

# Causal Inference and Counterfactual Reasoning

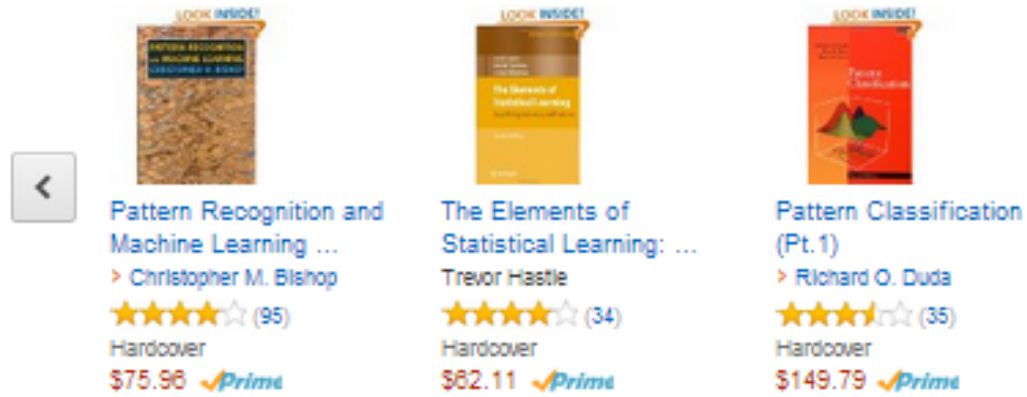
Emre Kıcıman and Amit Sharma

emrek@microsoft.com, amshar@microsoft.com

[Causal Inference and Counterfactual Reasoning](#) at Microsoft Research

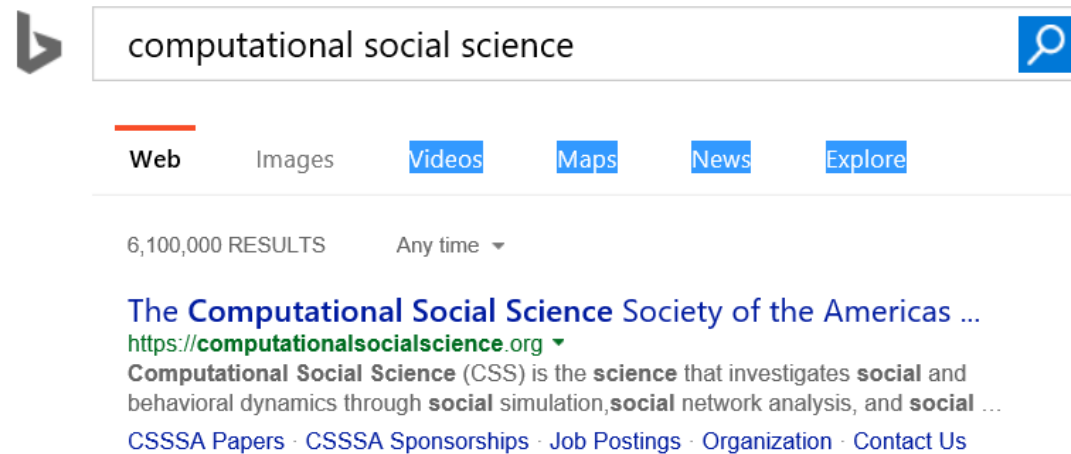
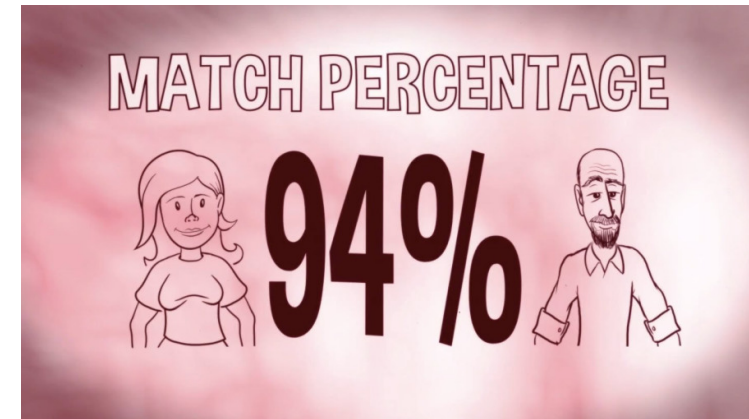
# Predictive systems are impacting our lives

## Customers Who Bought This Item Also Bought



A screenshot of an Amazon 'Customers Who Bought This Item Also Bought' section. It features three book covers with their titles, authors, star ratings, and prices. Each item has a 'Prime' logo next to its price.

Book Title	Author	Rating	Price
Pattern Recognition and Machine Learning ...	Christopher M. Bishop	★★★★☆ (95)	\$75.96
The Elements of Statistical Learning: ...	Trevor Hastie	★★★★☆ (34)	\$82.11
Pattern Classification (Pt. 1)	Richard O. Duda	★★★★☆ (35)	\$149.79



A screenshot of a search engine results page. The search bar contains the text 'computational social science'. Below the search bar are navigation tabs for 'Web', 'Images', 'Videos', 'Maps', 'News', and 'Explore'. The search results show '6,100,000 RESULTS' and a dropdown menu set to 'Any time'. The first result is for 'The Computational Social Science Society of the Americas ...' with a URL and a brief description of the field.

computational social science

Web Images Videos Maps News Explore

6,100,000 RESULTS Any time ▾

**The Computational Social Science Society of the Americas ...**  
<https://computationalsocialscience.org>  
Computational Social Science (CSS) is the science that investigates social and behavioral dynamics through social simulation, social network analysis, and social ...  
CSSSA Papers · CSSSA Sponsorships · Job Postings · Organization · Contact Us









# Why should we care about causality?

We have increasing amounts of data and highly accurate predictions.

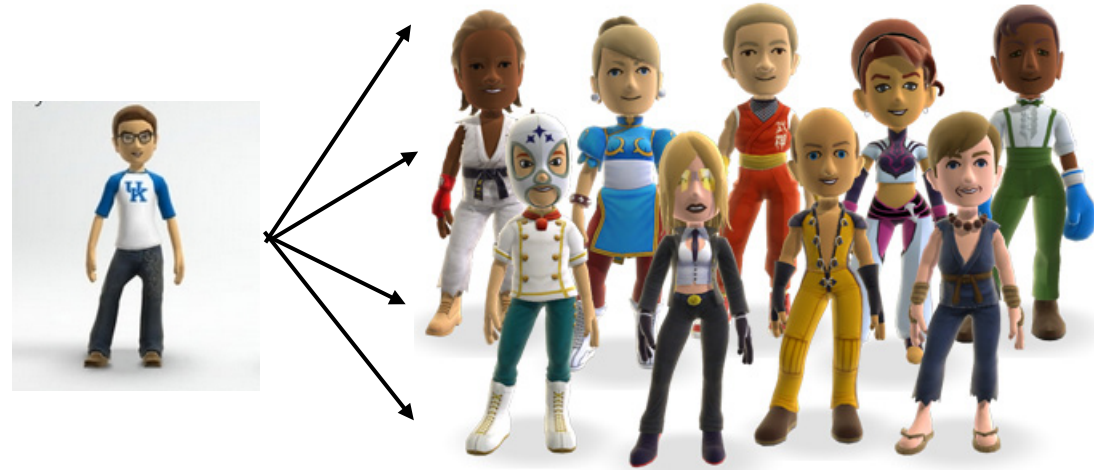
How is causal inference useful?

1) Do prediction models guide decision-making?



# From data to prediction

Can we predict a user's future activity based on exposure to their social feed?



Use the social feed to predict a user's future activity.

- Future Activity  $\rightarrow f(\text{items in social feed}) + \epsilon$

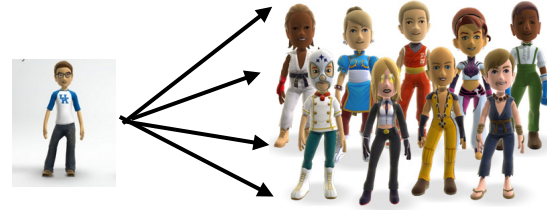
Highly predictive model.

Does it mean that feeds are influencing us significantly?

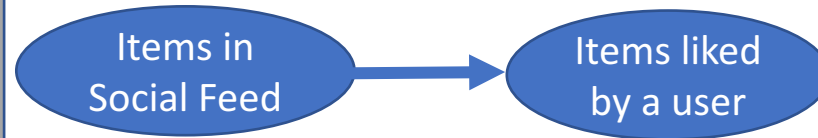
# From prediction to decision-making

Would changing what people see in the feed affect what a user likes?

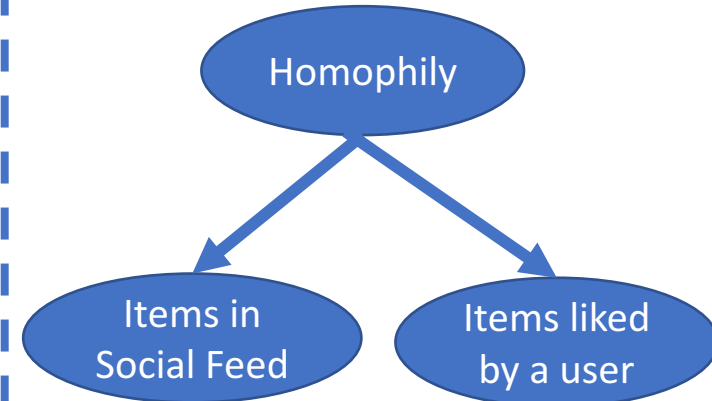
Maybe, maybe not (!)



Predictability due to **feed influence**



Predictability due to **homophily**



Friends' activity can predict a person's activity with high accuracy.  
But that tells us *nothing* about the effect of the social feed.

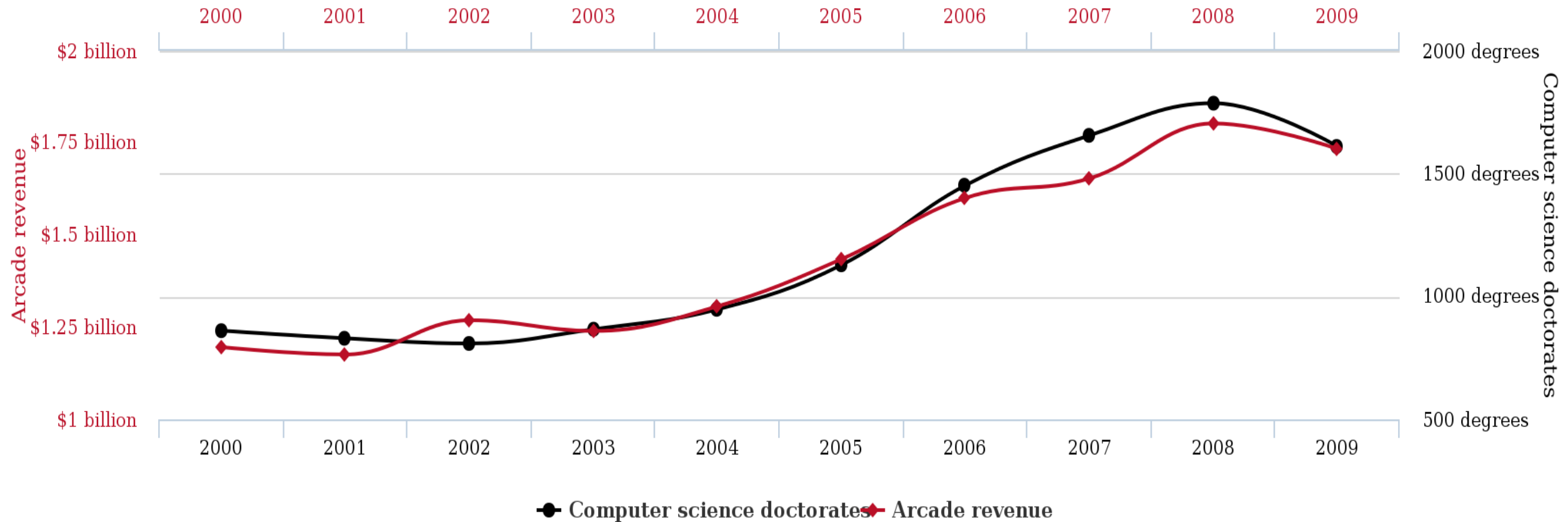


2) Will the predictions be robust tomorrow, or in new contexts?

