# 3.2 — DAGs

## ECON 480 • Econometrics • Fall 2022

Dr. Ryan Safner
Associate Professor of Economics

✈ safner@hood.edu
   ryansafner/metricsF22
🌐 metricsF22.classes.ryansafner.com

# Contents

# Causation and Correlation

# You Don't Need an RCT to Talk About Causality

- Statistics profession is obstinant that we cannot say anything about causality

- But you have to! It's how the human brain works!

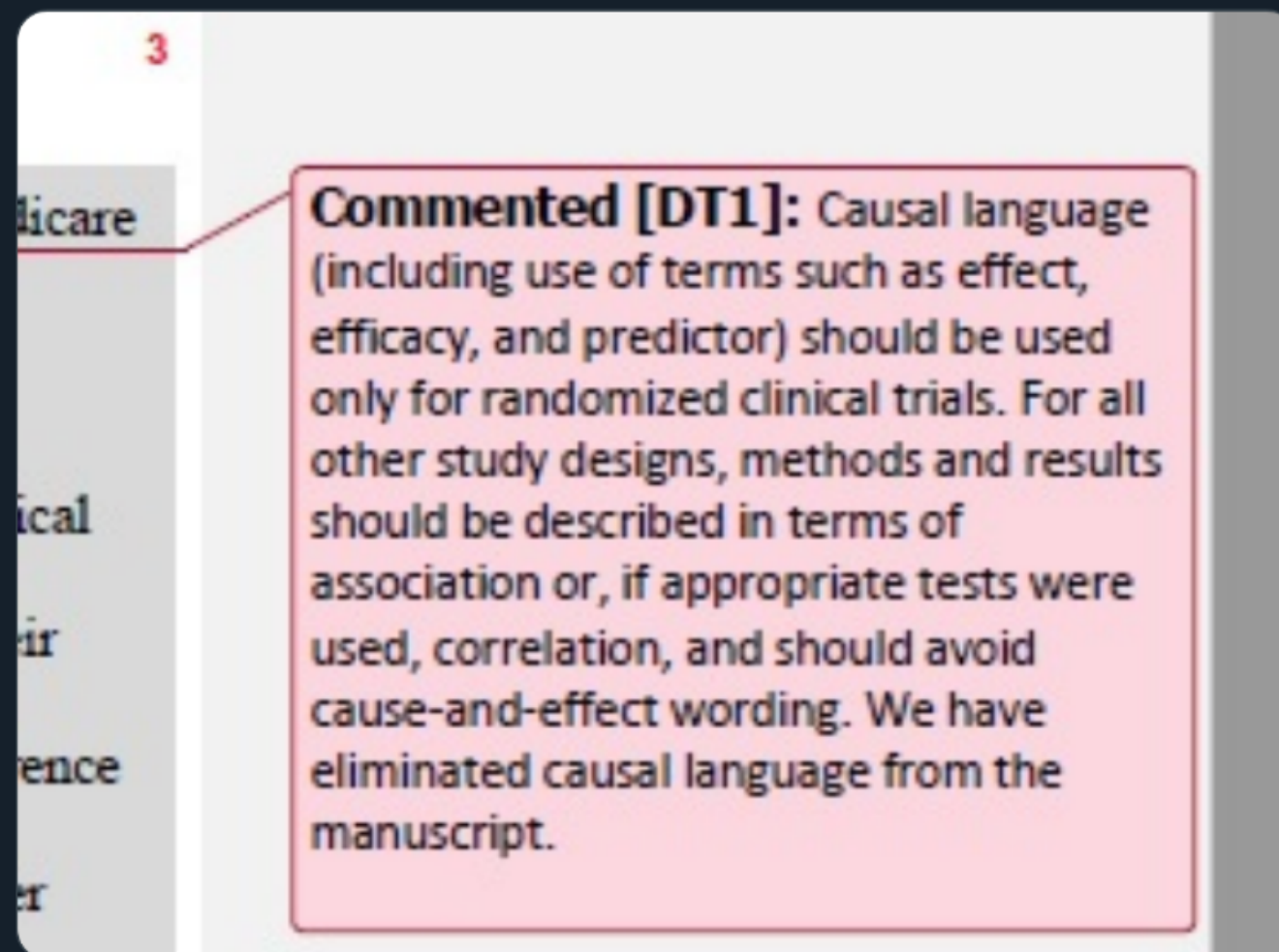- We can't conceive of (spurious) correlation without some causation

# The Causal Revolution

# RCTs and Evidence-Based Policy

## Should we *ONLY* base policies on the evidence from Randomized Controlled Trials

**Dr Ellie Murray, ScD** ✔
@EpiEllie · Follow

```
 _____
|               |
|  IF U DONT SMOKE,
|   U ALREADY
|    BELIEVE IN
| CAUSAL INFERENCE
|    WITHOUT
| RANDOMIZED TRIALS
|_____|
    (\__/)  ||
    (•ㅅ•)  ||
    /  　づ
```

#HistorianSignBunny #Epidemiology

12:13 AM · Jul 13, 2018

Read the full conversation on Twitter

♥ 910    💬 Reply    🔗 Copy link

**Read 33 replies**

---

**thebmj**   Research ⌄   Education ⌄   News & Views ⌄   Campaigns ⌄

**Research**    ⓒⓘⓢ Open access
Christmas 2018: Look Before You Leap

## Parachute use to prevent death and major trauma when jumping from aircraft: randomized controlled trial

*BMJ* 2018 ; 363  doi: https://doi.org/10.1136/bmj.k5094 (Published 13 December 2018)
Cite this as: *BMJ* 2018;363:k5094

🐦 Tweet    👍 Like 31K

See other articles in issue 8180

**Article tools**

**Linked opinion**
We jumped from planes without parachutes (and lived to tell the tale)

📄 PDF     💬 25 responses

✎ Respond to this article

➕ Data supplement

🖨 Print

📢 Alerts & updates ⌄

🌐 Citation tools ⌄

◎ Request permissions

👤 Author citation ⌄

➕ Add article to BMJ Portfolio

✉ Email to a friend ⌄

| Article | Related content | Metrics | Responses |

| Peer review |

This article has a correction. Please see:
Parachute use to prevent death and major trauma when jumping from aircraft: randomized controlled trial - December 18, 2018

**Topics**

*Robert W Yeh, associate professor*[1], *Linda R Valsdottir, research coordinator*[1], *Michael W Yeh, professor*[2], *Changyu Shen, director*[1], *Daniel B Kramer, assistant professor*[1], *Jordan B Strom, instructor*[1], *Eric A Secemsky, instructor*[1], *Joanne L Healy, administrative manager*[1], *Robert M Domeier, expert skydiver and clinical instructor*[3], *Dhruv S Kazi, associate director*[1],

Christmas

# RCTs and Evidence-Based Policy II

# RCTs and Evidence-Based Policy II

# Correlation vs. Causation I

David Robinson · Jun 22, 2017
@drob · Follow

Correlation implies causation, don't @ me

David Robinson
@drob · Follow

"Correlation implies casuation," the dean whispered as he handed me my PhD.

"But then why-"

"Because if they knew, they wouldn't need us."

3:46 PM · Jun 22, 2017 from Manhattan, NY

♥ 156      Reply      Copy link

**Read 4 replies**

# What Does Causation Mean?

- **"Correlation does not imply causation"**

  - this is exactly backwards!

  - this is just pointing out that **exogeneity is violated**

# What Does Causation Mean?

- **"Correlation does not imply causation"**

  - this is exactly backwards!

  - this is just pointing out that **exogeneity is violated**

- **"Correlation implies causation"**

  - for an association, there must be *some* causal chain that relates $X$ and $Y$

  - but not necessarily *merely* $X \rightarrow Y$

- **"Correlation plus exogeneity is causation."**

# Correlation and Causation

- **Correlation:**
  - Math & Statistics
  - Computers, AI, Machine learning can figure this out (better than humans)
- **Causation:**
  - Philosophy, Intuition, Theory
  - **Counterfactual thinking**, unique to humans (vs. animals or computers)
  - Computers *cannot* (yet) figure this out

# The Causal Revolution

# Causation Requires Counterfactual Thinking

IMAGINING

DOING

SEEING

**3. COUNTERFACTUALS**

ACTIVITY: Imagining, Retrospection, Understanding

QUESTIONS: *What if I had done …? Why?*
(Was it X that caused Y? What if X had not occurred? What if I had acted differently?)

EXAMPLES: Was it the aspirin that stopped my headache? Would Kennedy be alive if Oswald had not killed him? What if I had not smoked for the last 2 years?

**2. INTERVENTION**

ACTIVITY: Doing, Intervening

QUESTIONS: *What if I do …? How?*
(What would Y be if I do X? How can I make Y happen?)

EXAMPLES: If I take aspirin, will my headache be cured? What if we ban cigarettes?

**1. ASSOCIATION**

ACTIVITY: Seeing, Observing

QUESTIONS: *What if I see …?*
(How are the variables related? How would seeing X change my belief in Y?)

EXAMPLES: What does a symptom tell me about a disease? What does a survey tell us about the election results?

JUDEA PEARL
WINNER OF THE TURING AWARD
AND DANA MACKENZIE

THE
BOOK OF
WHY

α ⟶ β

THE NEW SCIENCE
OF CAUSE AND EFFECT

# Causal Inference

- We will seek to understand what causality *is* and how we can approach finding it

- We will also explore the different common **research designs** meant to **identify** causal relationships

- **These skills**, more than supply & demand, constrained optimization models, ISLM, etc, **are the tools and comparative advantage of a modern research economist**

  - Why all big companies (especially in tech) have entire economics departments in them

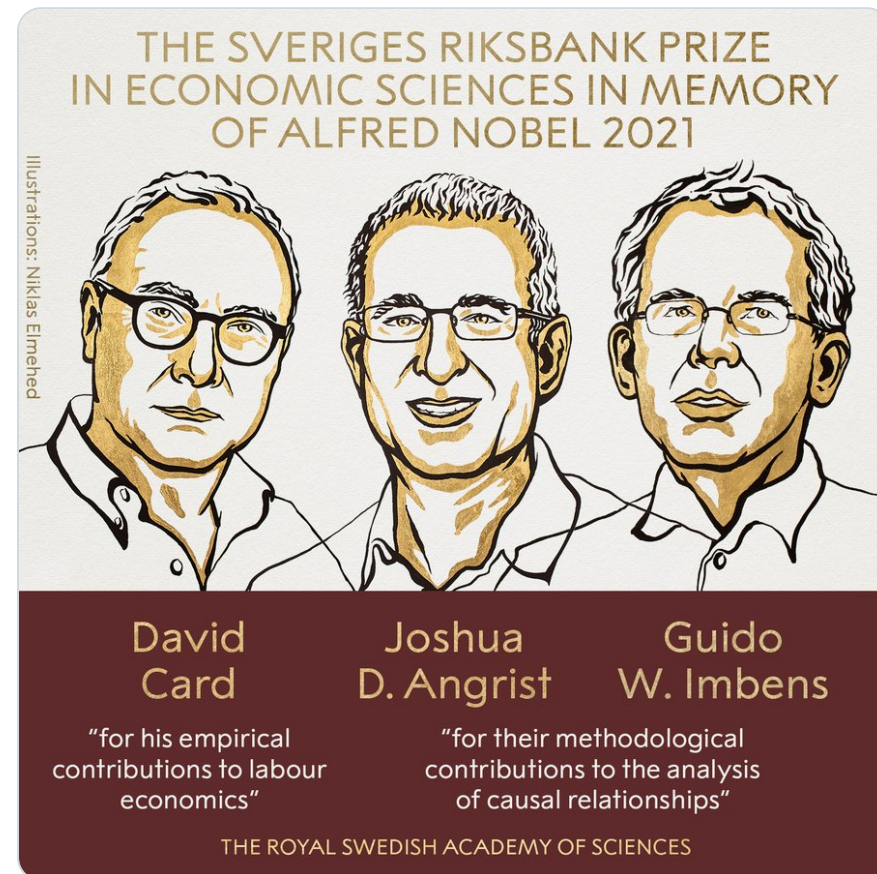# "The Credibility Revolution" in Econometrics



- Simultaneous **"credibility revolution"** in econometrics (c.1990s—2000s)

- Use clever research designs to approximate **natural experiments**

- Note: major disagreements between Pearl & Angrist/Imbens, etc.!

# Clever Research Designs Identify Causality

Diff-in-Diff

Fixed Effects

Multivariate Regression                                          RCTs

Differences                          Regression Discontinuity

Pre-Post              Matching              Natural Experiments

Correlation                                                        Causation

# Correlation and Causation

**John B. Holbein** · Apr 7, 2018
@JohnHolbein1 · **Follow**
Causality isn't binary; it's a continuum.

**John B. Holbein**
@JohnHolbein1 · **Follow**

Causality isn't achieved; it's approached.

11:05 AM · Apr 7, 2018

♥ 7    💬 **Reply**    🔗 **Copy link**

**Read 1 reply**

# What Then IS Causation?

# What Then IS Causation?

# Non-Causal Claims

- All of the following have non-zero correlations. Are they *causal*?

**Examples**

- Greater ice cream sales → more violent crime
- Rooster crows → the sun rises in the morning
- Taking Vitamin C → colds go away a few days later
- Political party $X$ in power → economy performs better/worse

# Counterfactuals

- The *sine qua non* of causal claims are **counterfactuals**: what would $Y$ have been if $X$ had been different?

- It is **impossible** to make a counterfactual claim from data alone!

- Need a (theoretical) **causal model** of the data-generating process!

