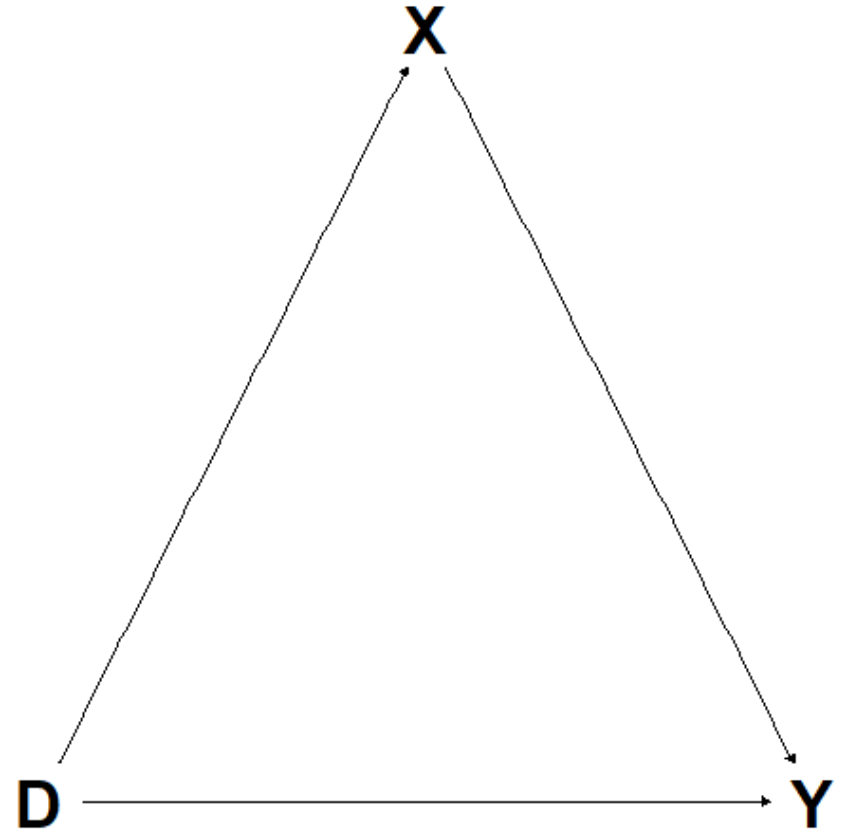
# Simulation III: Mediator

$y \sim d + x$

term	estimate	p.value
(Intercept)	-0.05	0.13
d	0.16	0.00
x	0.01	0.78

$y \sim d$

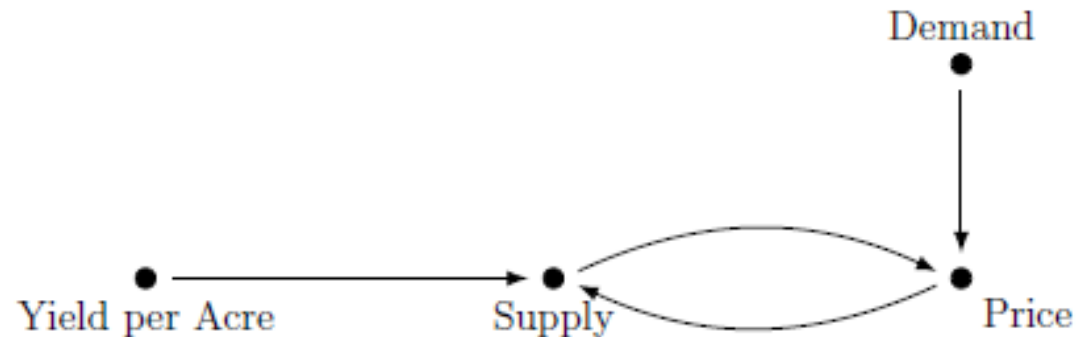
term	estimate	p.value
(Intercept)	-0.05	0.13
d	0.18	0.00



# Limitations of DAGs

- Hard to write down a DAG for complicated (econometric) structural models.
- Need to specify the entire DGP (it REALLY a limitation?)
- Simultaneity

*"In fact it is not immediately obvious to me how one would capture supply and demand models in a DAG" - Imbens (forthcoming)*



**Figure 11:** Based on Figure 7.10 in TBOW, p. 251.



# Recommended introductory level resources on DAGS

- [The Book of Why](#) by Pearl and Mackenzie.
- [Causal Inference in Machine Learning and AI](#) by Paul Hünermund.
- [Causal Inference: The Mixtape \(pp. 67-80\)](#) by Scott Cunningham.
- [Potential Outcome and Directed Acyclic Graph Approaches to Causality: Relevance for Empirical Practice in Economics](#) by Guido W. Imbens
- [A Crash Course in Good and Bad Controls](#) by Cinelli, Forney, and Pearl, J. (2020).

# Next time: Causal inference in high-dimensional setting

Consider again the standard "treatment effect regression":

$$Y_i = \alpha + \underbrace{\tau D_i}_{\text{low dimensional}} + \underbrace{\sum_{j=1}^k \beta_j X_{ij}}_{\text{high dimensional}} + \varepsilon_i, \quad \text{for } i = 1, \dots, n$$

Our object of interest is  $\hat{\tau}$ , the estimated *average treatment effect* (ATE).

In high-dimensional settings  $k \gg n$ .

```
slides %>% end()
```

 [Source code](#)

# Selected references

Hünermund, P., & Bareinboim, E. (2019). Causal Inference and Data-Fusion in Econometrics. arXiv preprint arXiv:1912.09104.

Imbens, W. G. (forthcoming). Potential Outcome and Directed Acyclic Graph Approaches to Causality: Relevance for Empirical Practice in Economics. *Journal of Economic Literature*.

Lewbel, A. (2019). The identification zoo: Meanings of identification in econometrics. *Journal of Economic Literature*, 57(4), 835-903.