# Total revenue generated by arcades

correlates with

## Computer science doctorates awarded in the US

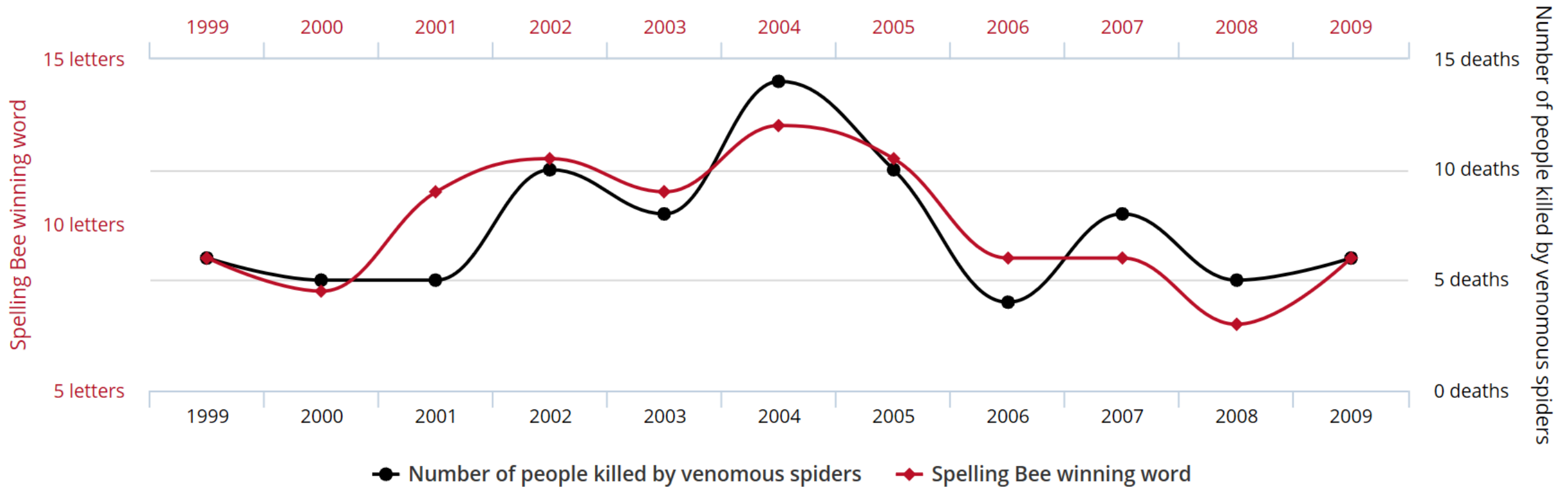




# Letters in Winning Word of Scripps National Spelling Bee correlates with

## Number of people killed by venomous spiders

Correlation: 80.57% (r=0.8057)



tylervigen.com

Data sources: National Spelling Bee and Centers for Disease Control & Prevention

# Story: London Taxi Drivers

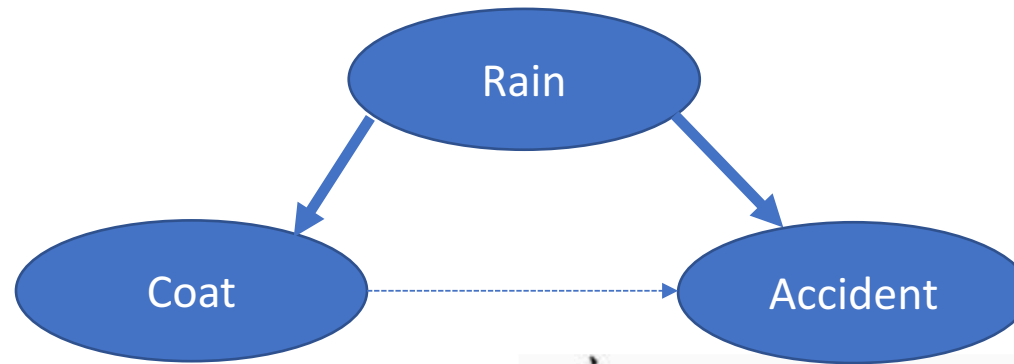
## ◆ Examples:

**London taxi drivers:** A survey has pointed out a positive and significant correlation between the number of accidents and wearing coats. They concluded that coats could hinder movements of drivers and be the cause of accidents. A new law was prepared to prohibit drivers from wearing coats when driving.



**Decision based on the causality ?**





◆ Examples:

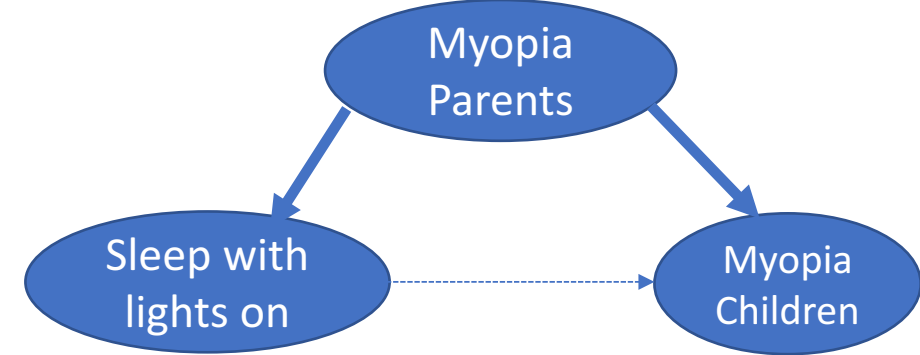
**London taxi drivers:** A survey has pointed out a positive and significant correlation between the number of accidents and wearing coats. They concluded that coats could hinder movements of drivers and be the cause of accidents. A new law was prepared to prohibit drivers from wearing coats when driving.



Finally another study pointed out that people wear coats when it rains...

Correlation is not causality  
Causality really matters

# Another example: Myopia study



- A study published in *Nature* made the causal conclusion that children who sleep with the light on are more likely to develop myopia later in life.

G. E. Quinn, C. H. Shin, M. G. Maguire, and R. A. Stone, “Myopia and ambient lighting at night,” *Nature*, vol. 399, no. 6732, pp. 113–113, 1999

- However, as it turns out, myopic parents tend to leave the light on more often, as well as pass their genetic predisposition to myopia to their children. Accounting for the confounding variable of parent’s myopia, the causal results were subsequently invalidated or substantially weakened.

**Gwiazda J**, Ong E, Held R, *et al.* Myopia and ambient night-time lighting. *Nature* 2000;**404**:144.

**Zadnik K**, Jones LA, Irvin BC, *et al.* Myopia and ambient night-time lighting. *Nature* 2000;**404**:143–4.

3) What if the prediction accuracy is really high?

# Interventions change the environment

- Train/test from same distribution in supervised learning
- No such guarantee in real life!
- Problematic: Acting on a prediction changes distribution!
  - Incl. critical domains: healthcare or adversarial scenarios.
- Connections to covariate shift, domain adaptation [Mansour et al. 2009, Ben-David 2007].





# Recap: Prediction is insufficient for choosing interventions

How often do they lead us to the right decision?

- Unclear, predictive algorithms provide no insight on effects of decisions

Will the predictions be robust tomorrow, or in new contexts?

- Correlations can change
- Causal mechanisms more robust

What if the prediction accuracy is really high? Does that help?

- Active interventions change correlations

PART I. Introduction to Counterfactual Reasoning

PART II. Methods for Causal Inference

PART III. Large-scale and Network Data

PART IV. Broader Landscape

# PART I. Introduction to Counterfactual Reasoning

What is causality?

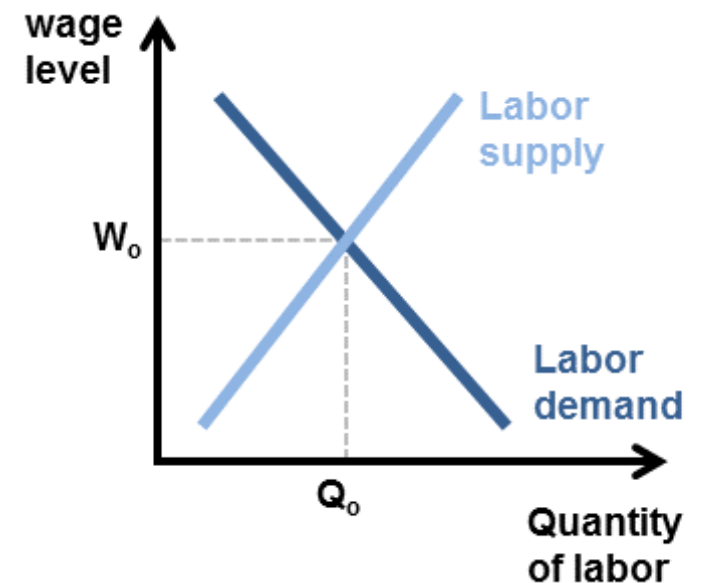
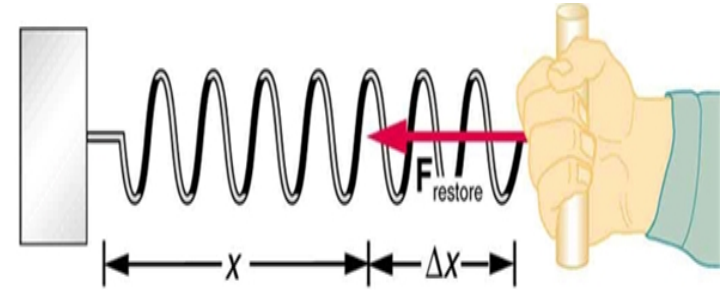
Potential Outcomes Framework