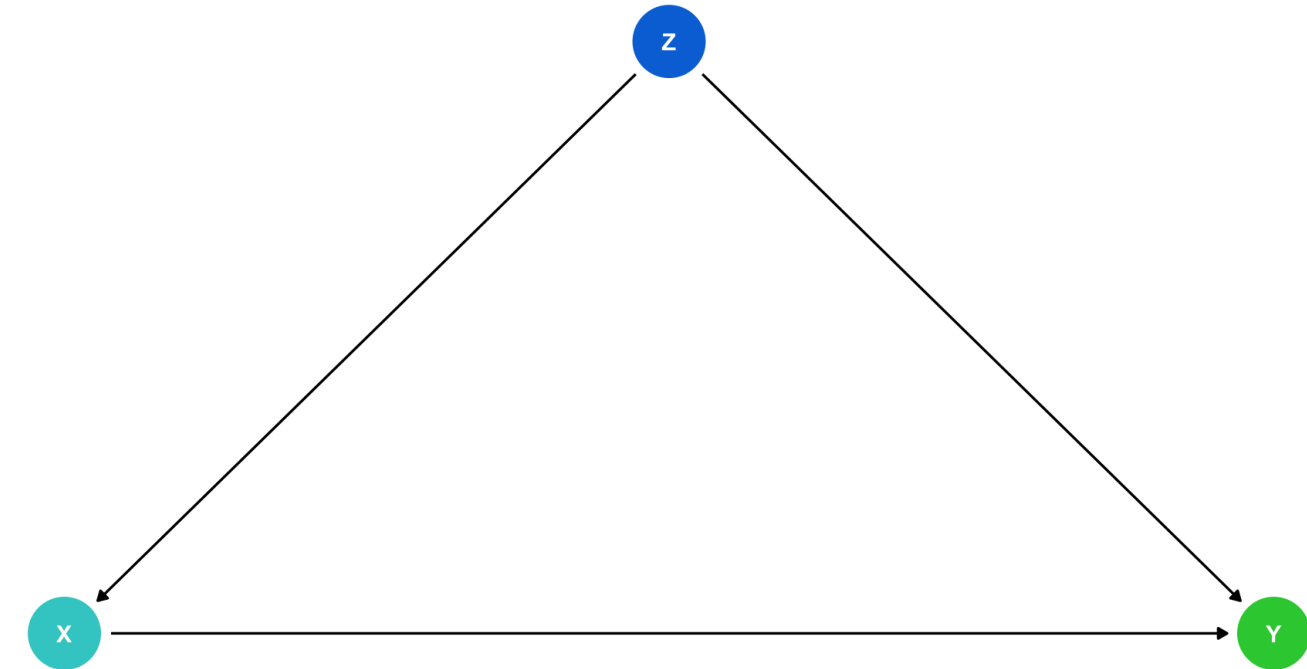# Counterfactuals and RCTs

# From RCTs to Causal Models

- RCTs are but the best-known method of a large, growing science of **causal inference**

- We need a **causal model** to describe the **data-generating process (DGP)**

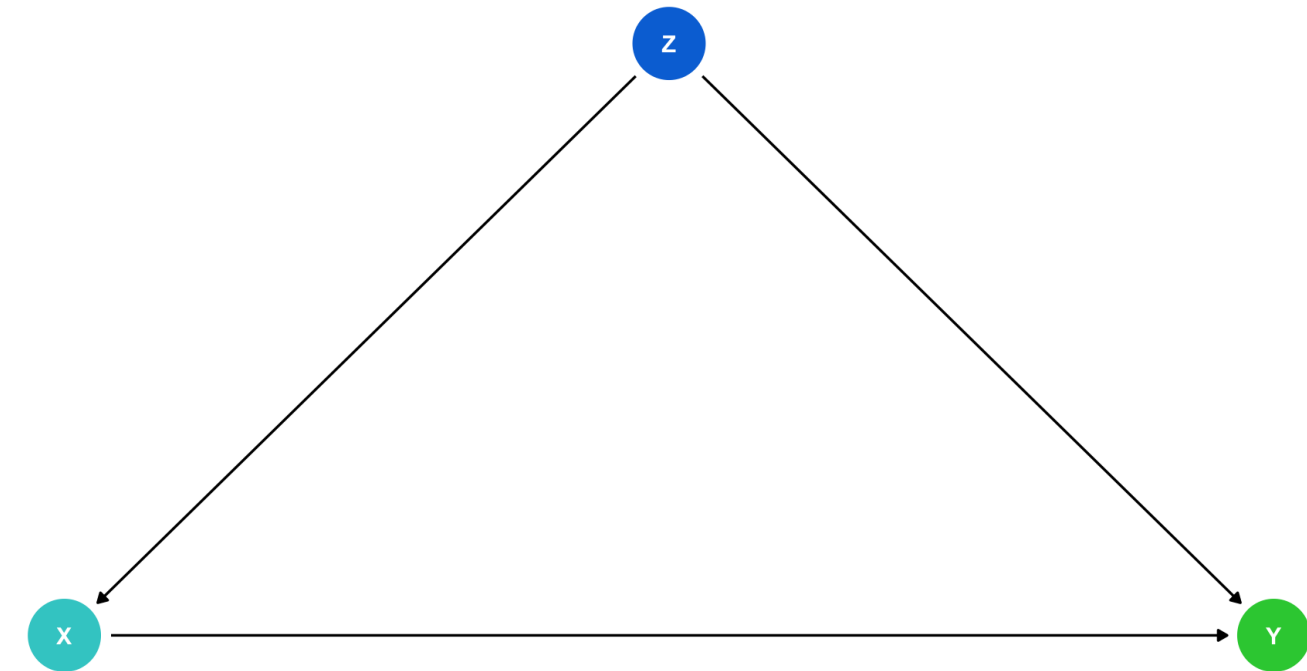- Requires us to make some **assumptions**

# Causal Diagrams

# Causal Diagrams/DAGs

- A surprisingly simple, yet rigorous and powerful method of modeling is using a **causal diagram** or **DAG**:

  - **Directed**: Each node has arrows that points only one direction

  - **Acyclic**: Arrows only have one direction, and cannot loop back
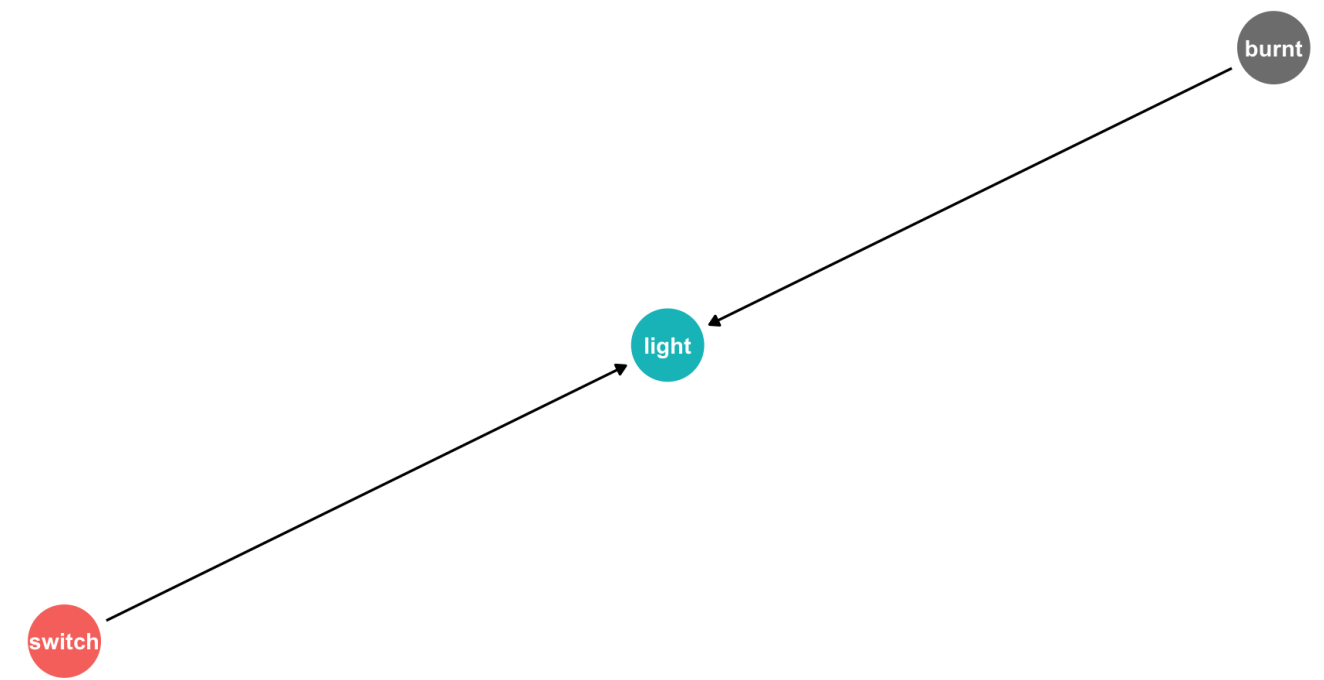
  - **Graph**

# Causal Diagrams/DAGs

- A visual model of the data-generating process, encodes our understanding of the causal relationships

- Requires some common sense/economic intuition

- Remember, all models are wrong, we just need them to be *useful*!
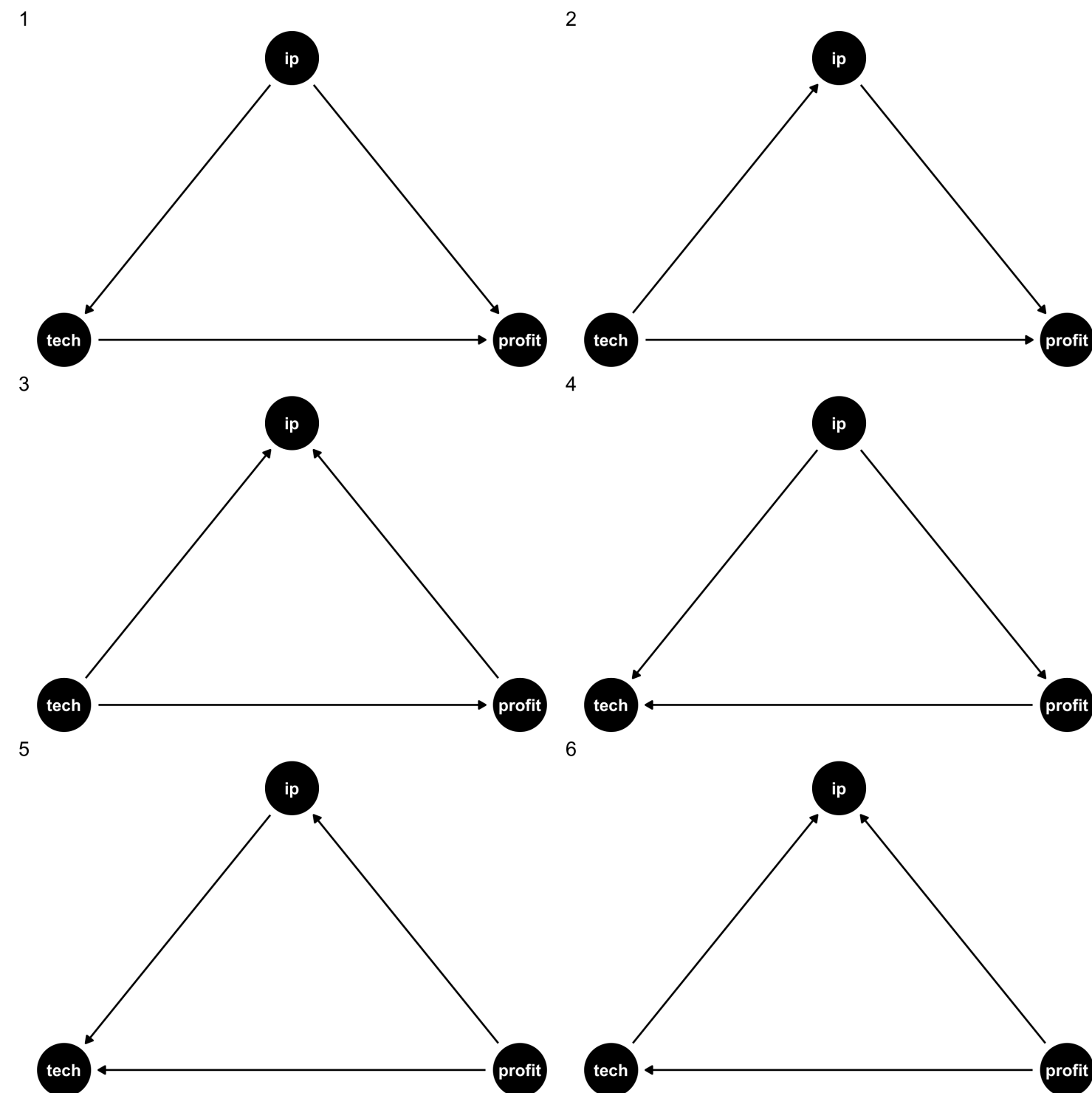
# Causal Diagrams/DAGs

- Our light switch example of causality

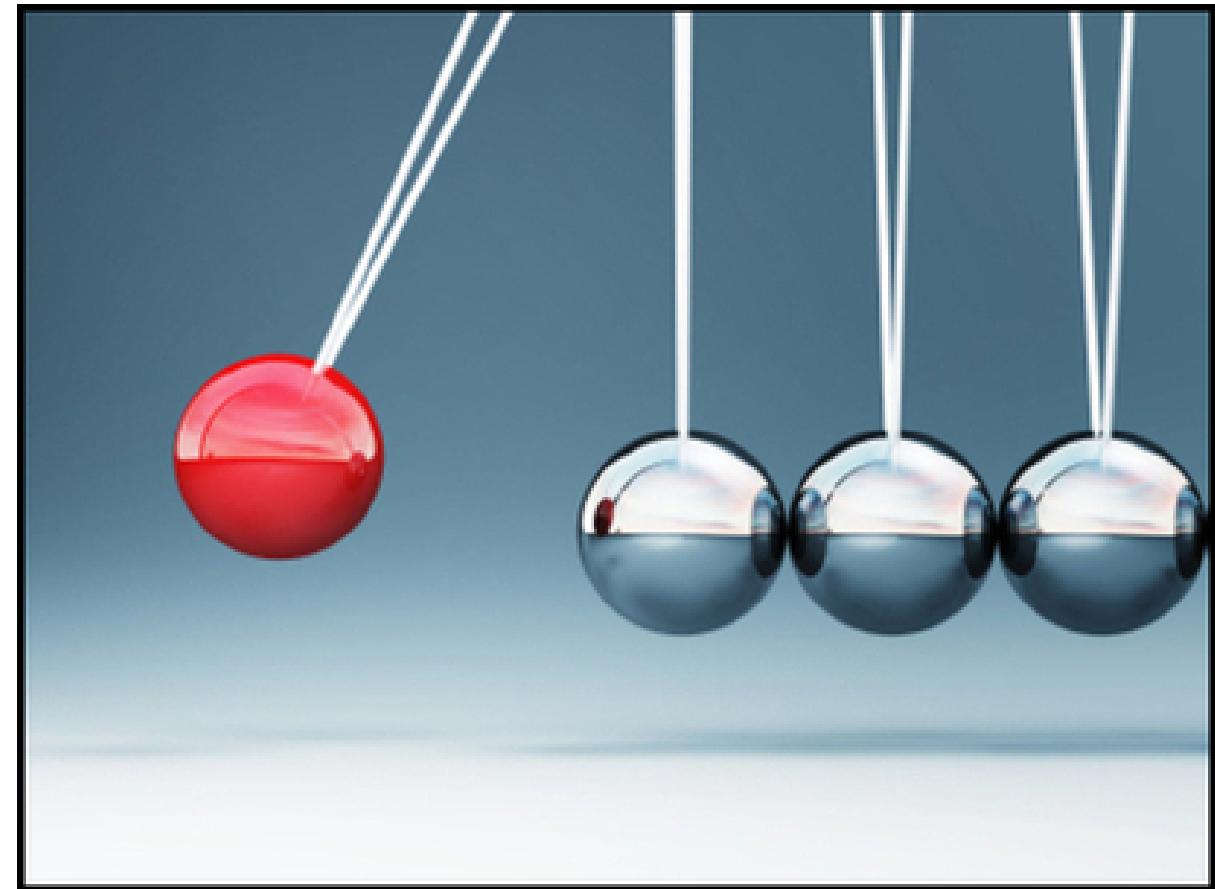# Drawing a DAG: Example

- Suppose we have data on three variables

  - IP: how much a firm spends on IP lawsuits

  - tech: whether a firm is in tech industry

  - profit: firm profits

- They are all correlated with each other, but what's are the causal relationships?

- We need our own **causal model** (from theory, intuition, etc) to sort

  - Data alone will not tell us!