Unobserved Confounds /  
Simpson's Paradox

Structural Causal Model  
Framework

# Cause and Effect

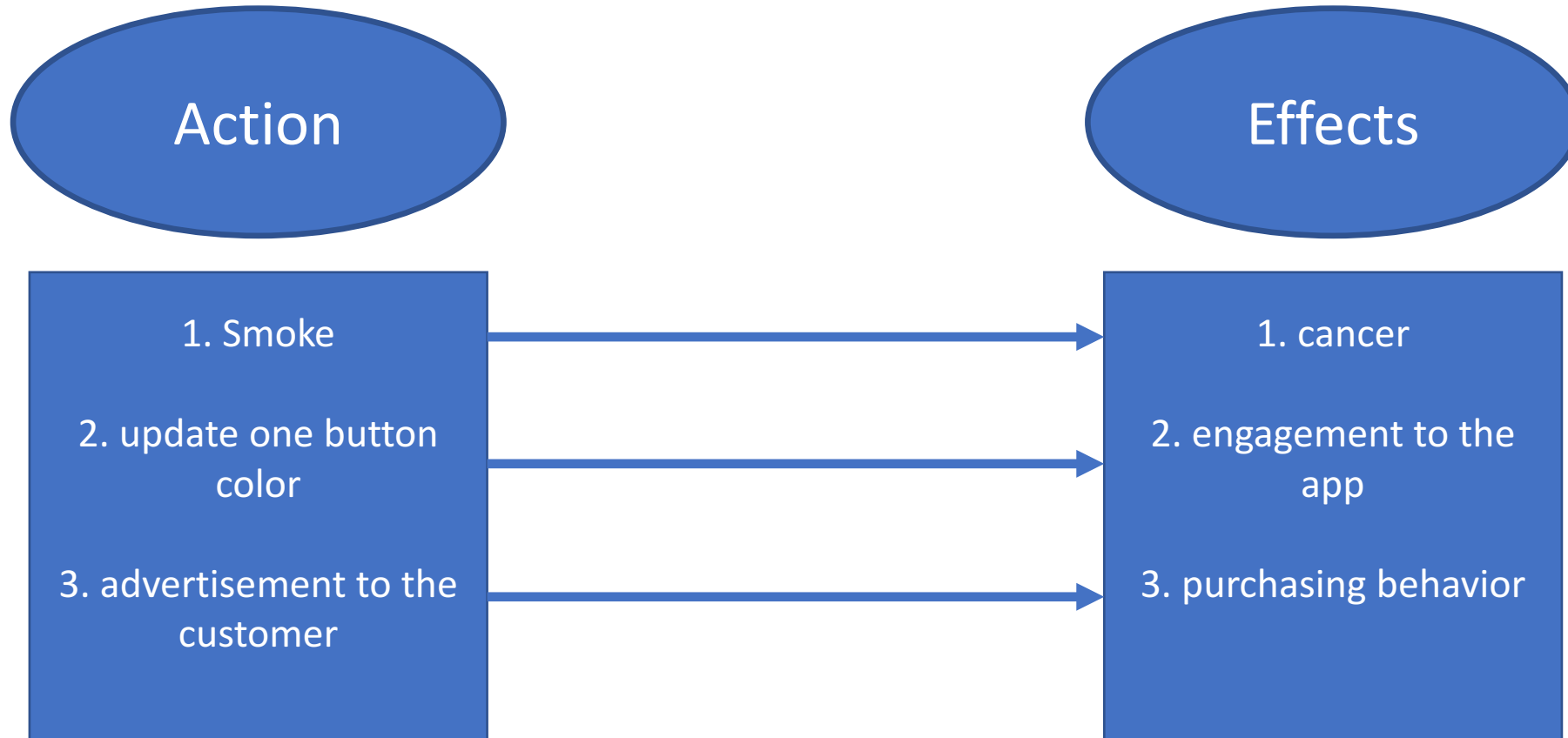
- Questions of cause and effect common in biomedical and social sciences
- Such questions form the basis of almost all scientific inquiry
  - Medicine: drug trials, effect of a drug
  - Social sciences: effect of a certain policy
  - Genetics: effect of genes on disease
- **So what is causality?**
- **What does it mean to *cause* something?**



# Causality examples (A causes B)

- Exposure/Action/Decision

Effects





# A big scholarly debate, from Aristotle to Russell





# What is causality?

- A fundamental question
- Surprisingly, until very recently---maybe the last 30+ years---we have not had a mathematical language of causation. We have not had an arithmetic for representing causal relationships.

*“More has been learned about causal inference in the last few decades than the sum total of everything that had been learned about it in all prior recorded history”*

--Gary King, Harvard University

# The Three Layer Causal Hierarchy

Pearl, Theoretical Impediments to Machine Learning with Seven Sparks from the Causal Revolution, arXiv:1801.04016v1. 11 Jan 2018

Level	Typical Activity	Typical Question	Examples
1. Association $P(y   x)$	Seeing	What is? How would seeing $X$ change my belief in $Y$ ?	What does a symptom tell me about a disease? What does a survey tell us about the election results?
2. Intervention $P(y   do(x), z)$	Doing, Intervening	What if? What if I do $X$ ?	What if I take aspirin, will my headache be cured? What if we ban cigarettes?
3. Counterfactuals $P(y_x   x', y')$	Imagining, Retrospection	Why? Was it $X$ that caused $Y$ ? What if I had acted differently?	Was it the aspirin that stopped my headache? Would Kennedy be alive had Oswald not shot him? What if I had not been smoking the past 2 years?

# A practical definition

**Definition:** T causes Y iff  
changing T leads to a change in Y,  
*keeping everything else constant.*

The **causal effect** is the magnitude by which Y is changed by a unit change in T.

Called the “interventionist” interpretation of causality.

\**Interventionist* definition [<http://plato.stanford.edu/entries/causation-mani/>]



# Keeping everything else constant: Imagine a *counterfactual* world

“What-if” questions

Reason about a world that does not exist.



- What if a system intervention was not done?
- What if an algorithm was changed?
- What if I gave a drug to a patient?

# PART I. Introduction to Counterfactual Reasoning

What is causality?

Potential Outcomes Framework

Unobserved Confounds /  
Simpson's Paradox

Structural Causal Model  
Framework

# Potential Outcomes framework

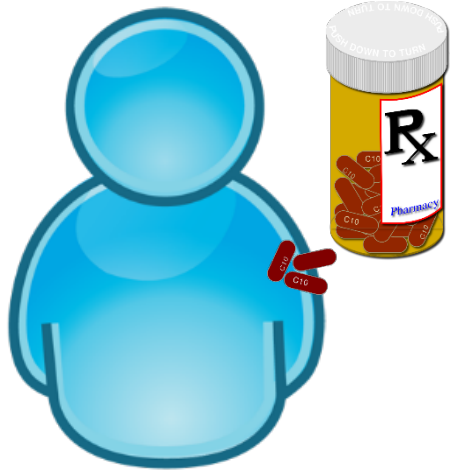


Alice



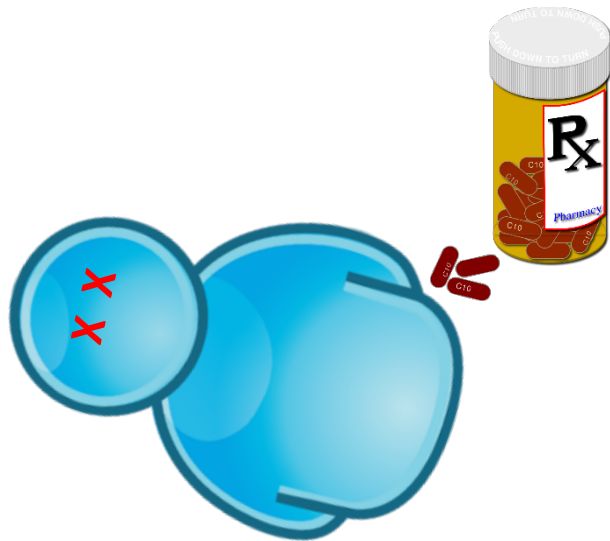
**Treatment**

# Potential Outcomes framework



Alice

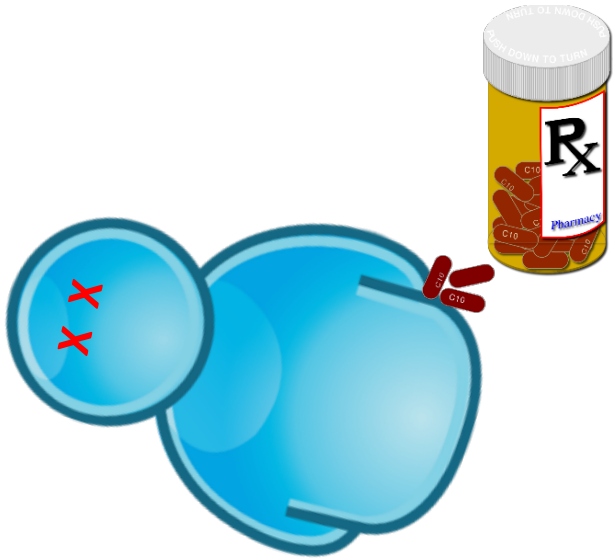
# Potential Outcomes framework



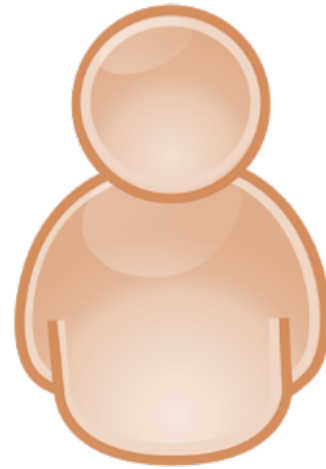
Alice



# Potential Outcomes framework: Introduce a counterfactual quantity



$Y_{T=1}$



$Y_{T=0}$



Causal effect of treatment =

$$E[Y_{T=1} - Y_{T=0}]$$

Causal inference is the problem of estimating the counterfactual  $Y_{t=\sim t}$

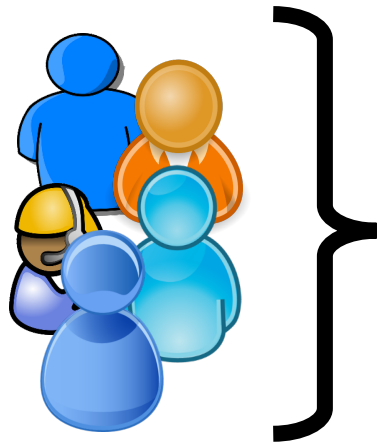
Person	T	$Y_{T=1}$	$Y_{T=0}$
P1	1	0.4	0.3
P2	0	0.8	0.6
P3	1	0.3	0.2
P4	0	0.3	0.1
P5	1	0.5	0.5
P6	0	0.6	0.5
P7	0	0.3	0.1

Causal effect:  $E[Y_{t=1} - Y_{t=0}]$

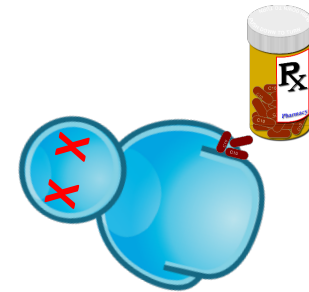
**Fundamental problem of causal inference:** For any person, observe only one: either  $Y_{t=1}$  or  $Y_{t=0}$

# Fundamental problem: counterfactual outcome is not observed

- “Missing data” problem
- Estimate missing data values using various methods
- $Y_{T=0}$  now becomes an estimated quantity, based on outcomes of other people who did not receive treatment



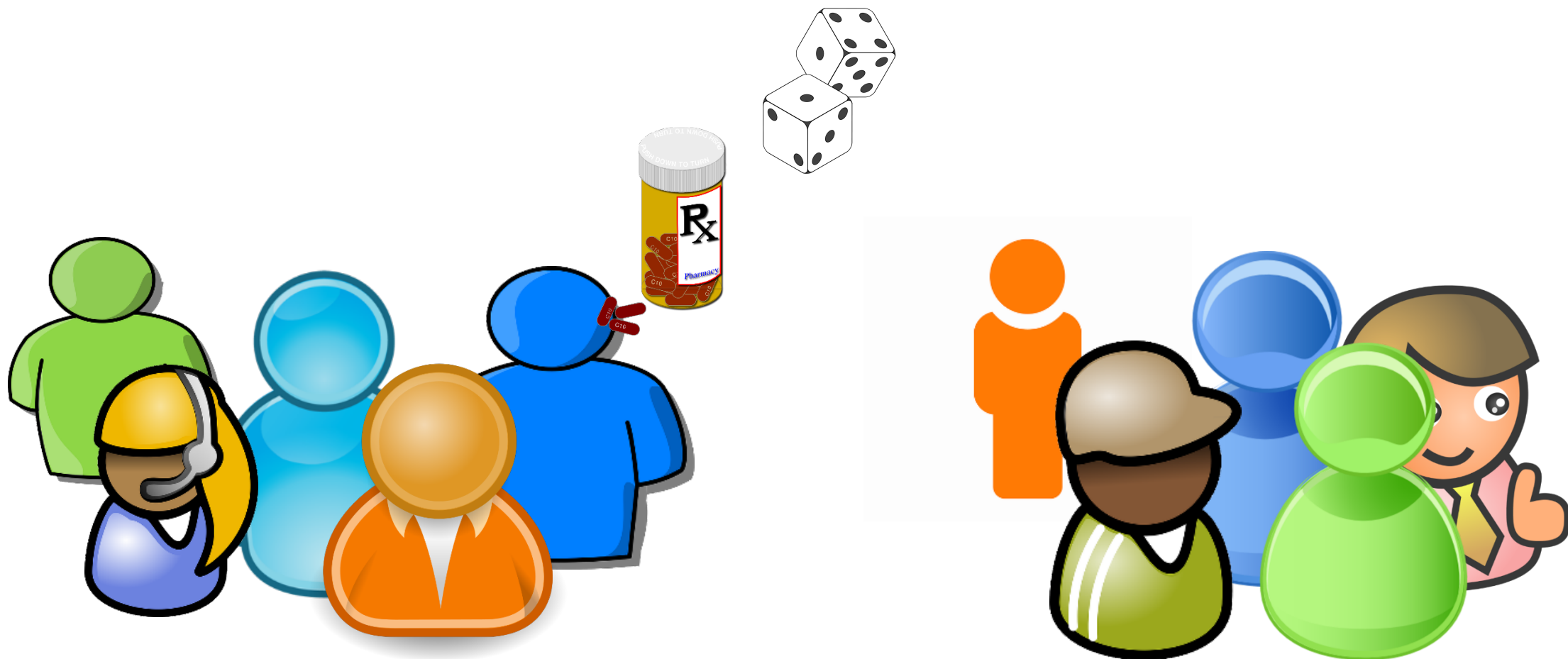
$$\hat{Y}^{T=0}$$



$$Y^{T=1}$$

# Randomized Experiments are the “gold standard”

One way to estimate counterfactual



# Cost: Possibly risky, unethical

Unethical to deny useful treatment or administer risky treatment.

Infeasible or costly in other situations.

What can we do when an experiment is not possible?  
Coming soon in Section 2

# Recap: Potential Outcomes Framework

- Potential outcomes reasons about causal effects by comparing outcome of treatment to outcome of no-treatment
- For any individual, we cannot observe both treatment and no-treatment.
- Randomized experiments are one solution
- We'll discuss others in tutorial Section 2



# PART I. Introduction to Counterfactual Reasoning

What is causality?

Potential Outcomes Framework

Unobserved Confounds /  
Simpson's Paradox

Structural Causal Model  
Framework

# Example: Auditing the effect of an algorithm

System changes algorithm from A to B at some point.

Is the new algorithm B better?

Say a feature that provides information or discount for a financial product.



Algorithm A

Success  
Rate= $\rho$



Algorithm B

?

# New algorithm increases overall success rate

Two algorithms, A (old) and B (new) running on the system.

From system logs, collect data for 1000 sessions for each.

Measure Success Rate (SR).

Old Algorithm (A)	New Algorithm (B)
50/1000 ( <b>5%</b> )	54/1000 ( <b>5.4%</b> )

New algorithm is better?

# Unobserved Confounds

What if there are unobserved features of audience that matter?



Old Algorithm (A)	New Algorithm (B)	Low-income Users
10/400 (2.5%)	4/200 (2%)	

Old Algorithm (A)	New Algorithm (B)	High-income Users
40/600 (6.6%)	50/800 (6.2%)	

The Simpson's paradox: New algorithm is better overall, but worse for each subgroup

	Old algorithm (A)	New Algorithm (B)
CTR for Low-income users	10/400 (2.5%)	4/200 (2%)
CTR for High-income users	40/600 (6.6%)	50/800 (6.2%)
<b>Total CTR</b>	<b>50/1000 (5%)</b>	<b>54/1000 (5.4%)</b>

So, which is better?

# From metrics to decision-making

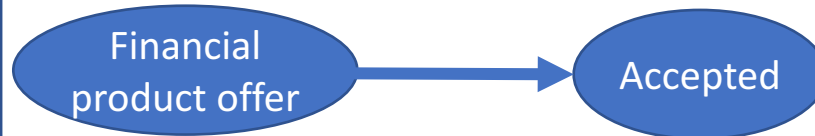
Did the change to new Algorithm increase success rate for the system?

Answer (as usual):

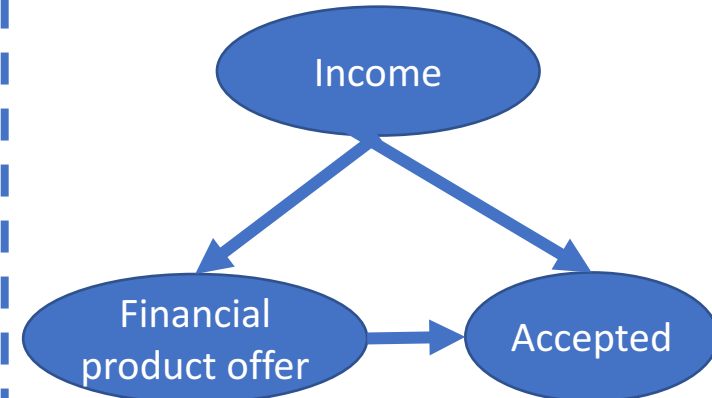
Maybe, maybe not (!)



Higher success rate due to  
**new algorithm**



Higher success rate due to  
**selection effects**



E.g., Algorithm B is shown at a different time than A.

There could be other hidden causal variations.

Not just theory. Differences in interpretations can attract lawsuits (UC Berkeley admissions, 1973)

# Simpson's Paradox in naturally generated data

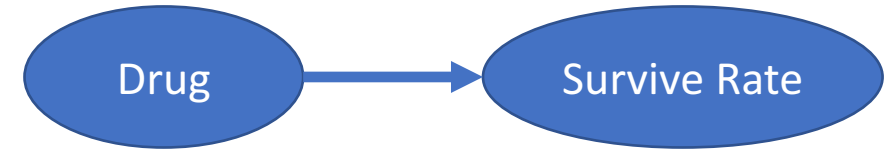


Table 1: Yule-Simpson's Paradox

Population			
	Survive	Die	Survive Rate
Treatment	20	20	50%
Control	16	24	40%
Male			
	Survive	Die	Survive Rate
Treatment	18	12	60%
Control	7	3	70%
Female			
	Survive	Die	Survive Rate
Treatment	2	8	20%
Control	9	21	30%

Treatment is better

Control is better

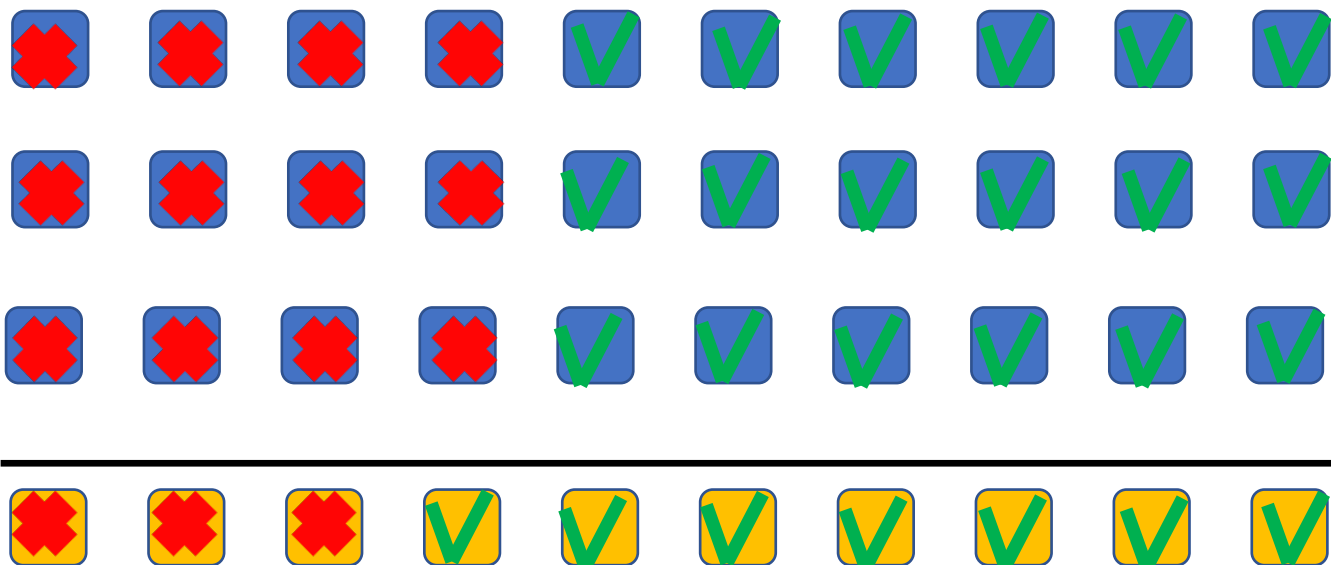
Control is better



# Simpson's Paradox

Table 1: Yule-Simpson's Paradox

Population	Survive	Die	Survive Rate
Treatment	20	20	50%
Control	16	24	40%
Male			
Treatment	18	12	60%
Control	7	3	70%
Female			
Treatment	2	8	20%
Control	9	21	30%



Male treatment

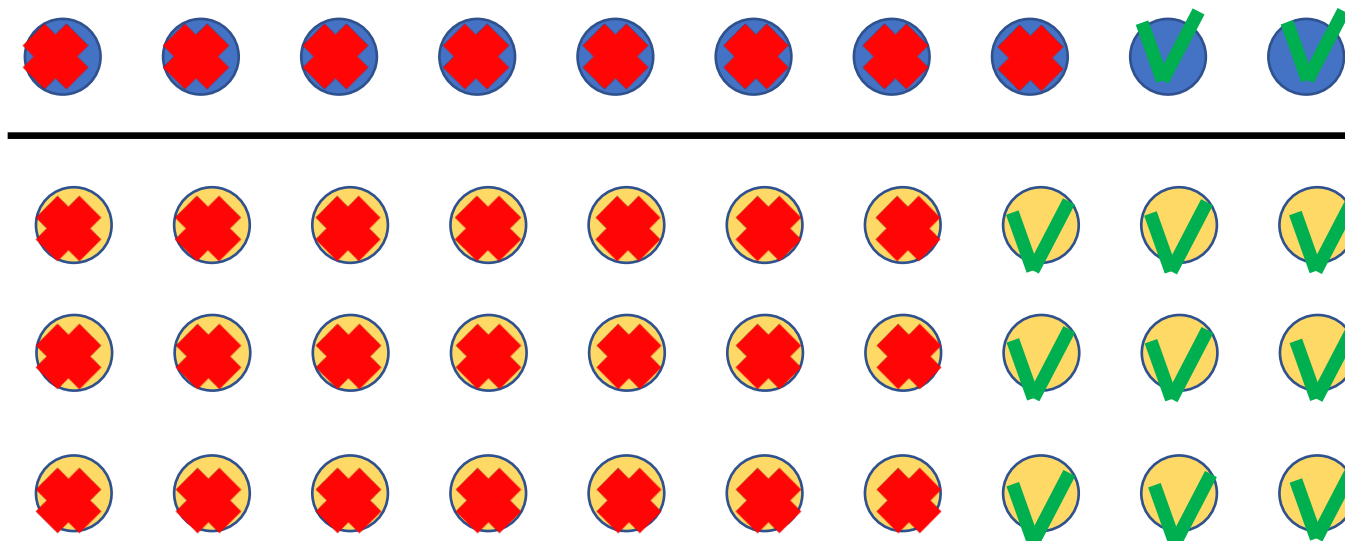


Male control

# Simpson's Paradox

Table 1: Yule-Simpson's Paradox

Population			
	Survive	Die	Survive Rate
Treatment	20	20	50%
Control	16	24	40%
Male			
	Survive	Die	Survive Rate
Treatment	18	12	60%
Control	7	3	70%
Female			
	Survive	Die	Survive Rate
Treatment	2	8	20%
Control	9	21	30%