# Drawing a DAG
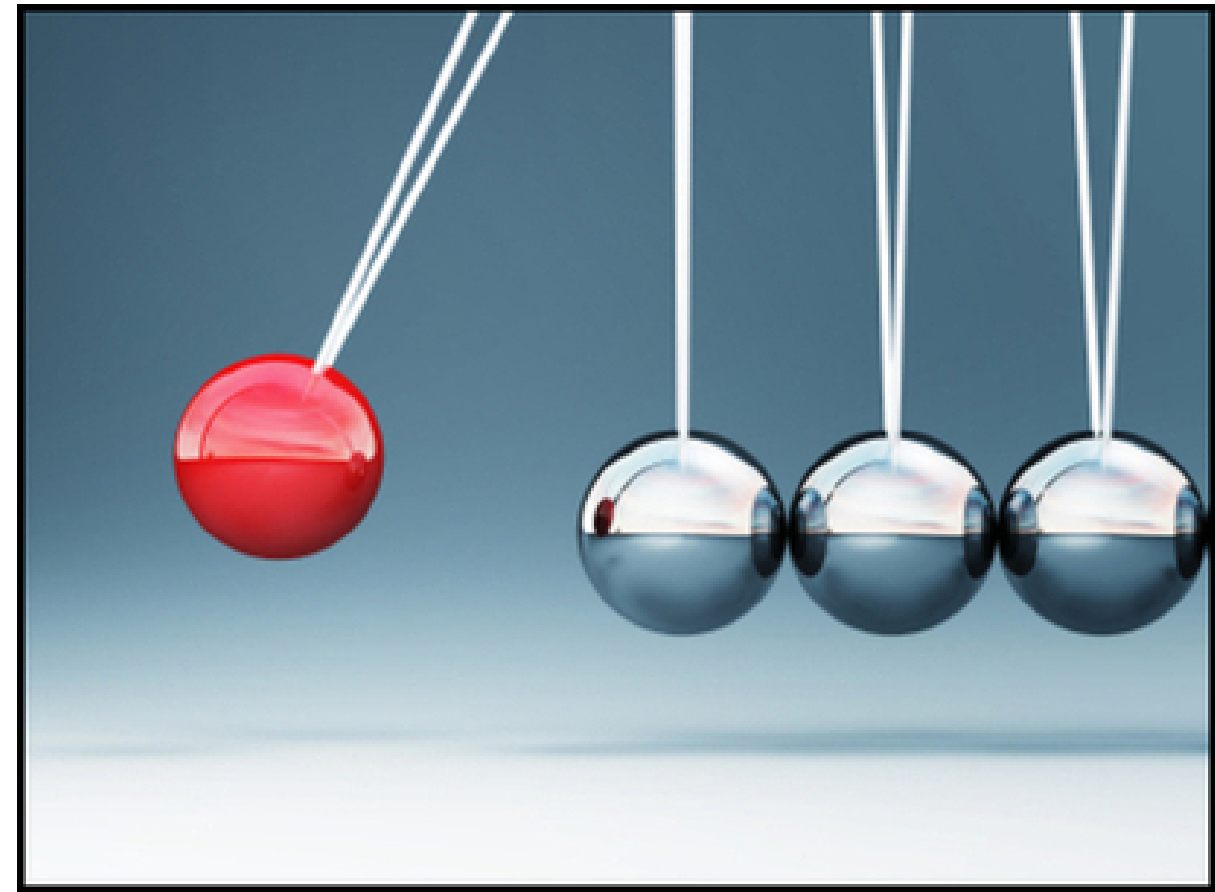
1. Consider all the variables likely to be important to the data-generating process (including variables we can't observe!)

2. For simplicity, combine some similar ones together or prune those that aren't very important

3. Consider which variables are likely to affect others, and draw arrows connecting them

4. Test some testable implications of the model (to see if we have a correct one!)

# Drawing a DAG

- Drawing an arrow requires a direction - making a statement about causality!

- *Omitting* an arrow makes an equally important statement too!

    - In fact, we will *need* omitted arrows to show causality!

- If two variables are correlated, but neither causes the other, likely they are both caused by another (perhaps **unobserved**) variable - add it!

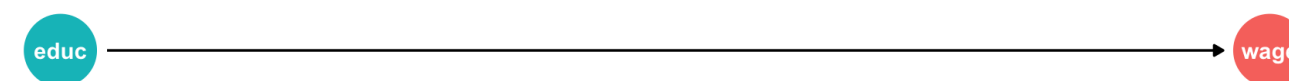- There should be no *cycles* or *loops* (if so, there's probably another missing variable, such as time)

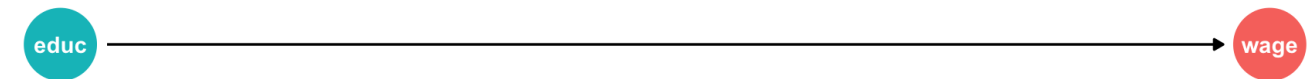# DAG Example I

> **Example**
>
> what is the effect of education on wages?
>
> - Education $X$, "treatment" or "exposure"
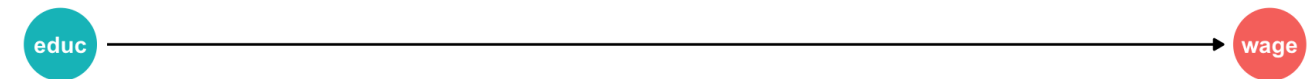> - Wages $Y$, "outcome" or "response"

# DAG Example I

- What other variables are important?
  - Ability
  - Socioeconomic status
  - Demographics
  - Phys. Ed. requirements
  - Year of birth
  - Location
  - Schooling laws
  - Job connections

# DAG Example I

- In social science and complex systems, 1000s of variables could plausibly be in DAG!

- So simplify:

  - Ignore trivial things (Phys. Ed. requirement)

  - Combine similar variables (Socioeconomic status, Demographics, Location) $\rightarrow$ Background

# DAG Example II

- Background, Year of birth, Location, Compulsory schooling, all cause education

- Background, year of birth, location, job connections probably cause wages
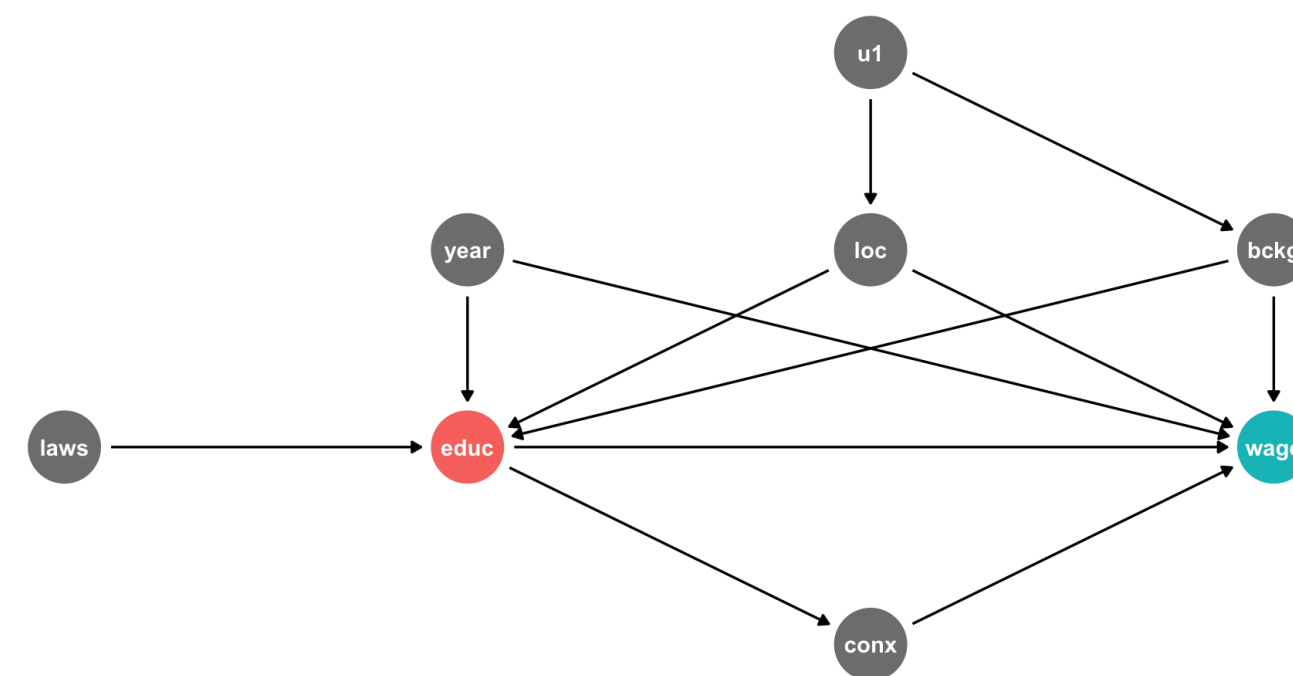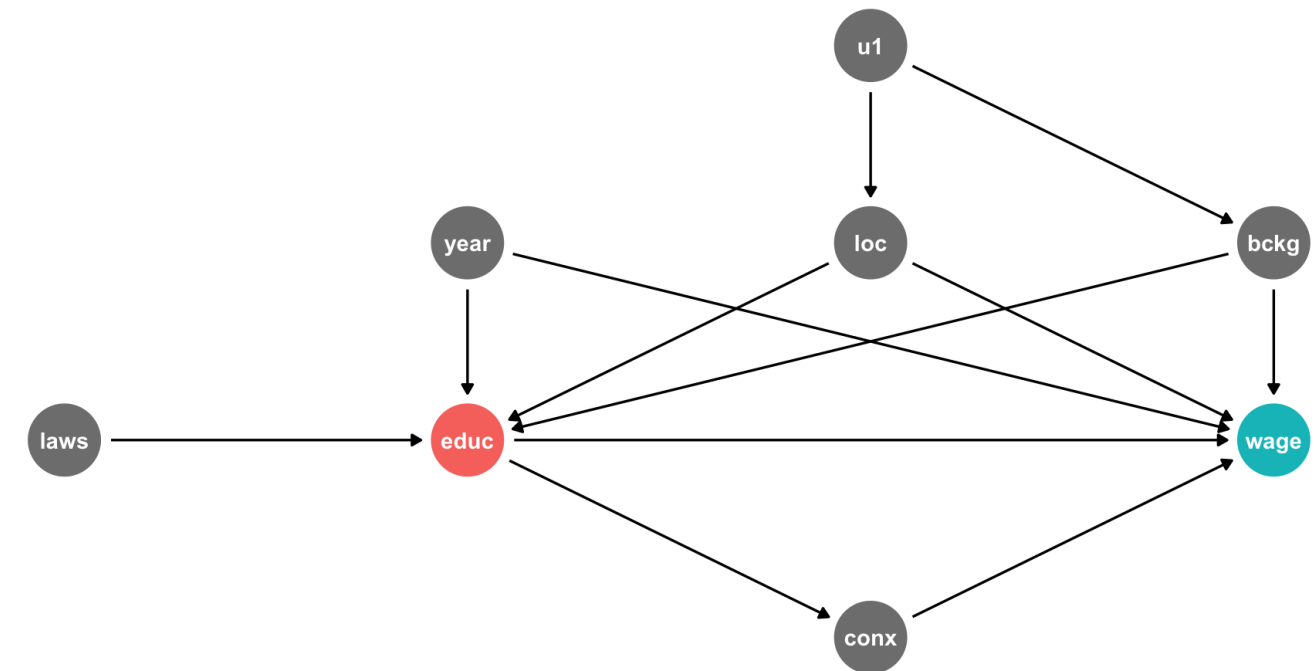
# DAG Example II

- Background, Year of birth, Location, Compulsory schooling, all cause education

- Background, year of birth, location, job connections probably cause wages

- Job connections in fact is probably caused by education!

- Location and background probably both caused by unobserved factor (u1)

# DAG Example II

- This is messy, but we have a causal model!

- Makes our assumptions **explicit**, and many of them are **testable**

- DAG suggests certain relationships that will *not* exist:

  - all relationships between `laws` and `conx` go through `educ`

  - so if we controlled for `educ`, then `cor(laws,conx)` should be zero!

# Let the Computer Do It: Dagitty.net I



- Dagitty.net is a great tool to make these and give you testable implications

- Click `Model -> New Model`

- Name your "exposure" variable ($X$ of interest) and "outcome" variable ($Y$)

# Let the Computer Do It: Dagitty.net II



- Click and drag to move nodes around

- Add a new variable by double-clicking

- Add an arrow by double-clicking one variable and then double-clicking on the target (do again to remove arrow)

# Let the Computer Do It: Dagitty.net II

# Let the Computer Do It: Dagitty.net III



- Tells you **how to identify your effect**! (upper right)

> **Minimal sufficient adjustment sets** containing background, location, year for estimating the total effect of educ on wage: background, location, year

# Let the Computer Do It: Dagitty.net III



- Tells you some **testable implications** of your model

- These are **(conditional) independencies**:

$$X \perp Y \mid Z$$

| "X is independent of Y, given Z"

- Implies that by *controlling for* $Z$, $X$ and $Y$ should have *no correlation*

# Let the Computer Do It: Dagitty.net III



- Tells you some **testable implications** of your model

- Example: look at the last one listed:

$$\texttt{job\_connections} \perp \texttt{year} \mid \texttt{educ}$$

> "Job connections are independent of year, controlling for education"

- Implies that by controlling for $\texttt{educ}$, there should be no correlation between $\texttt{job\_connections}$ and $\texttt{year}$ — can test this with data!

# Causal Effect



- If we control for `background`, `location`, and `year`, we can **identify the causal effect** of `educ` $\rightarrow$ `wage`.

# You Can Draw DAGs in R

- New package: ggdag

- Arrows are made with formula notation: Y ~ X + Z means "Y is caused by X and Z"

```
1   library(ggdag)
2   dagify(wage ~ educ + conx + year + bckg + loc,
3         educ ~ bckg + year + loc + laws,
4         conx ~ educ,
5         bckg ~ u1,
6         loc ~ u1,
7         exposure = "educ", # optional: define X
8         outcome = "wage" # optional: define Y
9         ) %>%
10  ggdag()+
11  theme_dag()
```

# You Can Draw DAGs in R II

- Or you can just copy the code from `dagitty.net`!

- Use `dagitty()` from the `dagitty` package, and paste the code in quotes

```
1  # install.packages("dagitty")
2  library(dagitty)
3  dagitty('dag {
4  bb="0,0,1,1"
5  background [pos="0.413,0.335"]
6  compulsory_schooling_laws [pos="0.544,0.076"]
7  educ [exposure,pos="0.185,0.121"]
8  job_connections [pos="0.302,0.510"]
9  location [pos="0.571,0.431"]
10 u1 [pos="0.539,0.206"]
11 wage [outcome,pos="0.552,0.761"]
12 year [pos="0.197,0.697"]
13 background -> educ
14 background -> wage
15 compulsory_schooling_laws -> educ
16 educ -> job_connections
17 educ -> wage
```

# You Can Draw DAGs In R

- It's not very pretty, but if you set `text = FALSE, use_labels = "name` inside `ggdag()`, makes it easier to read

```
1   dagitty('dag {
2   bb="0,0,1,1"
3   background [pos="0.413,0.335"]
4   compulsory_schooling_laws [pos="0.544,0.076"]
5   educ [exposure,pos="0.185,0.121"]
6   job_connections [pos="0.302,0.510"]
7   location [pos="0.571,0.431"]
8   u1 [pos="0.539,0.206"]
9   wage [outcome,pos="0.552,0.761"]
10  year [pos="0.197,0.697"]
11  background -> educ
12  background -> wage
13  compulsory_schooling_laws -> educ
14  educ -> job_connections
15  educ -> wage
16  job_connections -> wage
17  location -> educ
```

# ggdag: Additional Tools

- If you have defined X (`exposure`) and Y (`outcome`), you can use `ggdag_paths()` to have it show all possible paths between $X$ and $Y$!

```
1  dagify(wage ~ educ + conx + year + bckg + loc,
2         educ ~ bckg + year + loc + laws,
3         conx ~ educ,
4         bckg ~ u1,
5         loc ~ u1,
6         exposure = "educ",
7         outcome = "wage"
8         ) %>%
9   tidy_dagitty(seed = 2) %>%
10  ggdag_paths()+ #<<
11  theme_dag()
```
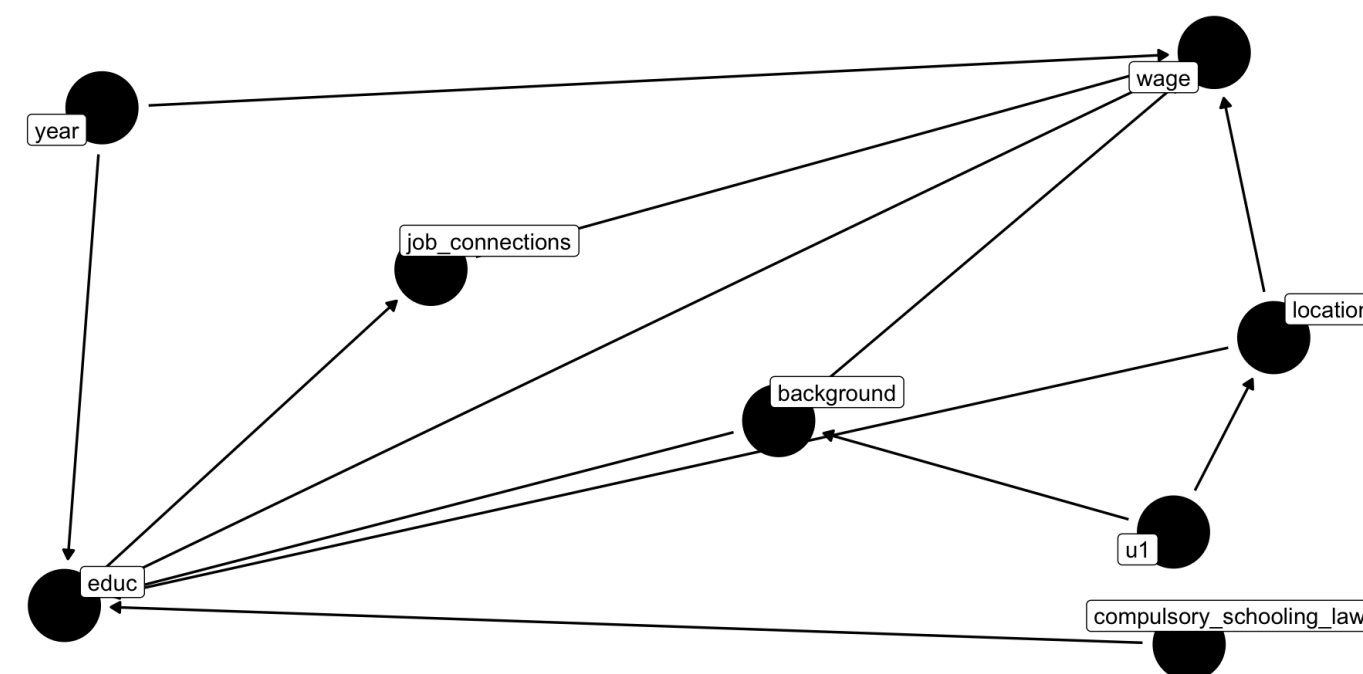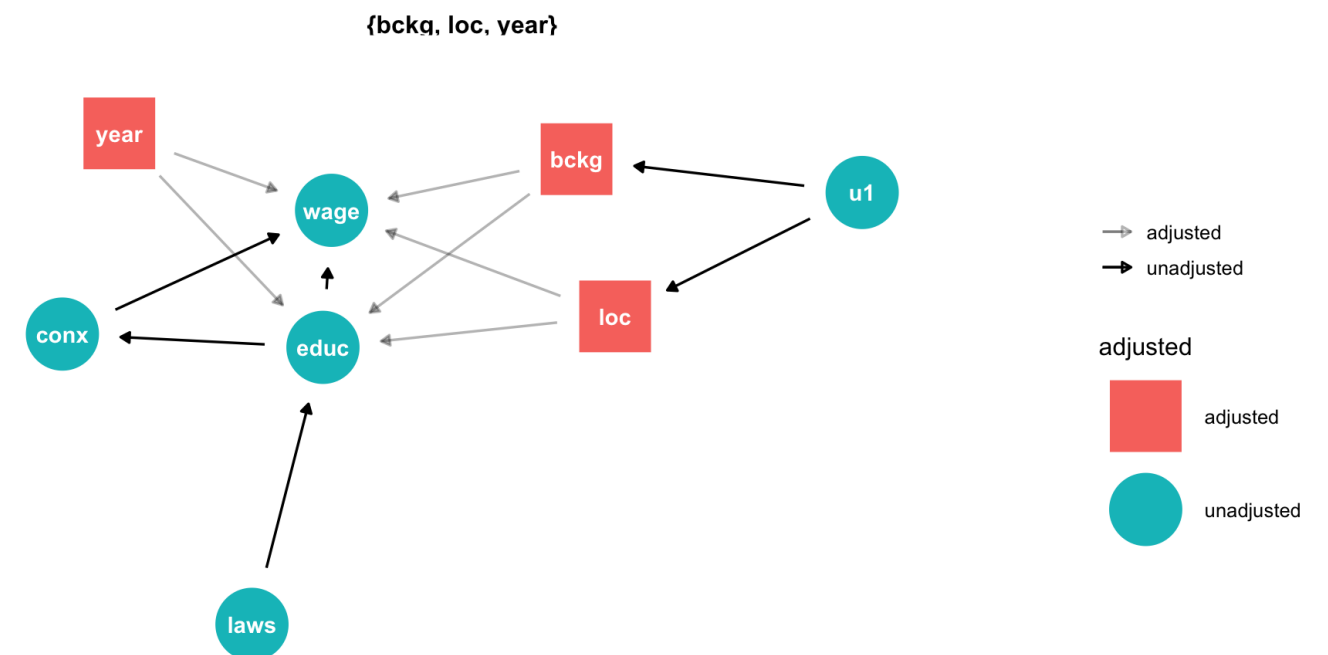
# ggdag: Additional Tools

- If you have defined X (`exposure`) and Y (`outcome`), you can use `ggdag_adjustment_set()` to have it show you what you need to control for in order to identify $X \rightarrow Y$!

```
1   dagify(wage ~ educ + conx + year + bckg + loc,
2          educ ~ bckg + year + loc + laws,
3          conx ~ educ,
4          bckg ~ u1,
5          loc ~ u1,
6          exposure = "educ",
7          outcome = "wage"
8          ) %>%
9   ggdag_adjustment_set(shadow = T)+ #<<
10  theme_dag()
```

# **ggdag**: **Additional Tools**

- You can also use `impliedConditionalIndependencies()` from the `dagitty` package to have it show the testable implications from dagitty.net

```
1  dagify(wage ~ educ + conx + year + bckg + loc,
2         educ ~ bckg + year + loc + laws,
3         conx ~ educ,
4         bckg ~ u1,
5         loc ~ u1,
6         exposure = "educ",
7         outcome = "wage"
8         ) %>%
9     impliedConditionalIndependencies() #<<
```

```
bckg _||_ conx | educ
bckg _||_ laws
bckg _||_ loc | u1
bckg _||_ year
conx _||_ laws | educ
conx _||_ loc | educ
conx _||_ u1 | bckg, loc
conx _||_ u1 | educ
conx _||_ year | educ
educ _||_ u1 | bckg, loc
laws _||_ loc
laws _||_ u1
laws _||_ wage | bckg, educ, loc, year
laws _||_ year
```

# DAG Rules

# DAG Rules

- How does dagitty.net and ggdag know how to identify effects, or what to control for, or what implications are testable?

- Comes from fancy math called "do-calculus"

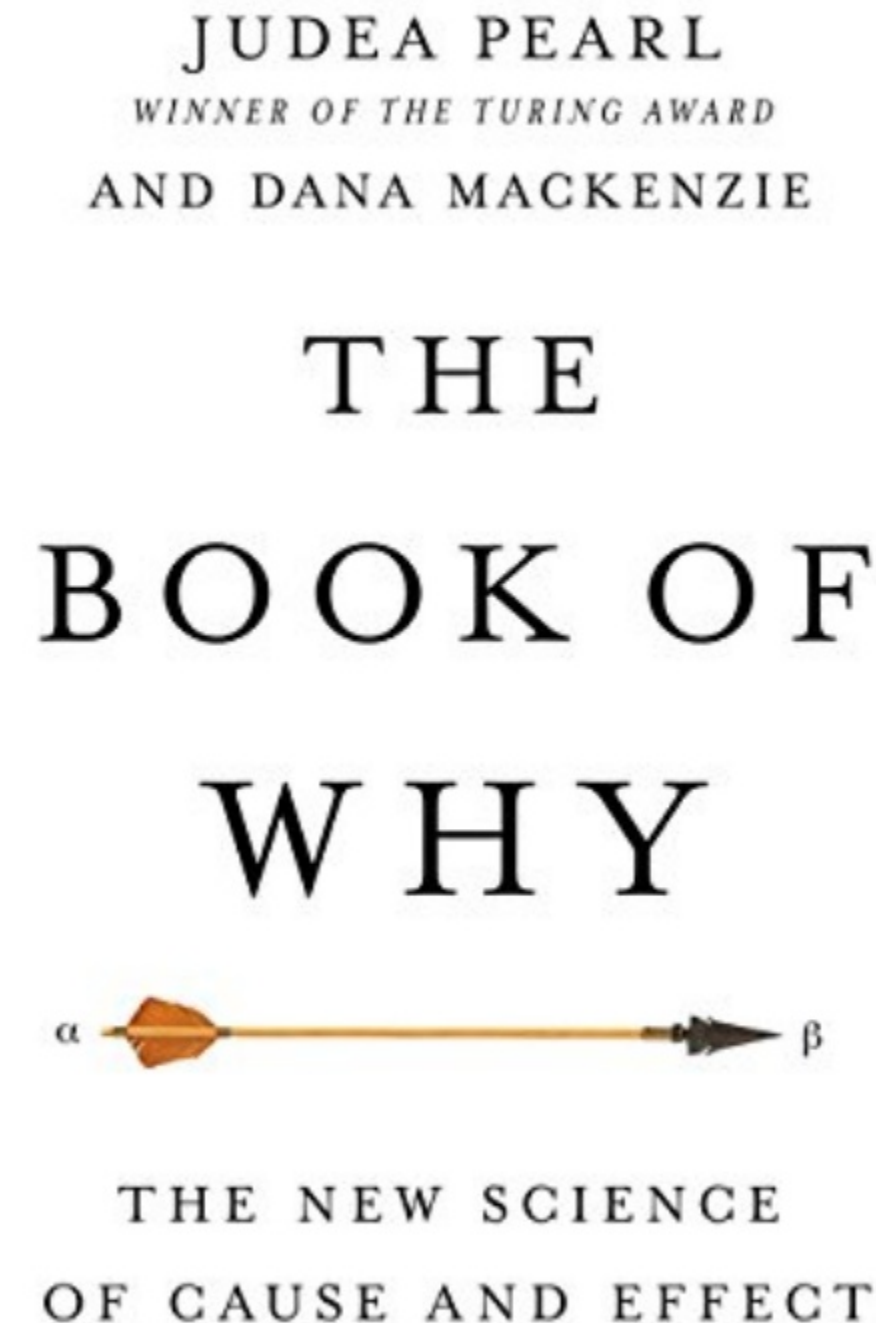**The do-calculus**  Let $G$ be a CGM, $G_{\overline{T}}$ represent $G$ post-intervention (i.e with all links into $T$ removed) and $G_{\underline{T}}$ represent $G$ with all links *out of* $T$ removed. Let $do(t)$ represent intervening to set a single variable $T$ to $t$,

**Rule 1:**  $\mathbb{P}(y|do(t), z, w) = \mathbb{P}(y|do(t), z)$ if $Y \perp\!\!\!\perp W|(Z, T)$ in $G_{\overline{T}}$

**Rule 2:**  $\mathbb{P}(y|do(t), z) = \mathbb{P}(y|t, z)$ if $Y \perp\!\!\!\perp T|Z$ in $G_{\underline{T}}$
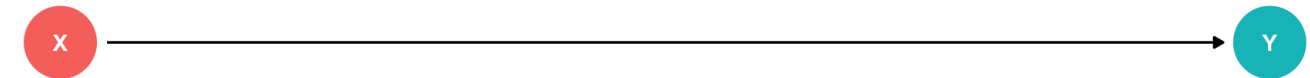
**Rule 3:**  $\mathbb{P}(y|do(t), z) = \mathbb{P}(y|z)$ if $Y \perp\!\!\!\perp T|Z$ in $G_{\overline{T}}$, and $Z$ is not a decedent of $T$.

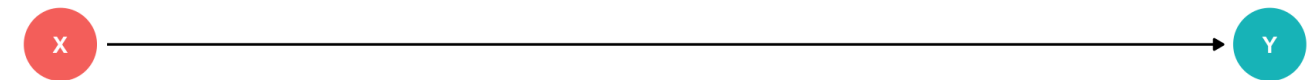- Fortunately, these amount to a few simple rules that we can see on a DAG

# DAGs I

- Typical notation:
- $X$ is independent variable of interest
  - Epidemiology: **"intervention"** or **"exposure"**
- $Y$ is dependent or **"response"** variable
- Other variables use other letters
- You can of course use words instead of letters!

# DAGs and Causal Effects

- Arrows indicate causal effect (& direction)

- Two types of causal effect:
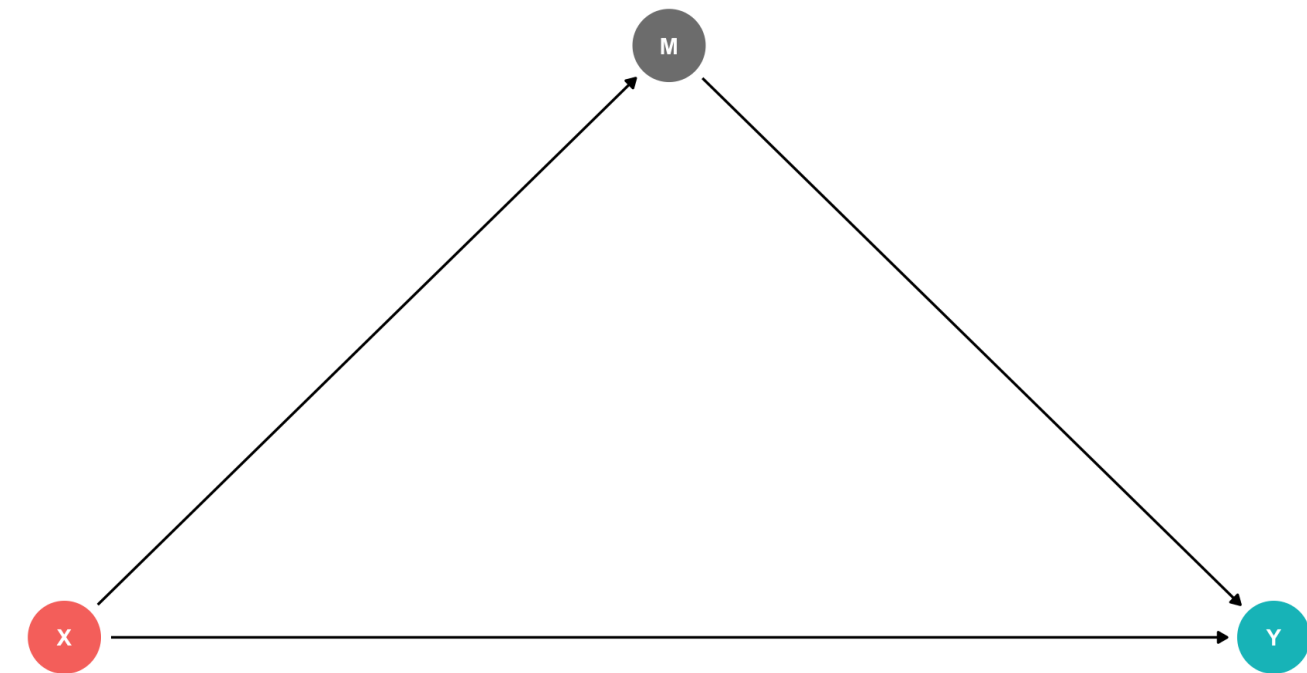
1. Direct effects: $X \rightarrow Y$

# DAGs and Causal Effects

- Arrows indicate causal effect (& direction)

- Two types of causal effect:

1. Direct effects: $X \rightarrow Y$

2. Indirect effects: $X \rightarrow M \rightarrow Y$

- $M$ is a **"mediator"** variable, the **mechanism** by which $X$ affects $Y$
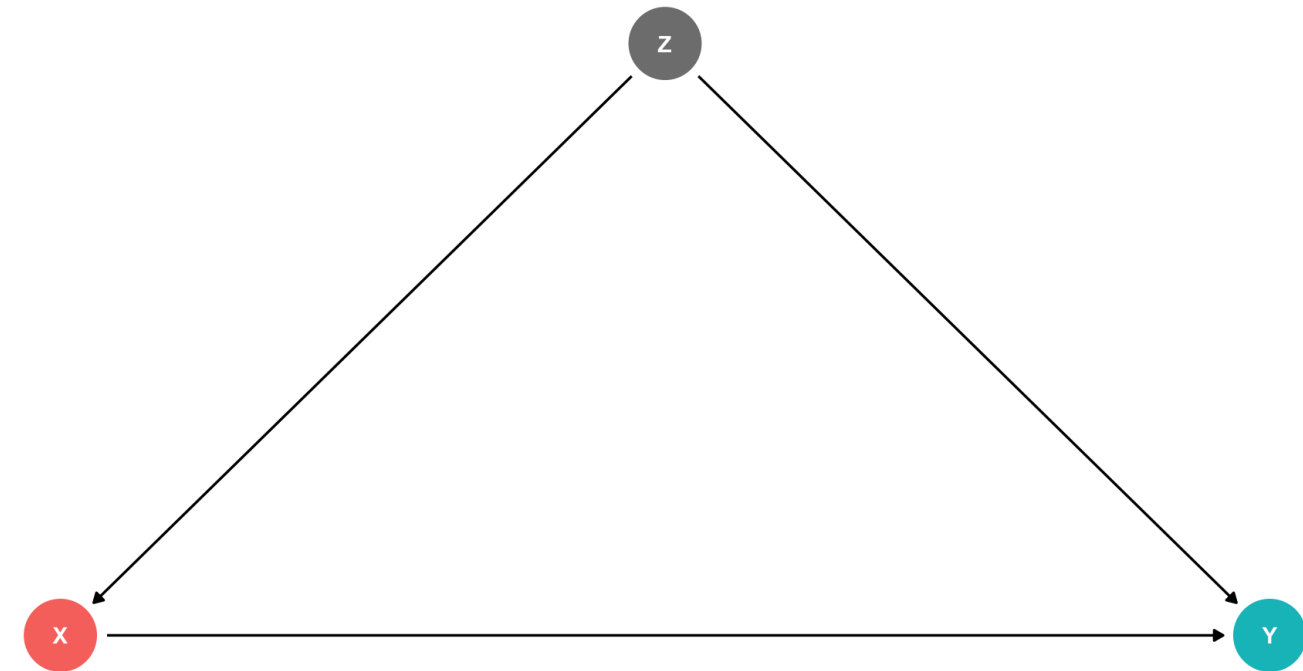
# DAGs and Causal Effects

- Arrows indicate causal effect (& direction)

- Two types of causal effect:

1. Direct effects: $X \rightarrow Y$

2. Indirect effects: $X \rightarrow M \rightarrow Y$

- $M$ is a **"mediator"** variable, the **mechanism** by which $X$ affects $Y$

3. You of course might have both!

# Confounders

- $Z$ is a **"confounder"**: it causes *both* $X$ and $Y$

- $cor(X, Y)$ is made up of two parts:

  1. Causal effect of $(X \rightarrow Y)$ 👍

  2. $Z$ causing both the values of $X$ and $Y$ 👎

- Failing to control for $Z$ will **bias** our estimate of the causal effect of $X \rightarrow Y$!
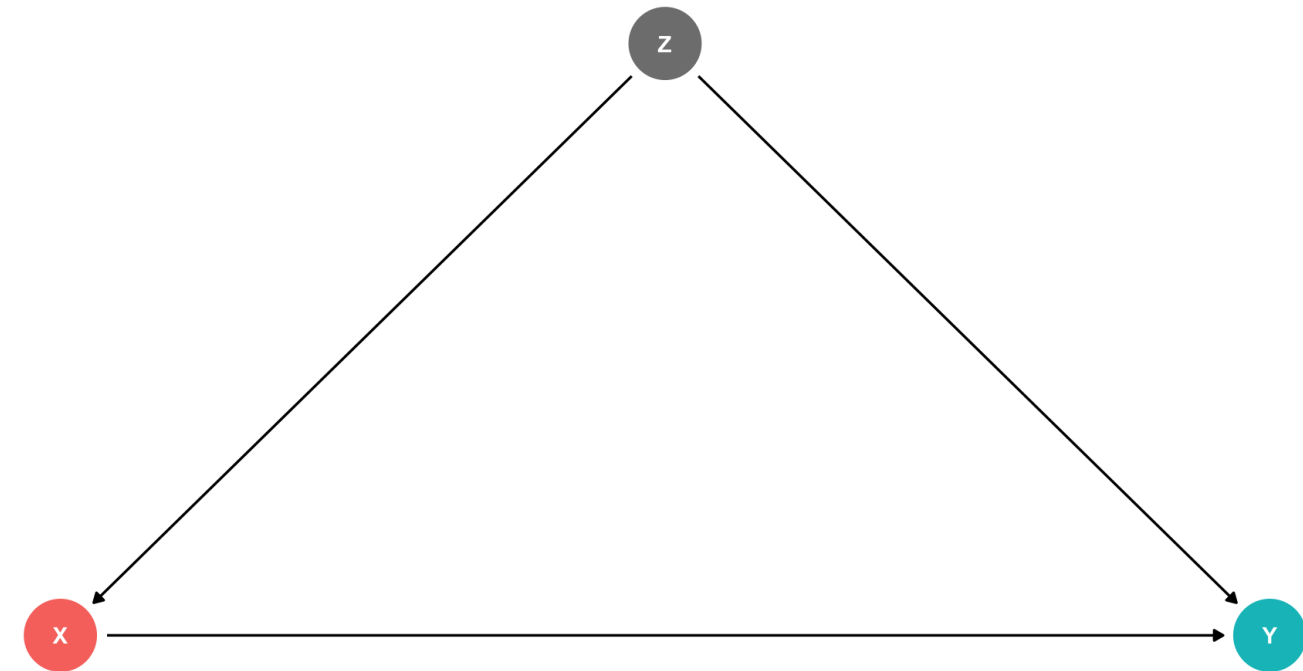
# Confounders

- Confounders are the DAG-equivalent of **omitted variable bias** (next class)

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

- By leaving out $Z_i$, this regression is **biased**

- $\hat{\beta}_1$ picks up *both*:

  1. $X \rightarrow Y$

  2. $X \leftarrow Z \rightarrow Y$

# "Front Doors" and "Back Doors"

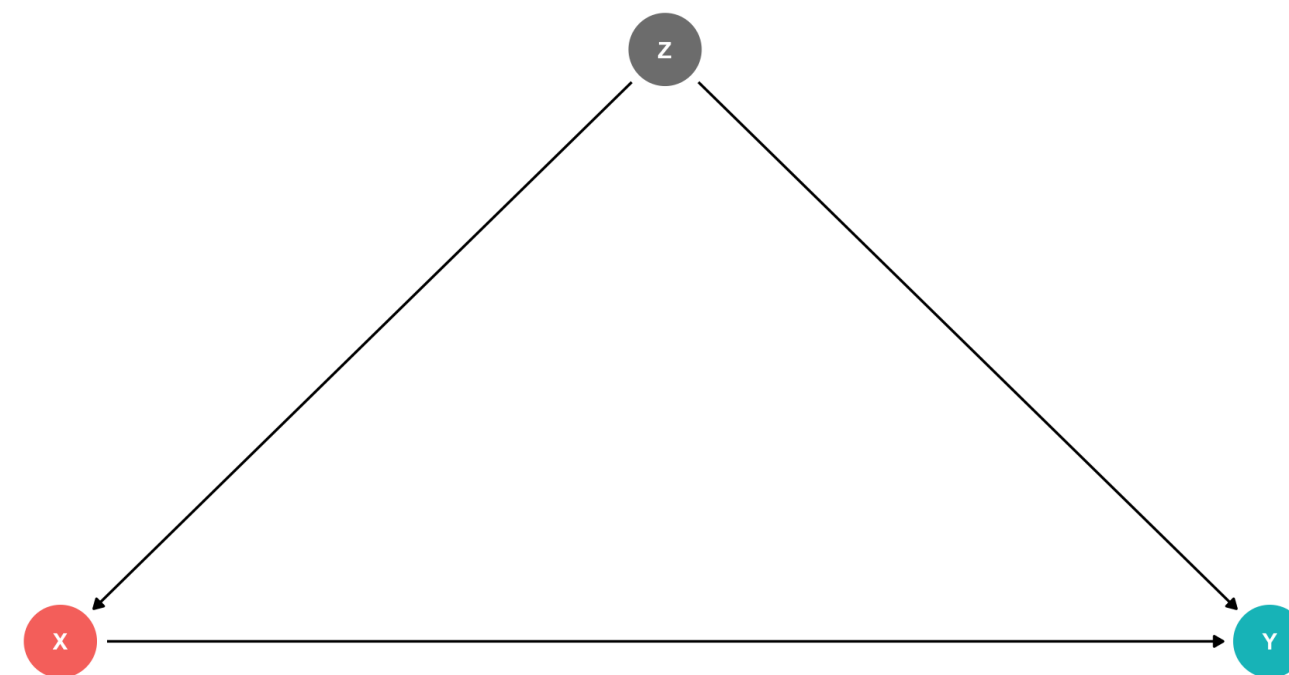- With this DAG, there are 2 paths that connect $X$ and $Y$[1]:

1. A **causal "front-door" path**: $X \rightarrow Y$

- 👍 what we want to measure

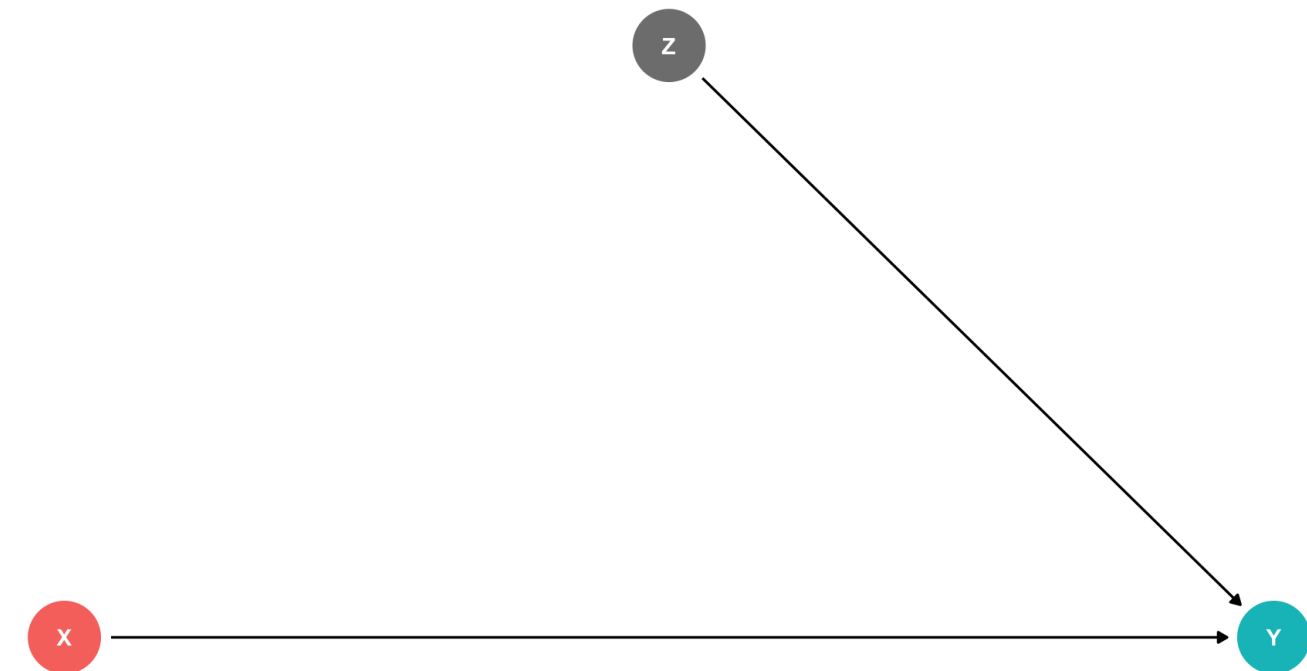2. A **non-causal "back-door" path**: $X \leftarrow Z \rightarrow Y$

- At least one causal arrow runs in the opposite direction

- 👎 adds a confounding bias

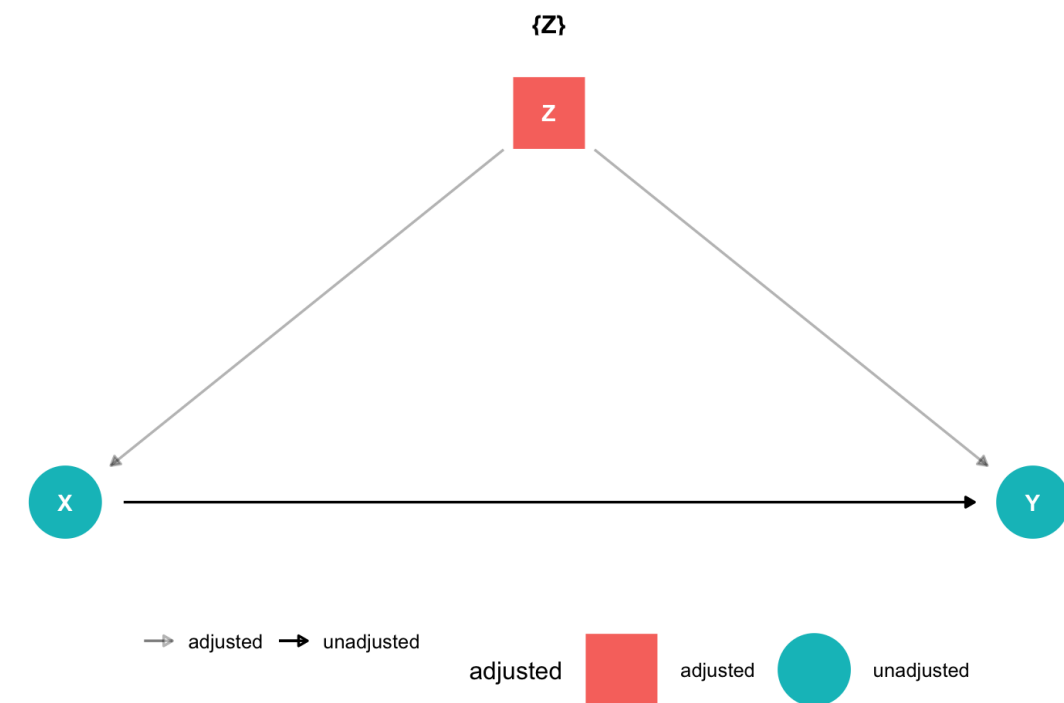1. Regardless of the *directions* of the arrows!