Female treatment

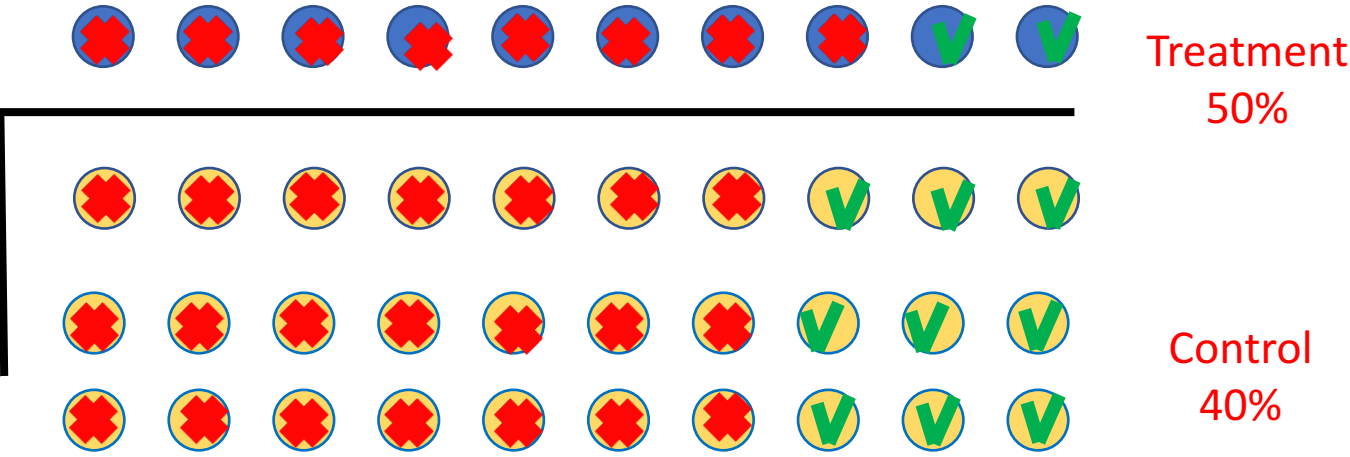
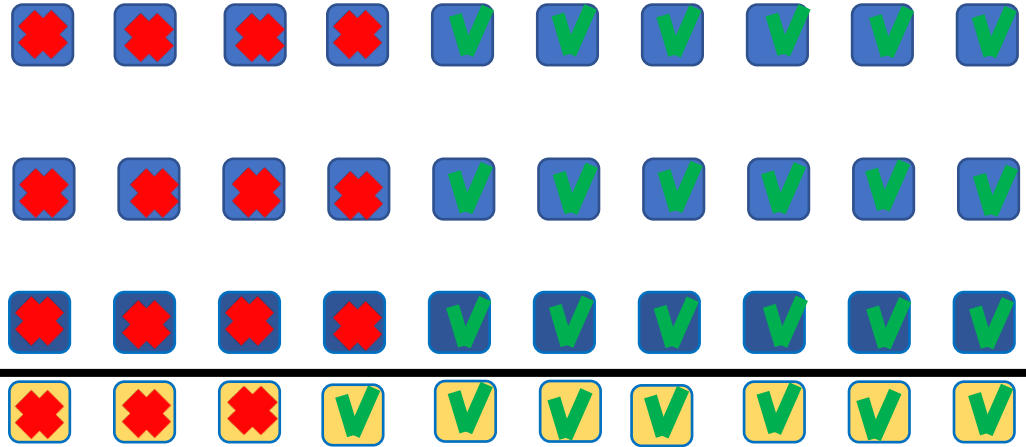




Female control



# Simpson's Paradox

Table 1: Yule-Simpson's Paradox

Population			
	Survive	Die	Survive Rate
Treatment	20	20	50%
Control	16	24	40%
Male			
	Survive	Die	Survive Rate
Treatment	18	12	60%
Control	7	3	70%
Female			
	Survive	Die	Survive Rate
Treatment	2	8	20%
Control	9	21	30%



 Male treatment  
 Male control

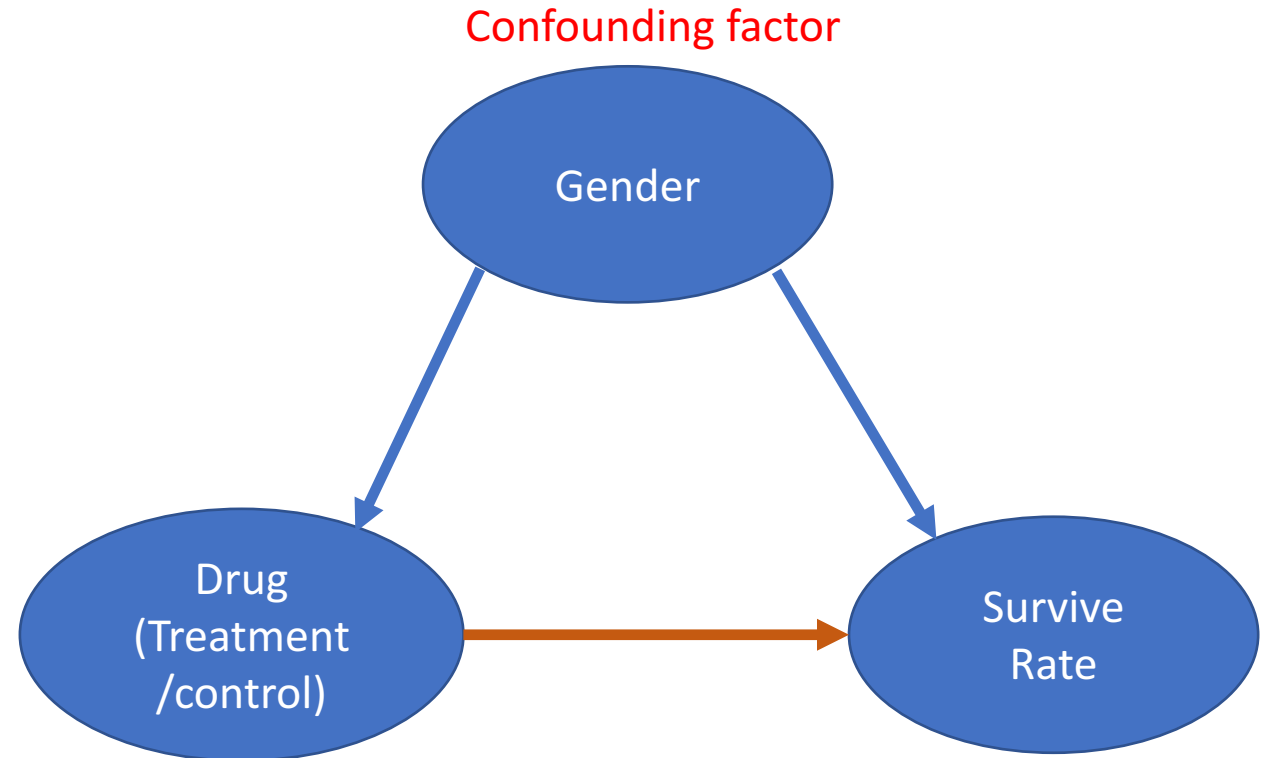
 Female treatment  
 Female control

Treatment  
50%  
  
 Control  
40%

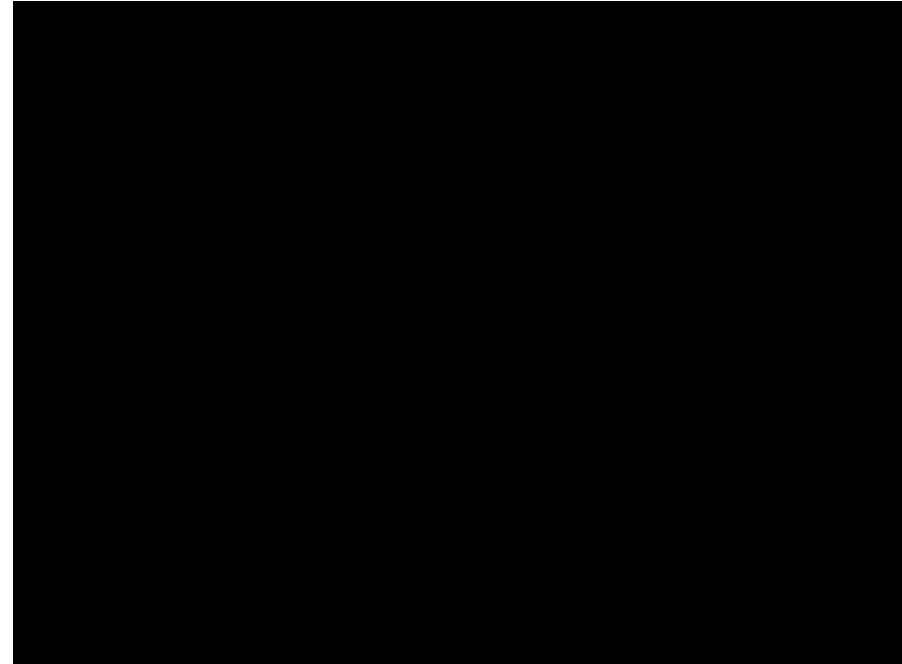
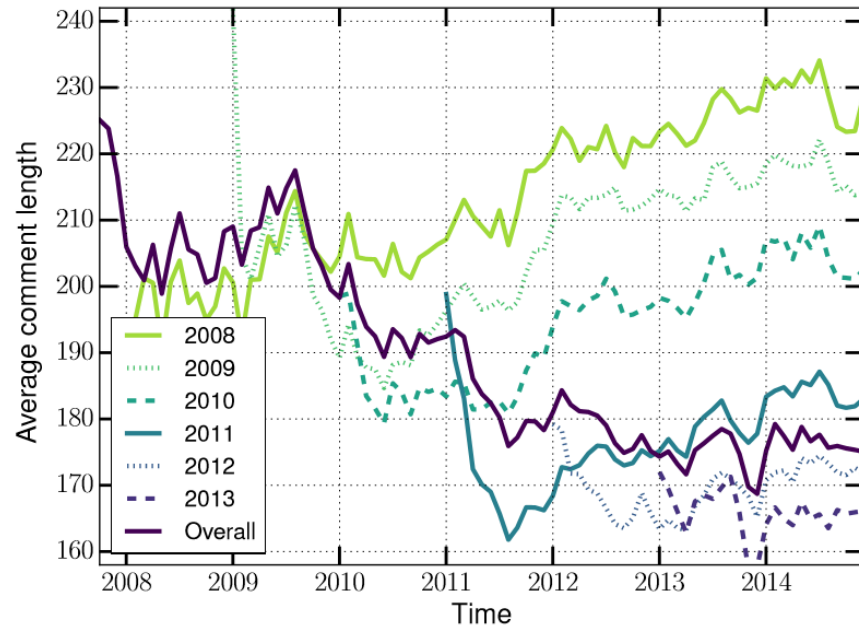
# Confounding factor: Gender

Table 1: Yule-Simpson's Paradox

Population	Survive	Die	Survive Rate
Treatment	20	20	50%
Control	16	24	40%
Male	Survive	Die	Survive Rate
Treatment	18	12	60%
Control	7	3	70%
Female	Survive	Die	Survive Rate
Treatment	2	8	20%
Control	9	21	30%



# Example: Simpson's paradox in Reddit



Average comment length decreases over time.