# Controlling I

- Ideally, if we ran a **randomized control trial** and randomly assigned different values of $X$ to different individuals, this would delete the arrow between $Z$ and $X$

  - Individuals' values of $Z$ do not affect whether or not they are treated ($X$)

- This would only leave the front-door, $X \rightarrow Y$
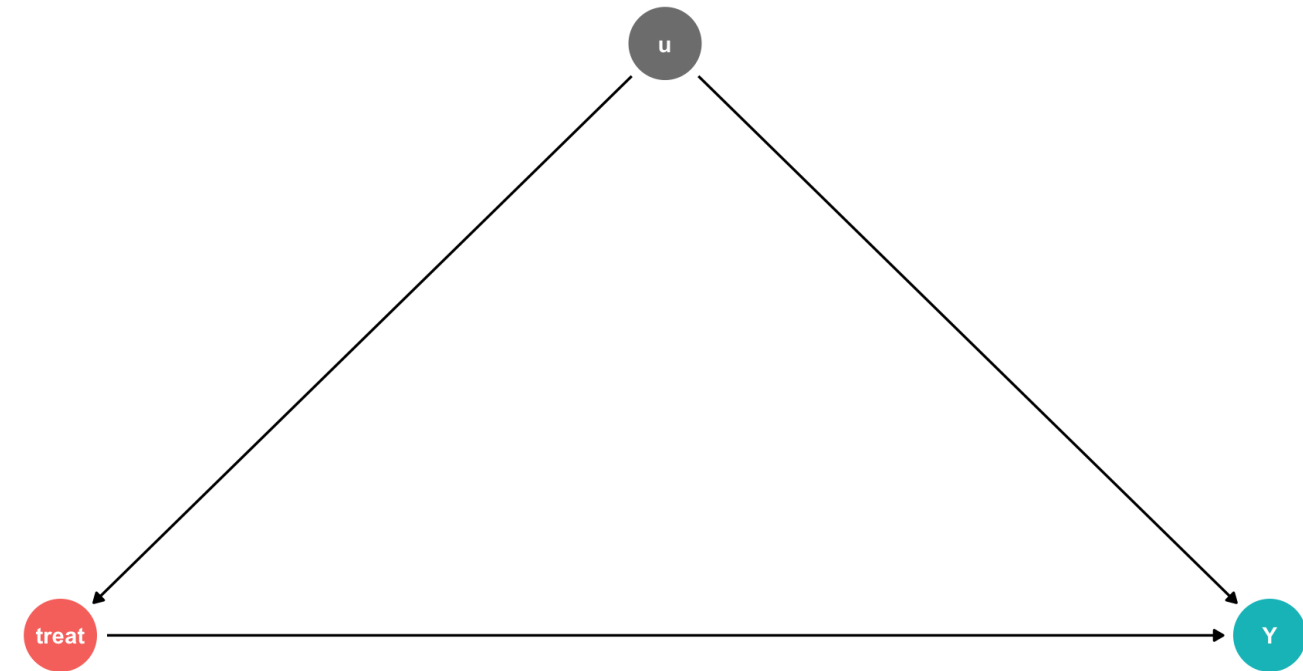
- But we can rarely run an ideal RCT

# Controlling II

- Instead of an RCT, if we can just **"adjust for"** or **"control for"** $Z$, we can *block* the back-door path $X \leftarrow Z \rightarrow Y$

- This would only leave the front-door path open, $X \rightarrow Y$

- "As good as" an RCT!

# Controlling II

- Using our terminology from last class, we have an outcome $(Y)$, and some treatment

- But there are **unobserved factors** $(u)$

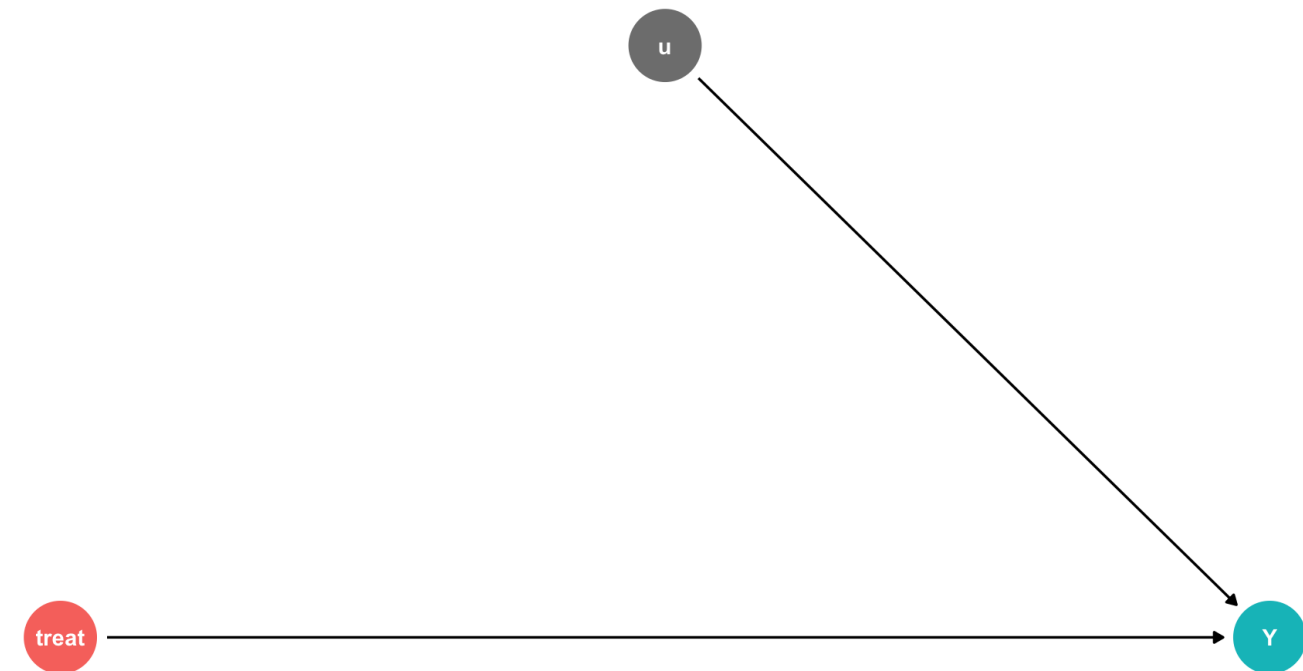$$Y_i = \beta_0 + \beta_1 Treatment + u_i$$

# Controlling II

- Using our terminology from last class, we have an outcome $(Y)$, and some treatment

- But there are **unobserved factors** $(u)$

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

- If we can *randomly* assign treatment, this makes treatment exogenous:
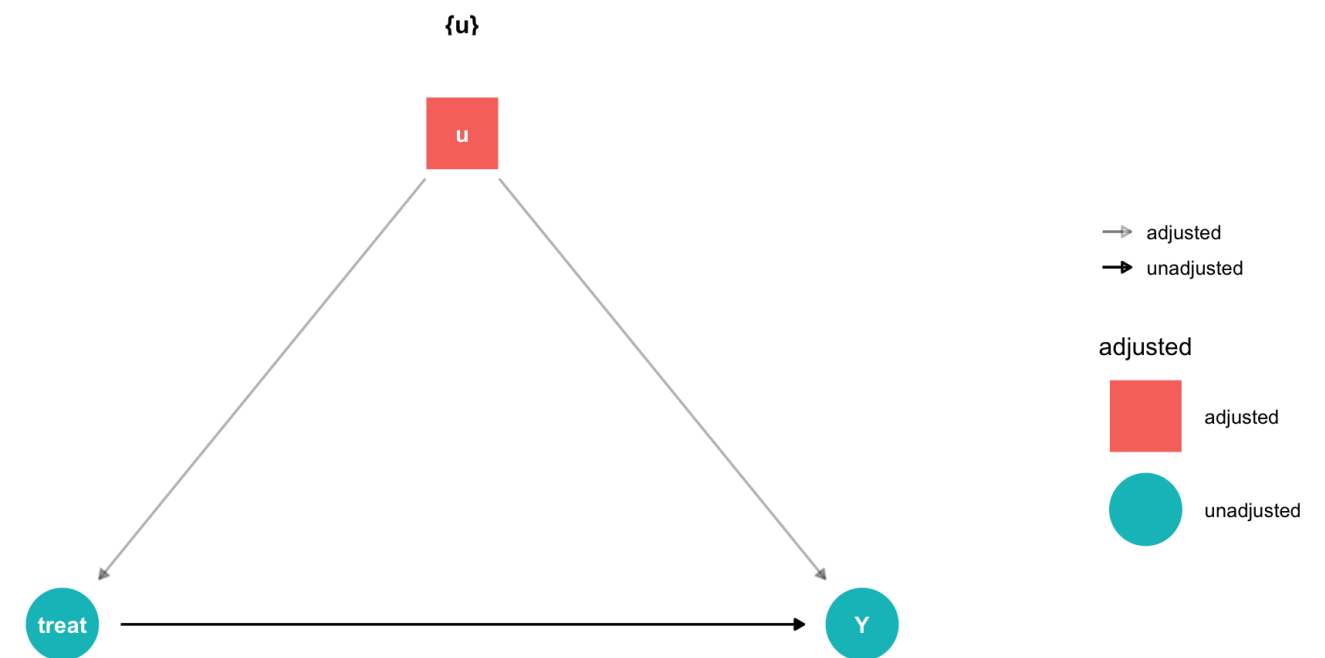
$$cor(treatment, u) = 0$$

# Controlling II

- Using our terminology from last class, we have an outcome $(Y)$, and some treatment

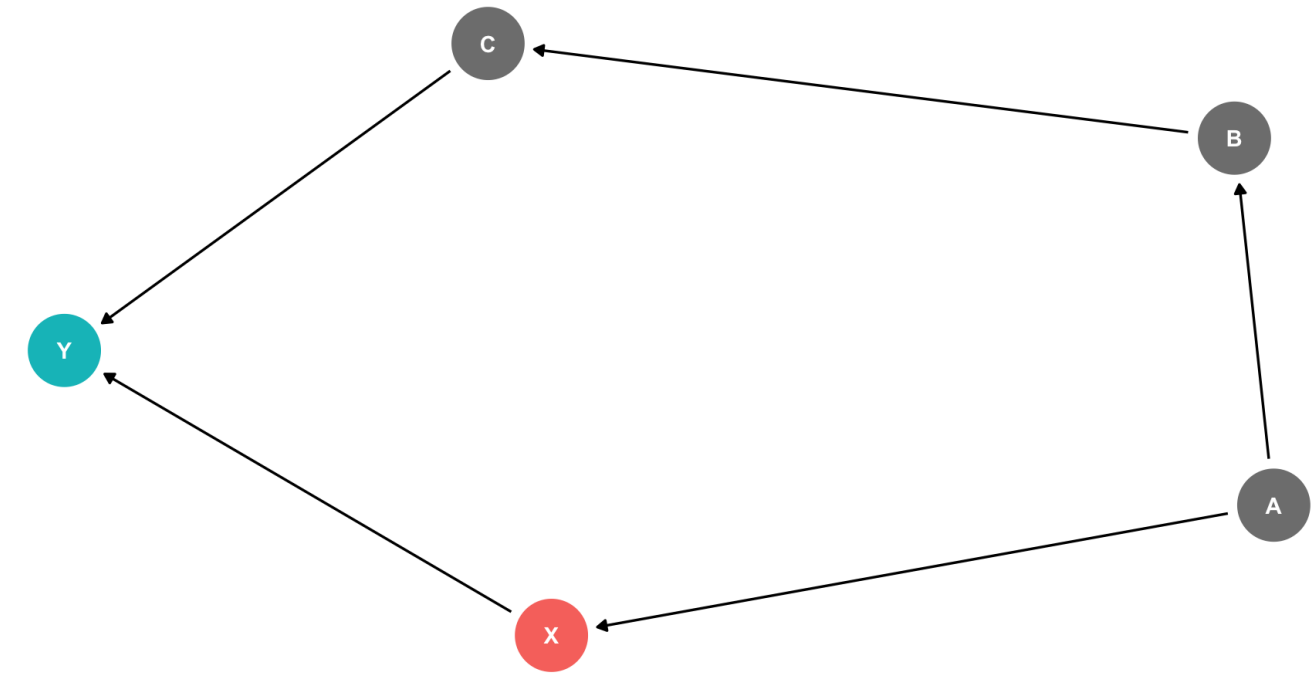- But there are **unobserved factors** $(u)$

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

- When we (often) can't randomly assign treatment, we have to find another way to control for measurable things in $u$

# Controlling II

- Controlling for a single variable along a long causal path is sufficient to block that path!

- Causal path: $X \rightarrow Y$

- Backdoor path: $X \leftarrow A \rightarrow B \rightarrow C \rightarrow Y$

- It is sufficient to block this backdoor by controlling **either** $A$ **or** $B$ **or** $C$!
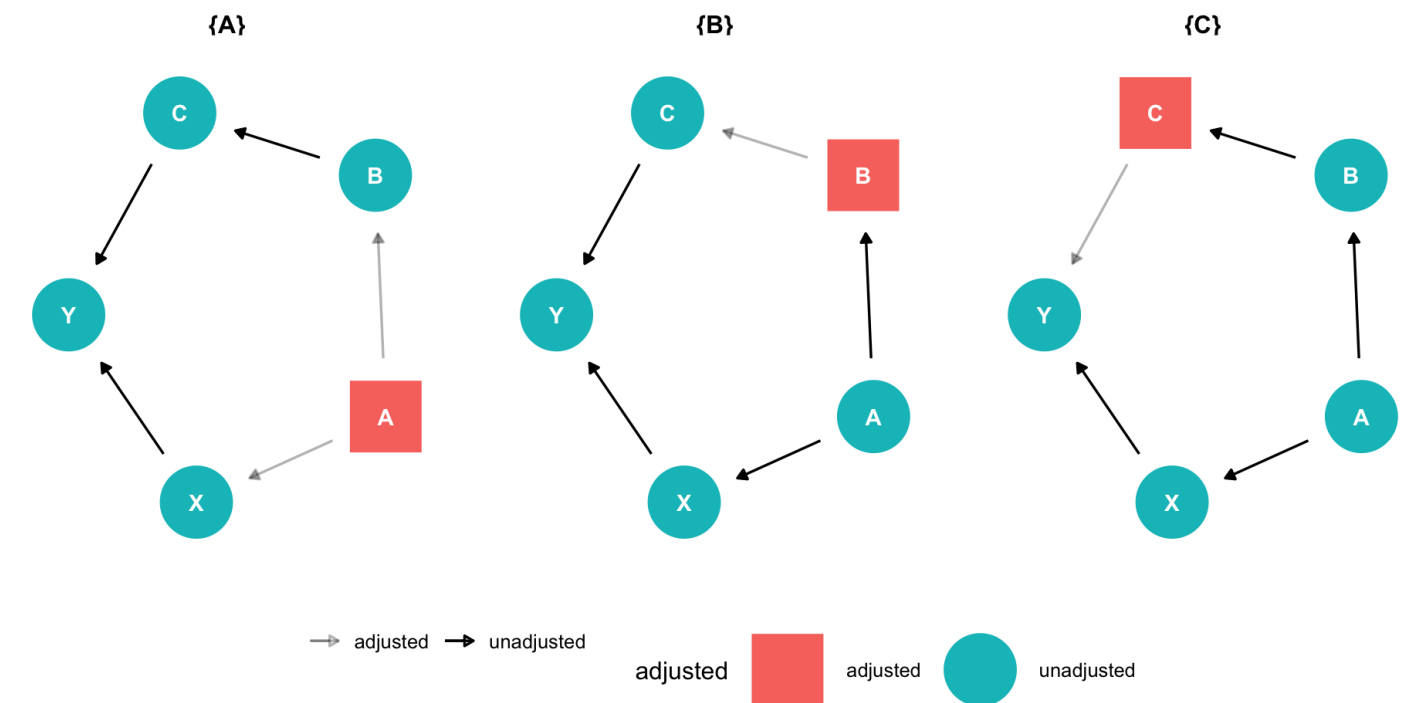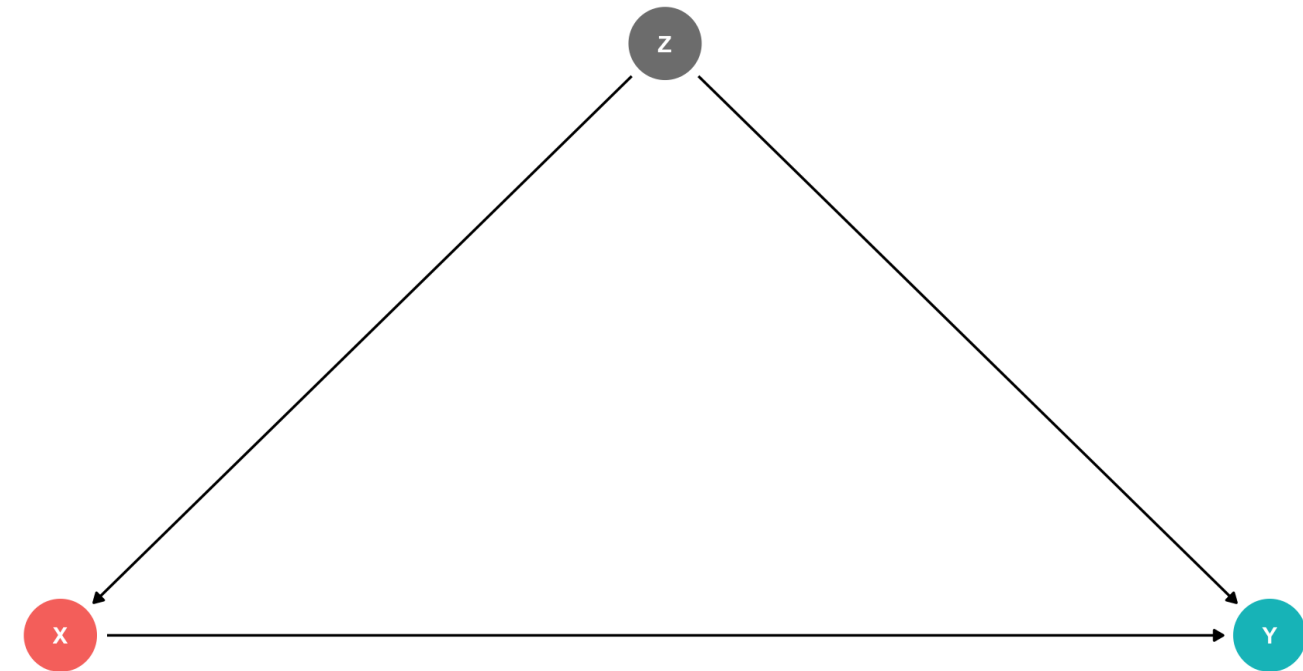
# Controlling II

- Controlling for a single variable along a long causal path is sufficient to block that path!

- Causal path: $X \rightarrow Y$

- Backdoor path: $X \leftarrow A \rightarrow B \rightarrow C \rightarrow Y$

- It is sufficient to block this backdoor by controlling **either** $A$ **or** $B$ **or** $C$!

# The Back Door Criterion

- To **identify** the causal effect of $X \rightarrow Y$:

- **"Back-door criterion"**: control for the minimal amount of variables sufficient to ensure that **no open back-door exists** between $X$ and $Y$

- Example: in this DAG, control for $Z$

# The Back Door Criterion

- Implications of the Back-door criterion:

1. You *only* need to control for the variables that keep a back-door open, *not all other variables!*

**Example**:

- $X \rightarrow Y$ (front-door)
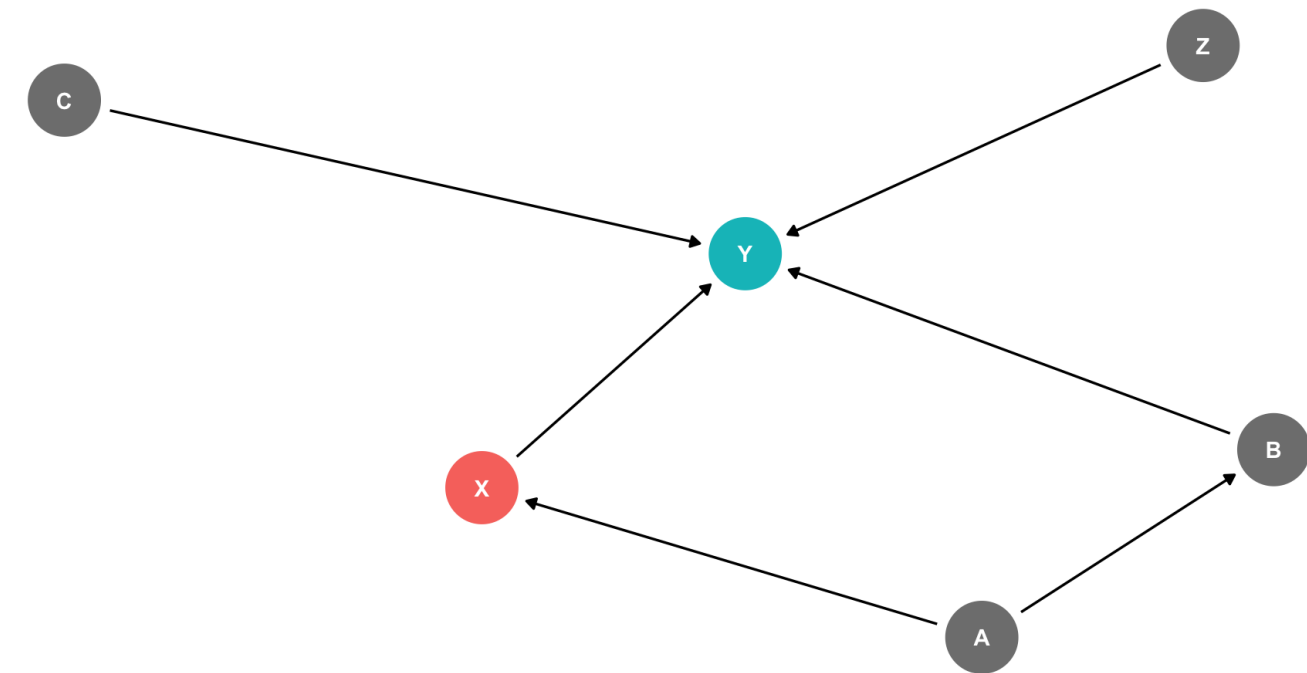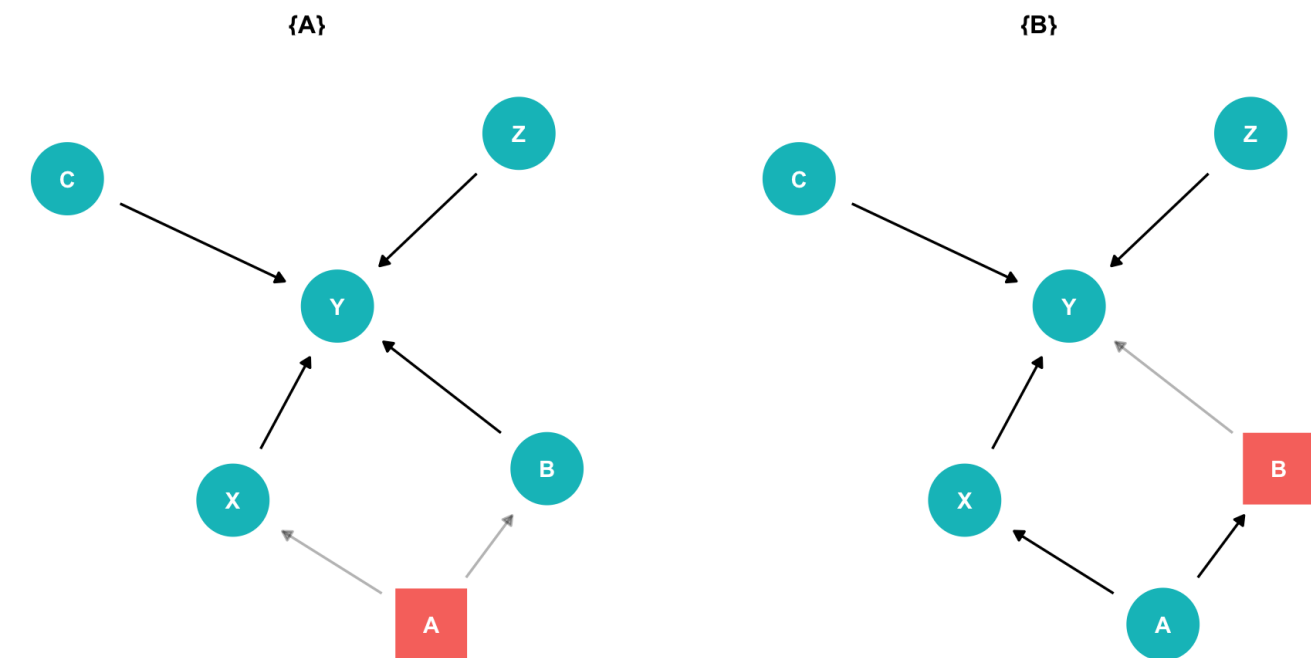- $X \leftarrow A \rightarrow B \rightarrow Y$ (back-door)

# The Back Door Criterion

- Implications of the Back-door criterion:

1. You *only* need to control for the variables that keep a back-door open, *not all other variables!*
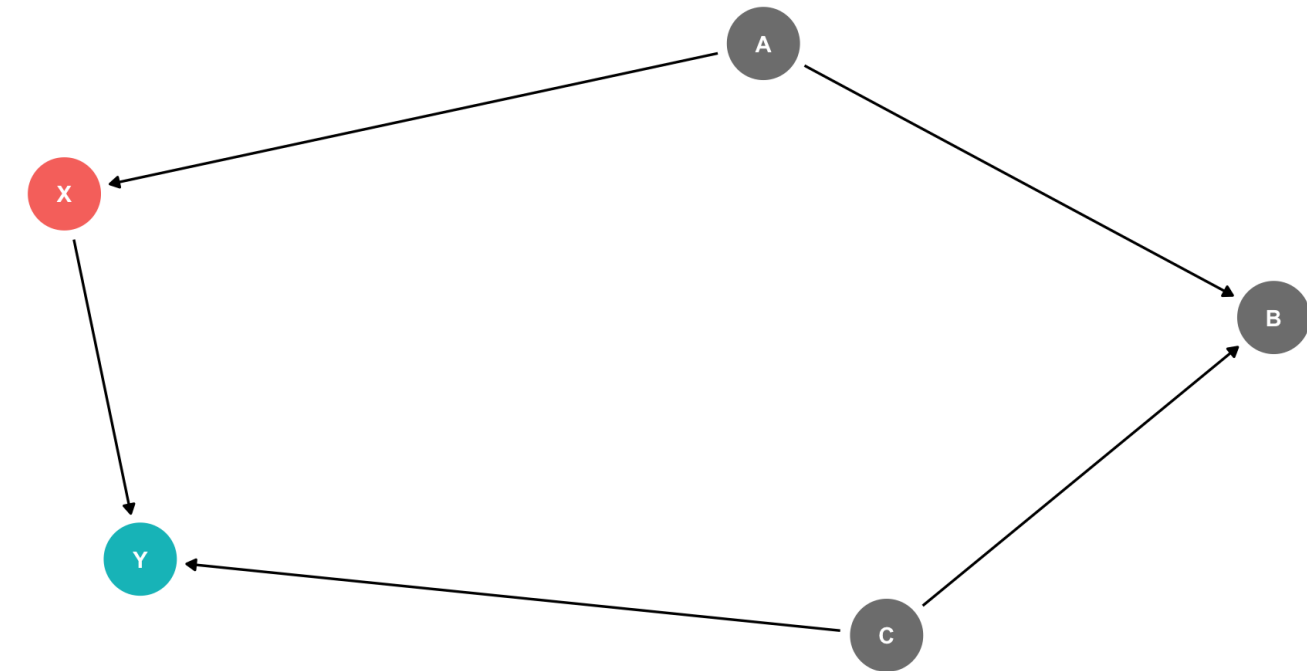
**Example**:

- $X \rightarrow Y$ (front-door)

- $X \leftarrow A \rightarrow B \rightarrow Y$ (back-door)

- Need only control for $A$ *or* $B$ to block the back-door path

- $C$ and $Z$ have no effect on $X$, and therefore we don't need to control for them!

# The Back Door Criterion: Colliders

2. Exception: the case of a **"collider"**

- If arrows "collide" at a node, **that node is automatically blocking the pathway**, do not control for it!

- Controlling for a collider would *open* the path and **add bias!**

**Example**:

- $X \rightarrow Y$ (front-door)
- $X \leftarrow A \rightarrow B \leftarrow C \rightarrow Y$ (back-door, but **blocked by B!**)

# The Back Door Criterion: Colliders

2. Exception: the case of a **"collider"**

- If arrows "collide" at a node, **that node is automatically blocking the pathway**, do not control for it!

- Controlling for a collider would *open* the path and **add bias!**

**Example**:

- $X \rightarrow Y$ (front-door)

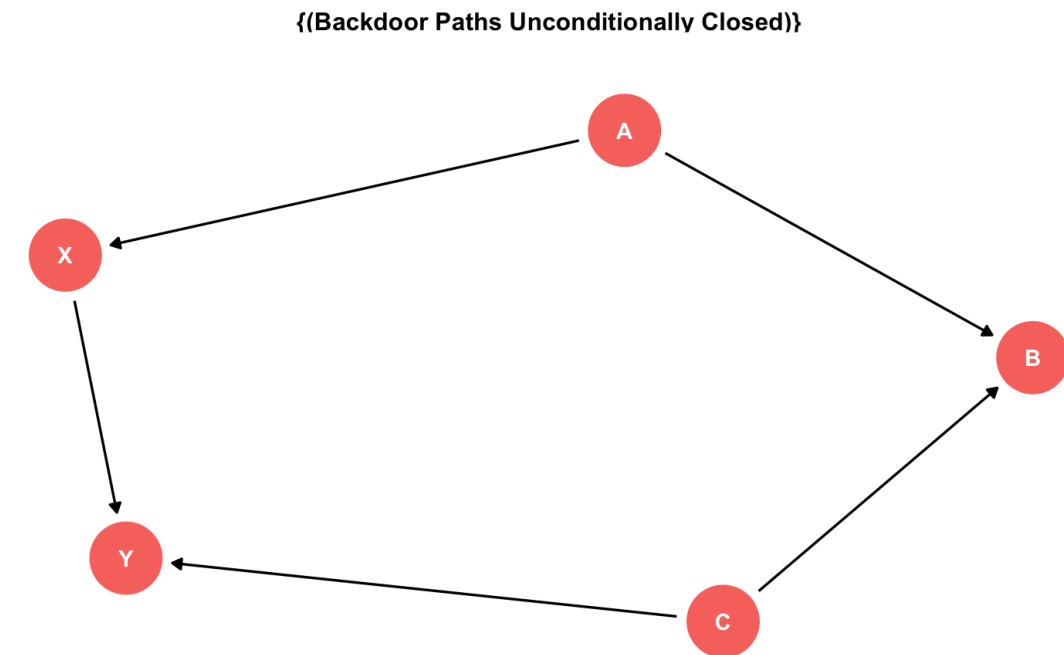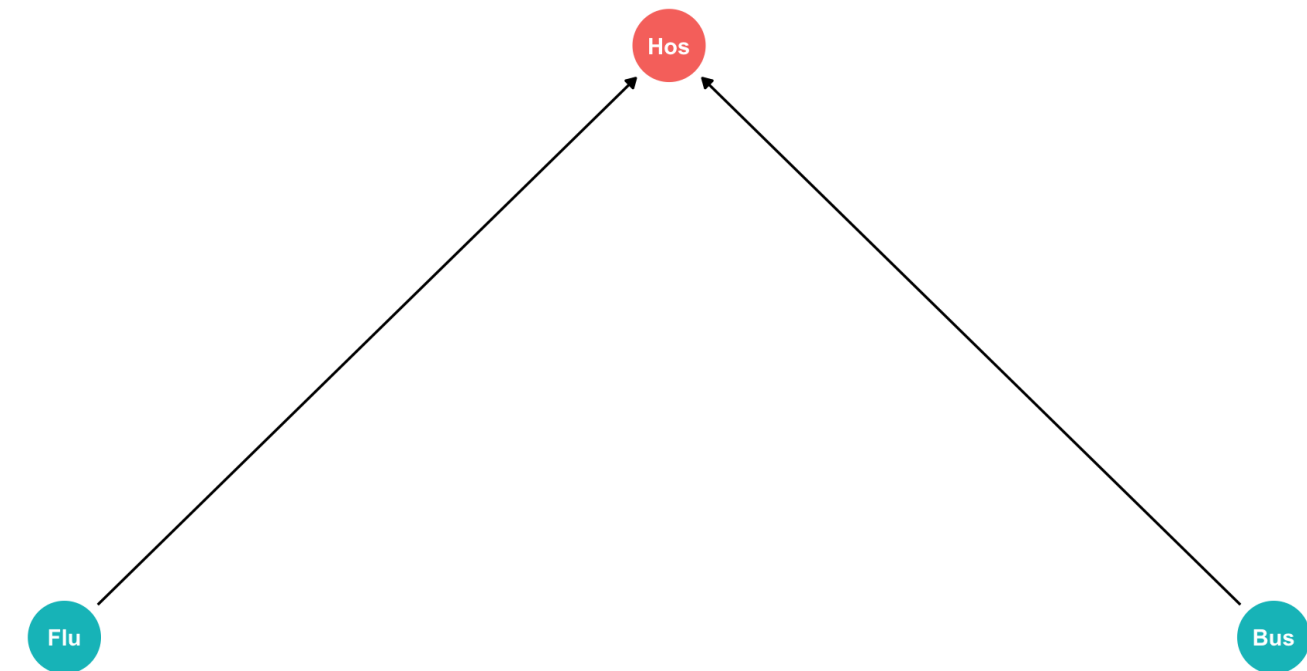- $X \leftarrow A \rightarrow B \leftarrow C \rightarrow Y$ (back-door, but **blocked by B!**)

- Don't need to control for anything here!