Making sense of such data can be too complex.

**D'oh!**



Not Simpson's Paradox

# Recap: Unobserved Confounds

- Unobserved confounds are a threat to causal reasoning

# PART I. Introduction to Counterfactual Reasoning

What is causality?

Potential Outcomes Framework

Unobserved Confounds /  
Simpson's Paradox

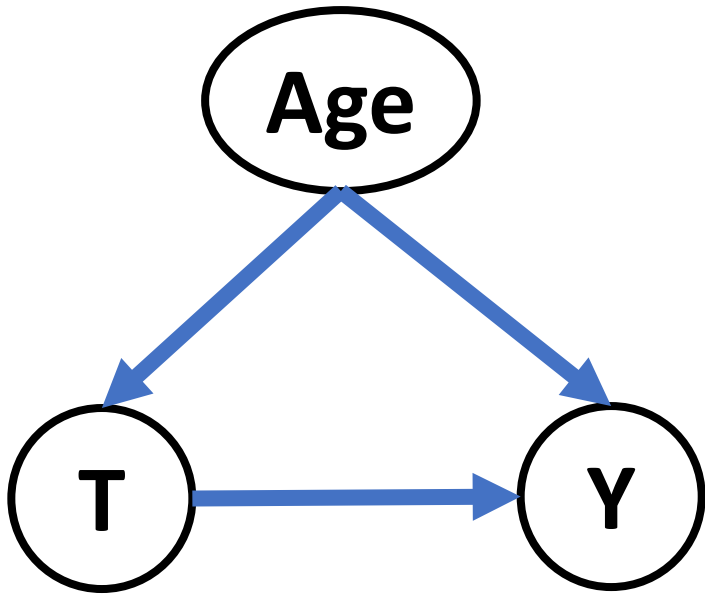
Structural Causal Model  
Framework



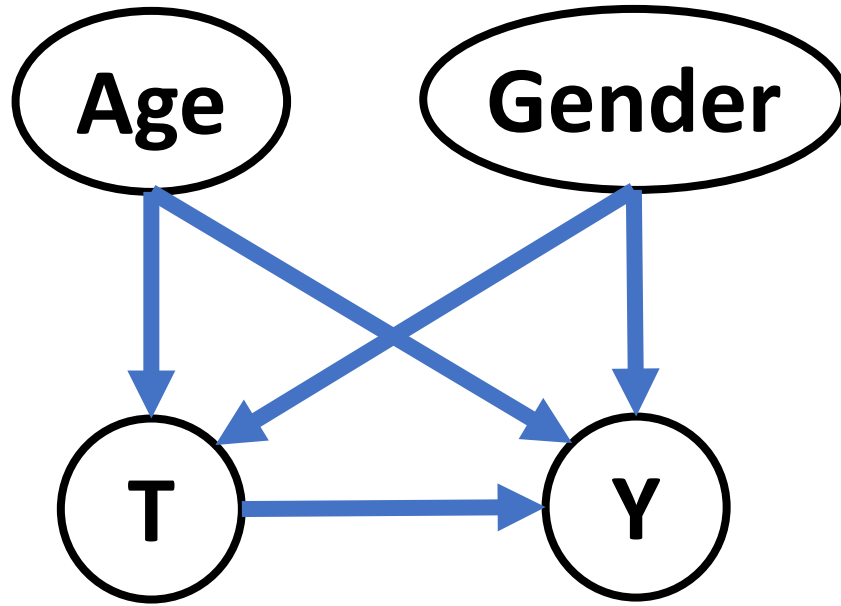
# Real world is complicated

- People may have inter-related characteristics
  - How are these characteristics associated with each other?
- Other factors can influence the observed outcome
  - How do they affect treatment and outcome?
  - Which ones to include?
- How to identify the causal effect in such cases?
- When is it possible to find a causal effect?
  - We can use graphical model framework to answer this

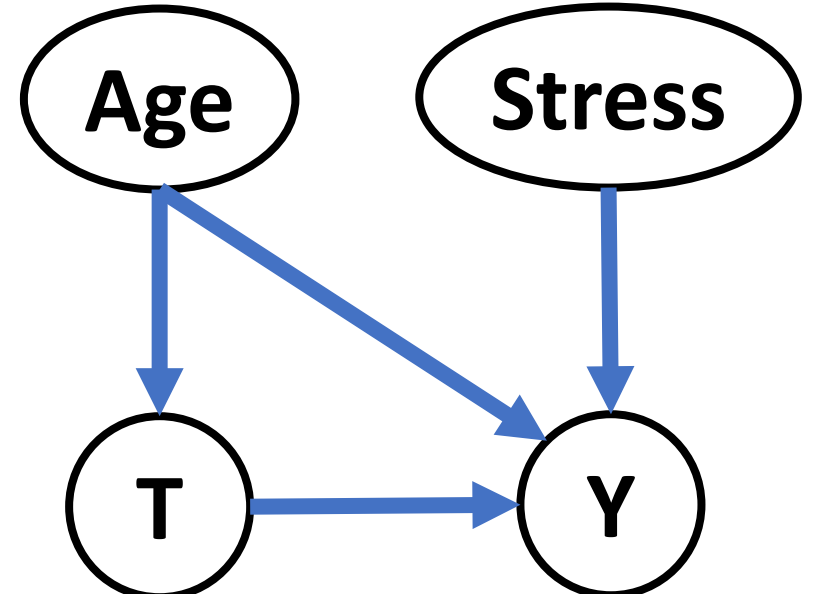
Which variables to condition on?



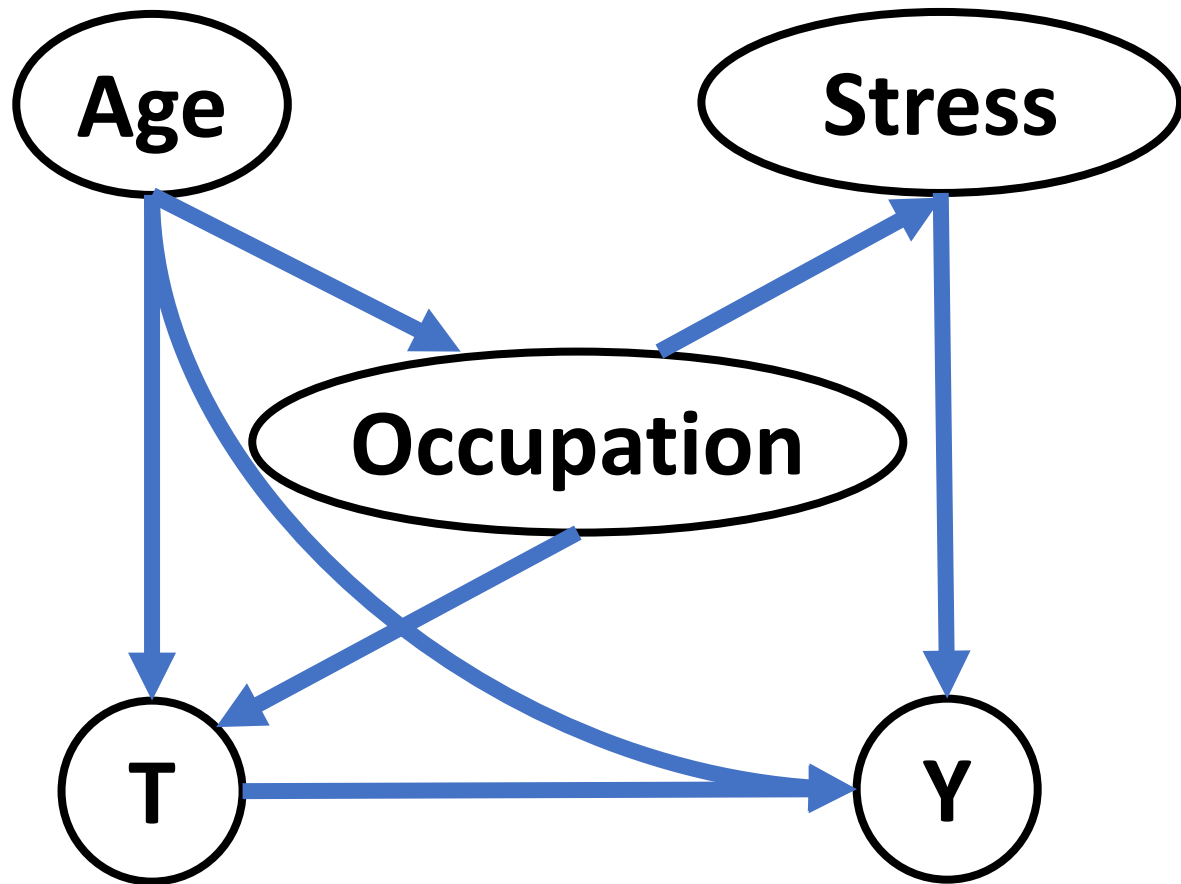
$X = \{Age\}$



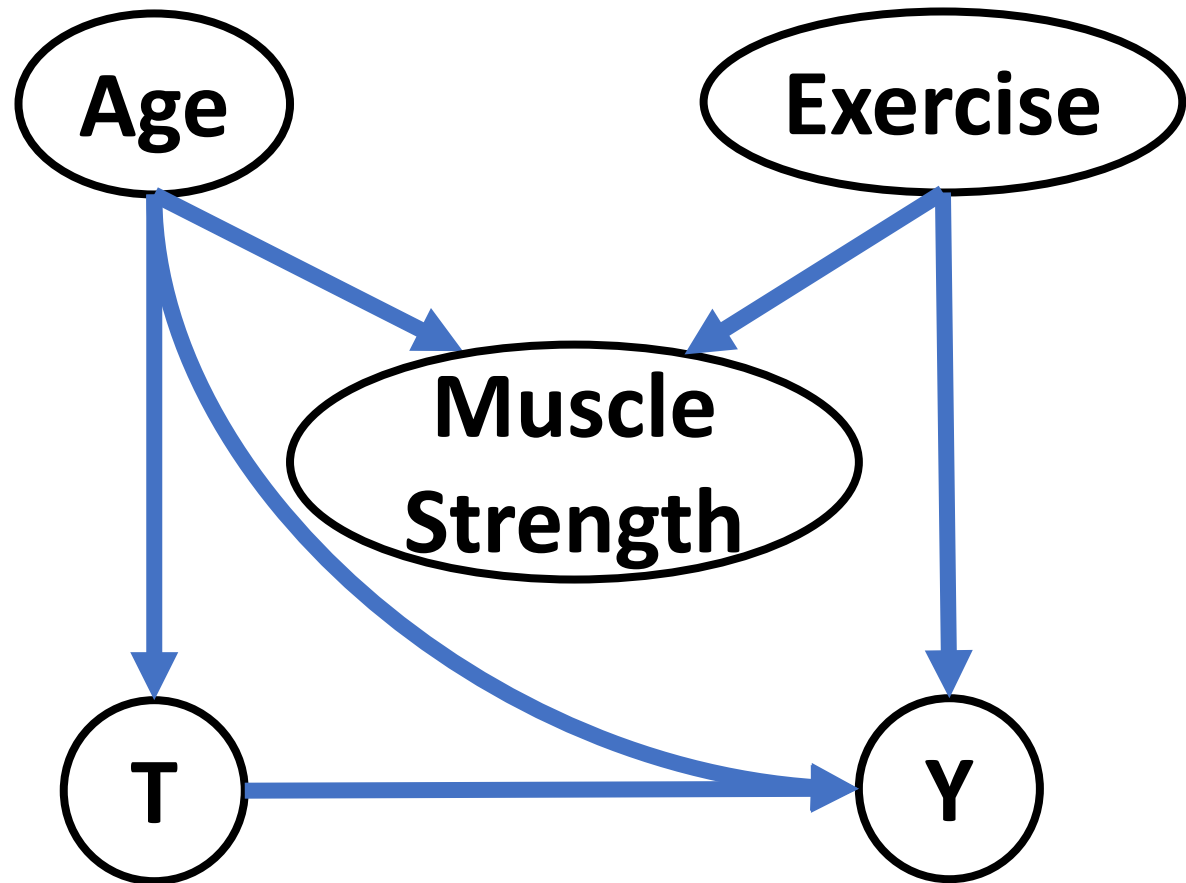
$X = \{Age, Gender\}$



$X = \{Age\}$

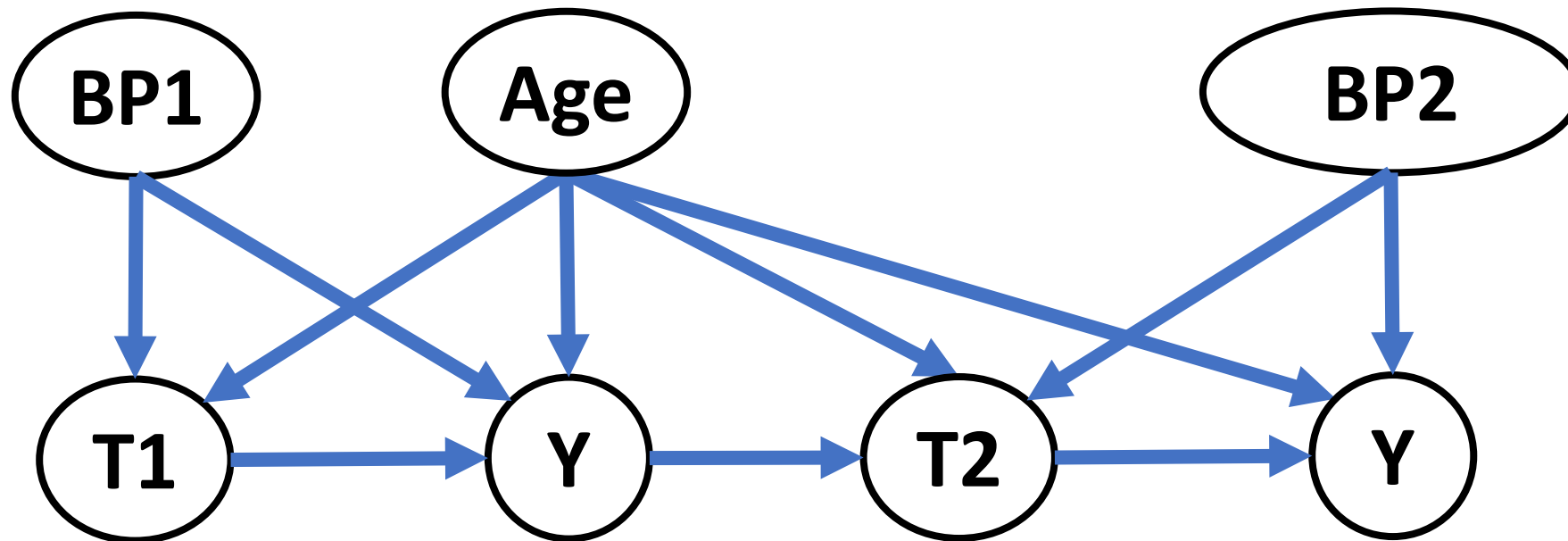


$X = ?$



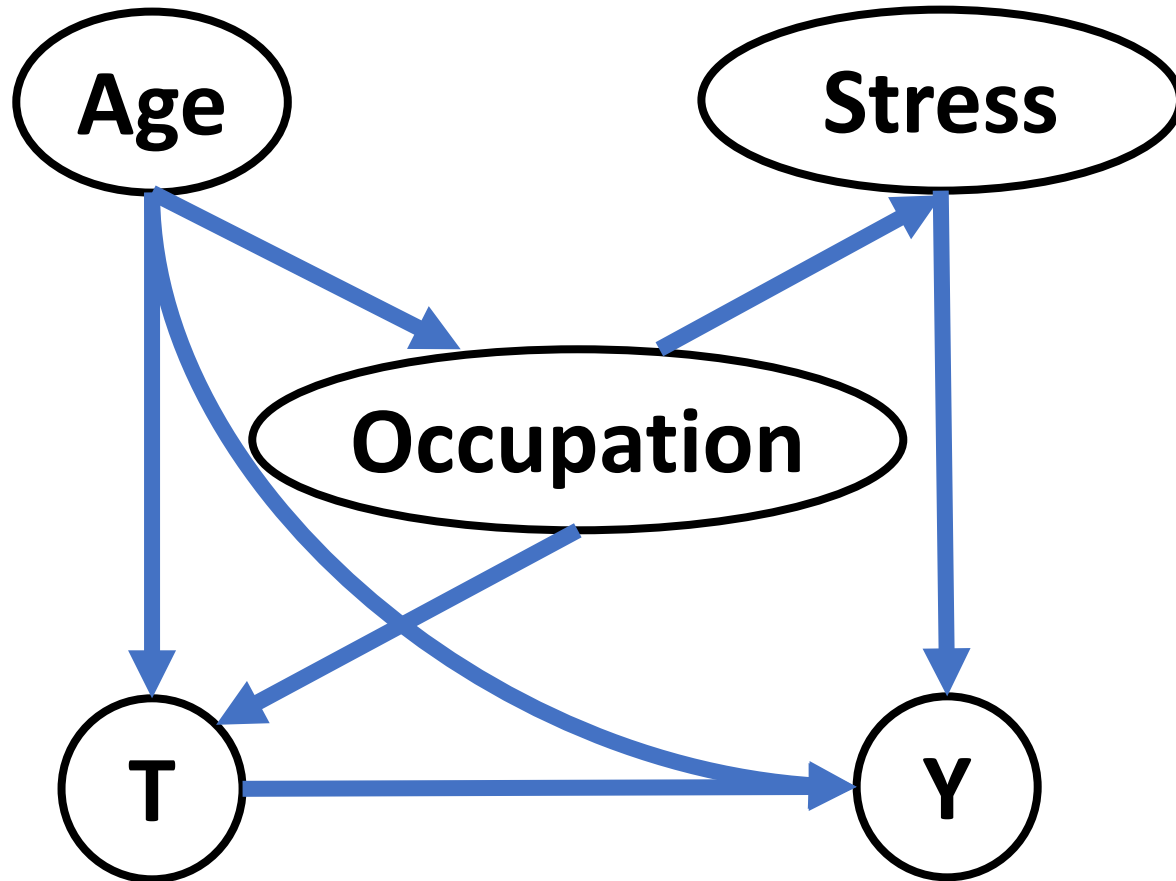
$X = ?$

Another example: Repeated treatment (!)



**How to reason about causal effects in such cases?**

# Structural Causal Model: A framework for expressing complex causal relationships



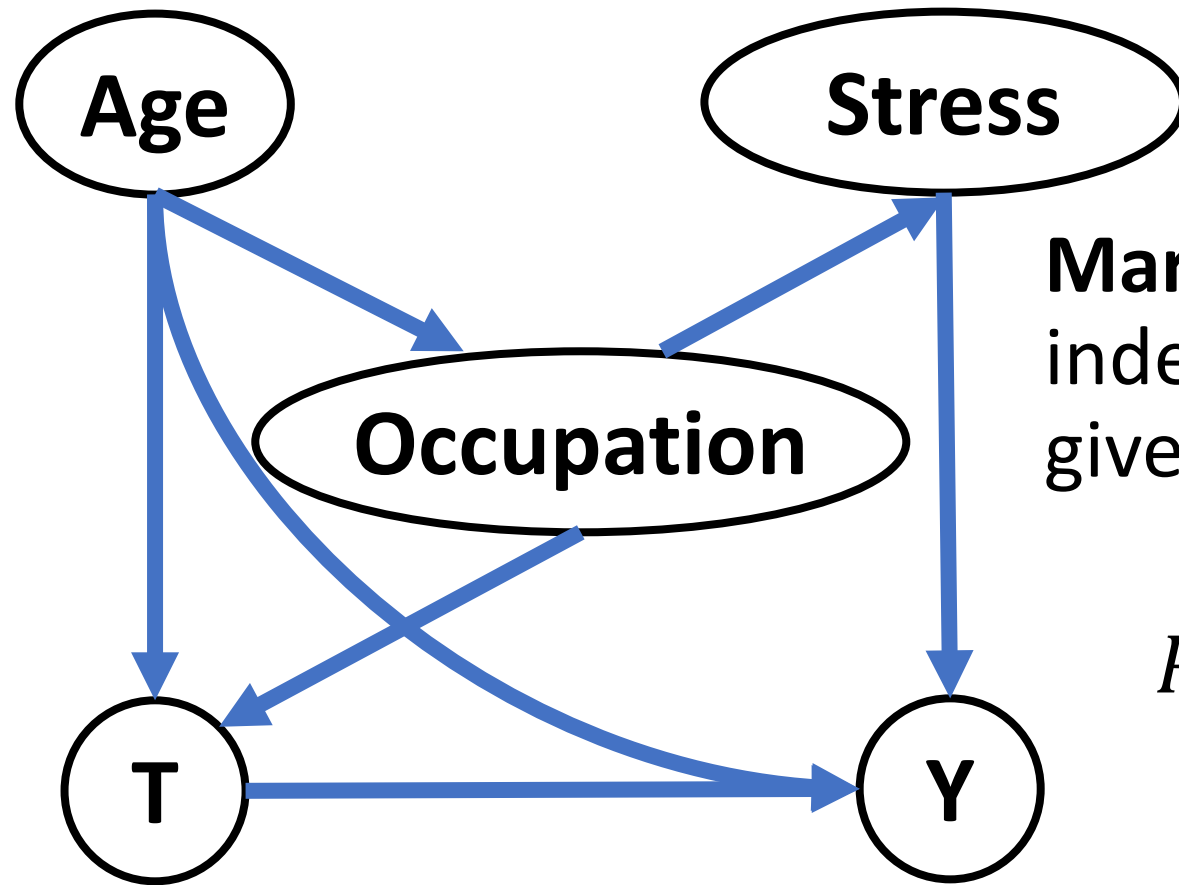
$$\begin{aligned} \text{Occupation} &= h(\text{Age}, u_o) \\ \text{Stress} &= k(\text{Occupation}, u_s) \end{aligned}$$

$$\begin{aligned} T &= g(\text{Age}, \text{Occupation}, u_t) \\ Y &= f(T, \text{Age}, \text{Stress}, u_y) \end{aligned}$$

Edges represent *direct* causes.

Directed paths represent *indirect* causes.