{(Backdoor Paths Unconditionally Closed)}
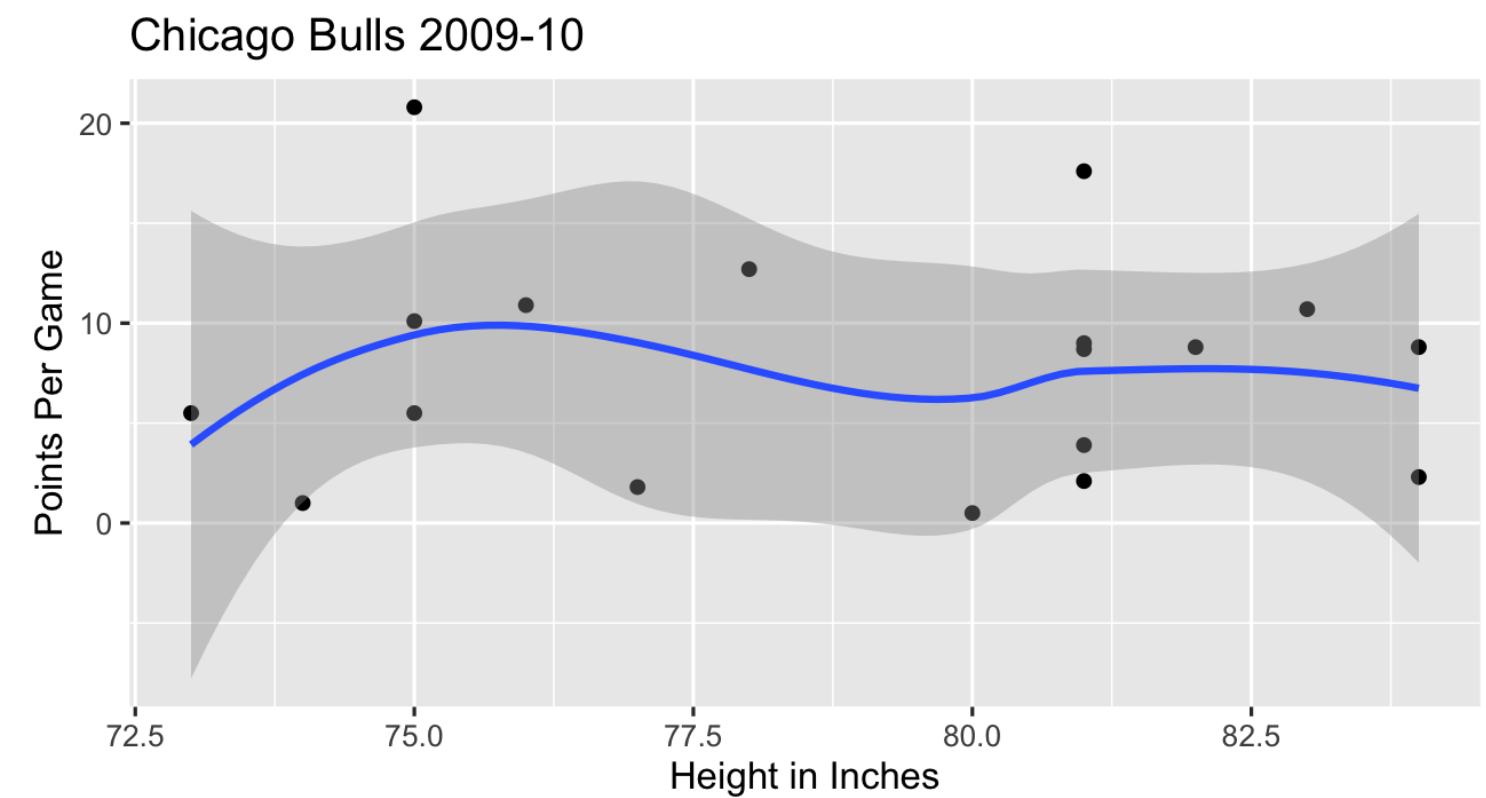
# The Back Door Criterion: Colliders Example

Are you less likely to get the flu if you are hit by
a bus?

- $Flu$: getting the flu

- $Bus$: being hit by a bus

- $Hos$: being in the hospital

- Both $Flu$ and $Bus$ send you to $Hos$ (arrows)

- Conditional on being in $Hos$, negative
  correlation between $Flu$ and $Bus$ (spurious!)

# The Back Door Criterion: Colliders Example

- In the NBA, apparently players' height has no relationship to points scored?
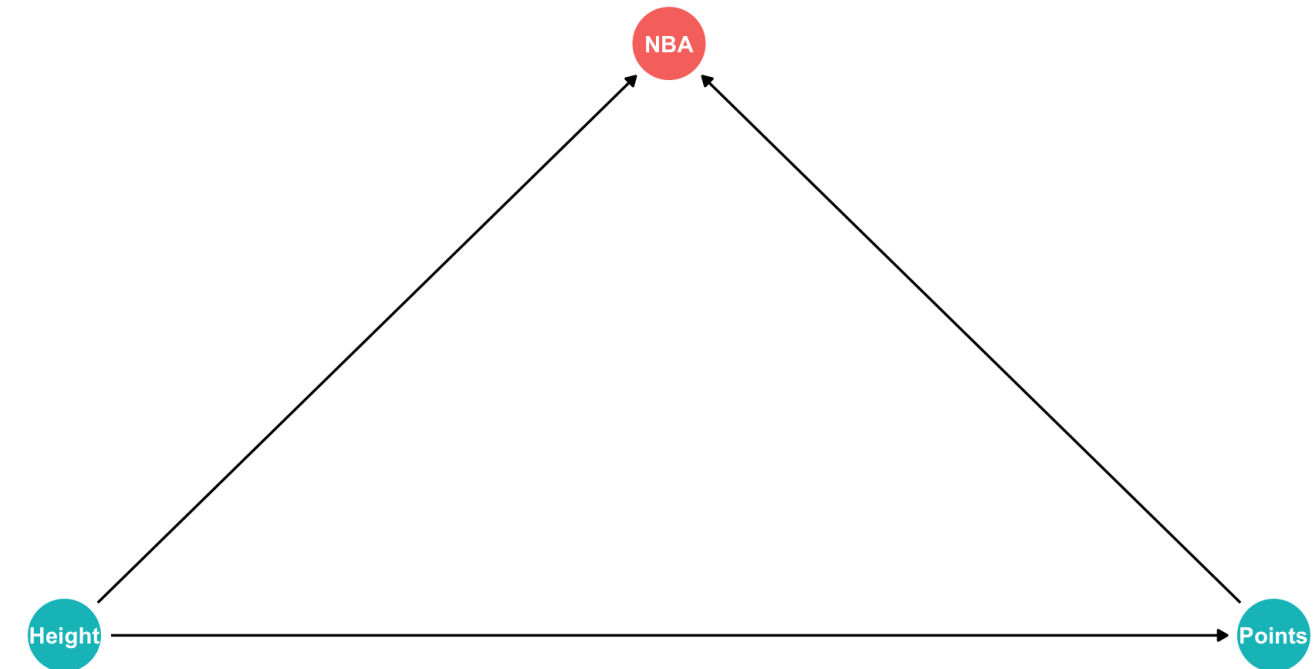


Chicago Bulls 2009-10

# The Back Door Criterion: Colliders Example

- **In the NBA**, apparently players' height has no relationship to points scored?
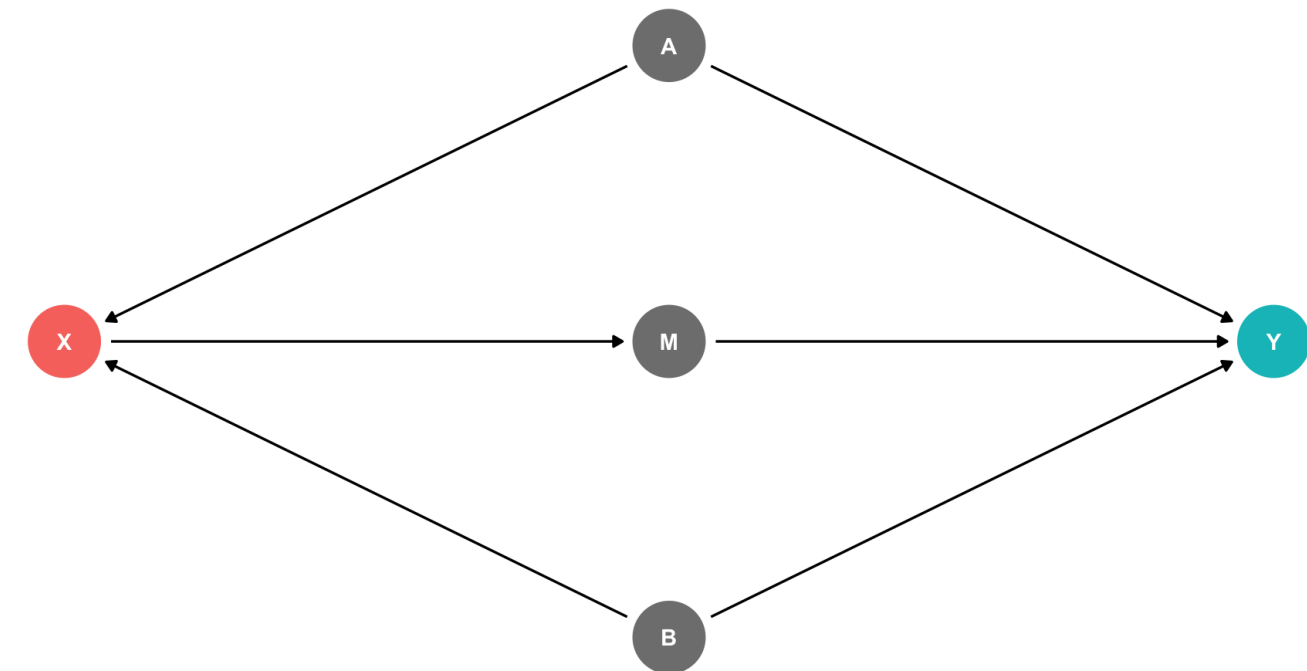
# The Front Door Criterion: Mediators I

- Another case where controlling for a variable actually *adds bias* is if that variable is known as a **"mediator"**.

Example

- $X \rightarrow M \rightarrow Y$ (front-door)

- $X \leftarrow A \rightarrow Y$ (back-door)

- $X \leftarrow B \rightarrow Y$ (back-door)

- Should we control for $M$?

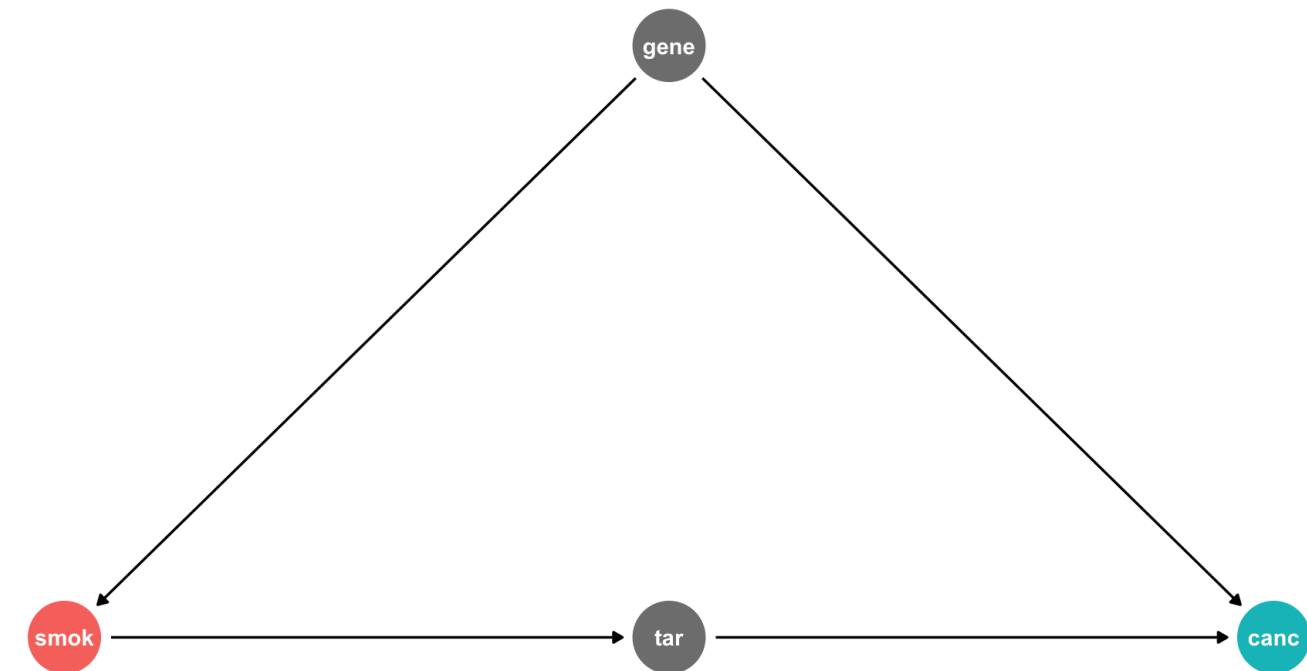- If we did, this would block the front-door!

# The Front Door Criterion: Mediators II

# The Front Door Criterion: Mediators III

- Tobacco industry claimed that $cor(smoking, cancer)$ could be spurious due to a confounding gene that affects both!

  - Smoking gene is unobservable

- Suppose smoking causes `tar` buildup in lungs, which cause `cancer`

- We should *not* control for `tar`, it's on the **front-door path**

  - This is how scientific studies can relate smoking to cancer

# Summary: DAG Rules for Causal Identification

Thus, to achieve **causal identification**, control for the minimal amount of variables such that:

1. Ensure **no back-door path remains open**

- Close back-door paths by *controlling* for any one variable along that path

- Colliders along a path *automatically* close that path

2. Ensure **no front-door path is closed**

- Do not control for mediators