# Structural Causal Model: A framework for expressing complex causal relationships

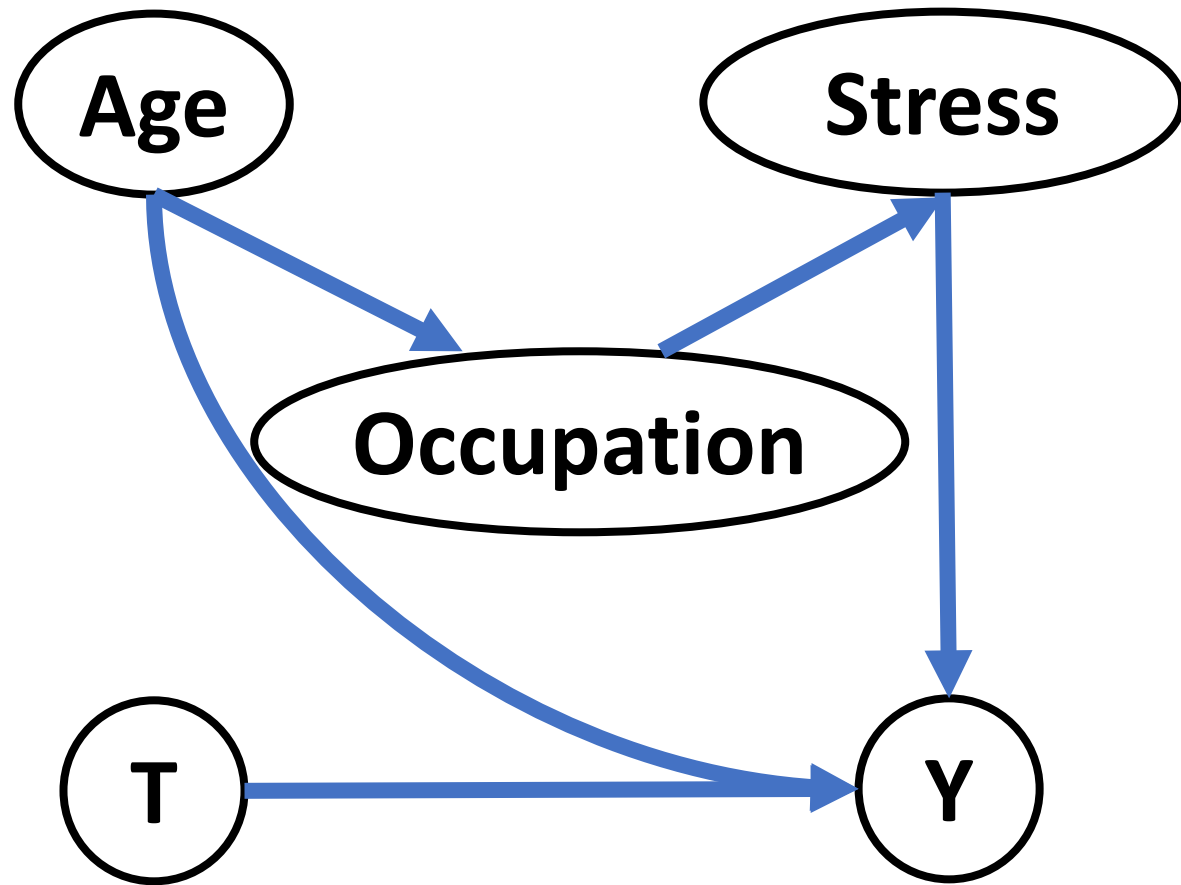


**Markov assumption:** A node is independent of all its non-descendants given its parents.

$$P(T|Occ., Stress) = P(T|Occ.)$$

$$P(G) = P(Age)P(Occ.|Age)P(Stress|Occ.)P(T|Age, Occ.)P(Y|T, Age, Stress)$$

Structural Causal Model: Causal effect is represented by the intervention distribution



**Counterfactual (Intervention) world:**

All edges to Treatment T removed, *keeping everything else the same.*

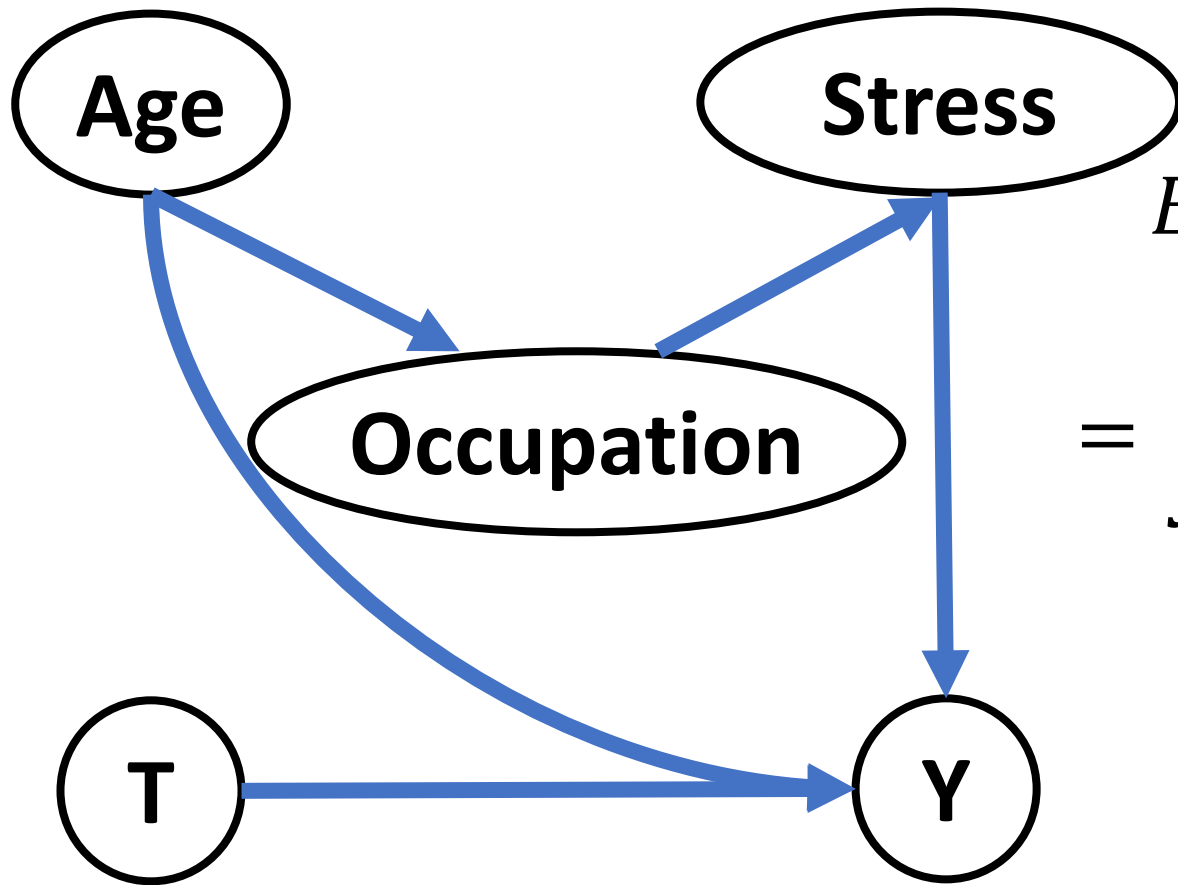
Observed correlation =  $P(Y|T)$

Causal Effect =  $P^*(Y|T)$

$$P^*(\Phi) = P(\text{Age})P(\text{Occ.}|\text{Age})P(\text{Stress}|\text{Occ.})P^*(T|\text{Age}, \text{Occ.})P(Y|T, \text{Age}, \text{Stress})$$



Structural Causal Model: Causal effect is represented by the intervention distribution

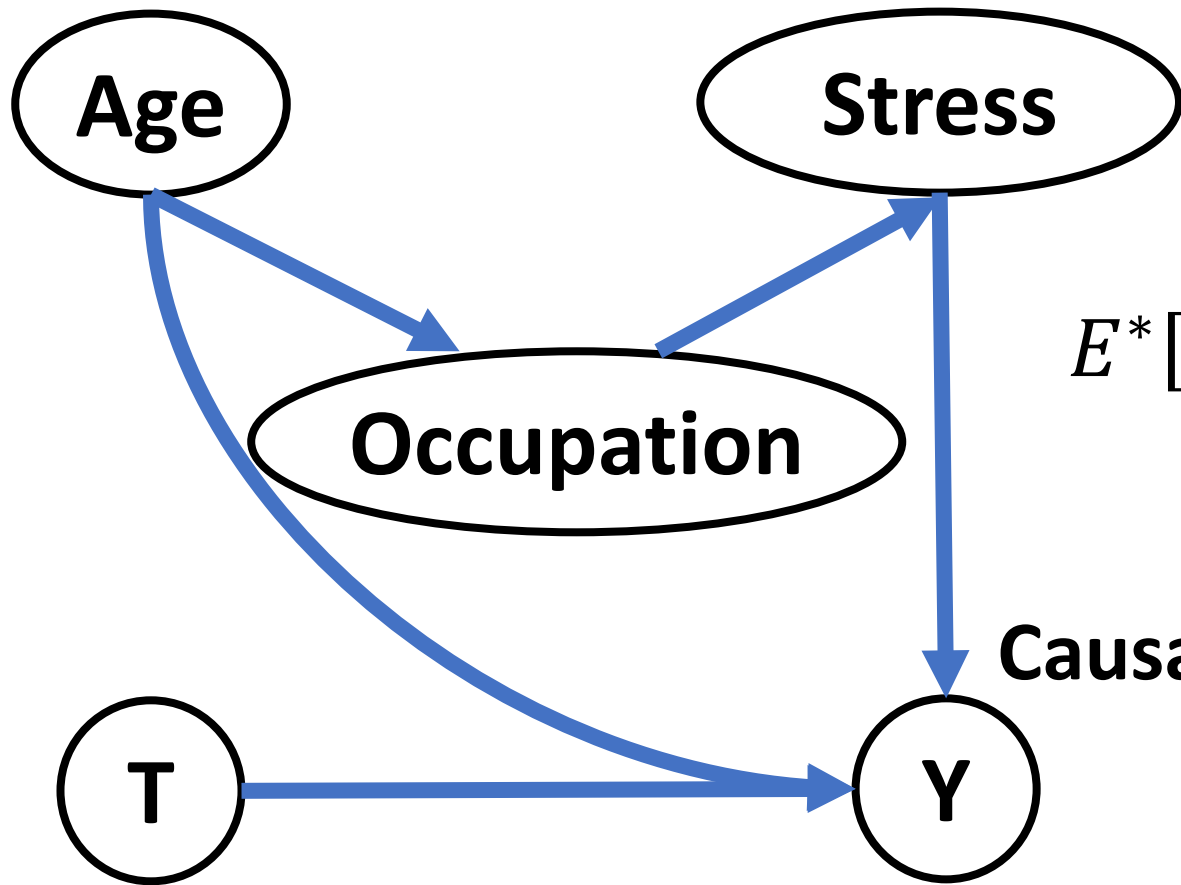


$$\begin{aligned}
 E^*[Y] &= E_{\phi \sim P^*(\Phi)}[y] = \int_{\phi} y P^*(\phi) \\
 &= \int_{\phi} y \frac{P(\phi)}{P(\phi)} P^*(\phi) = \int_{\phi} y \frac{P^*(\phi)}{P(\phi)} P(\phi) \\
 &= \int_{\phi} y \left[ \frac{P^*(T|Age, Occ.)}{P(T|Age, Occ.)} \right] P(\phi)
 \end{aligned}$$

$$P^*(\Phi)$$

$$= P(Age)P(Occ. | Age)P(Stress | Occ.)P^*(T | Age, Occ.)P(Y | T, Age, Stress) \quad 60$$

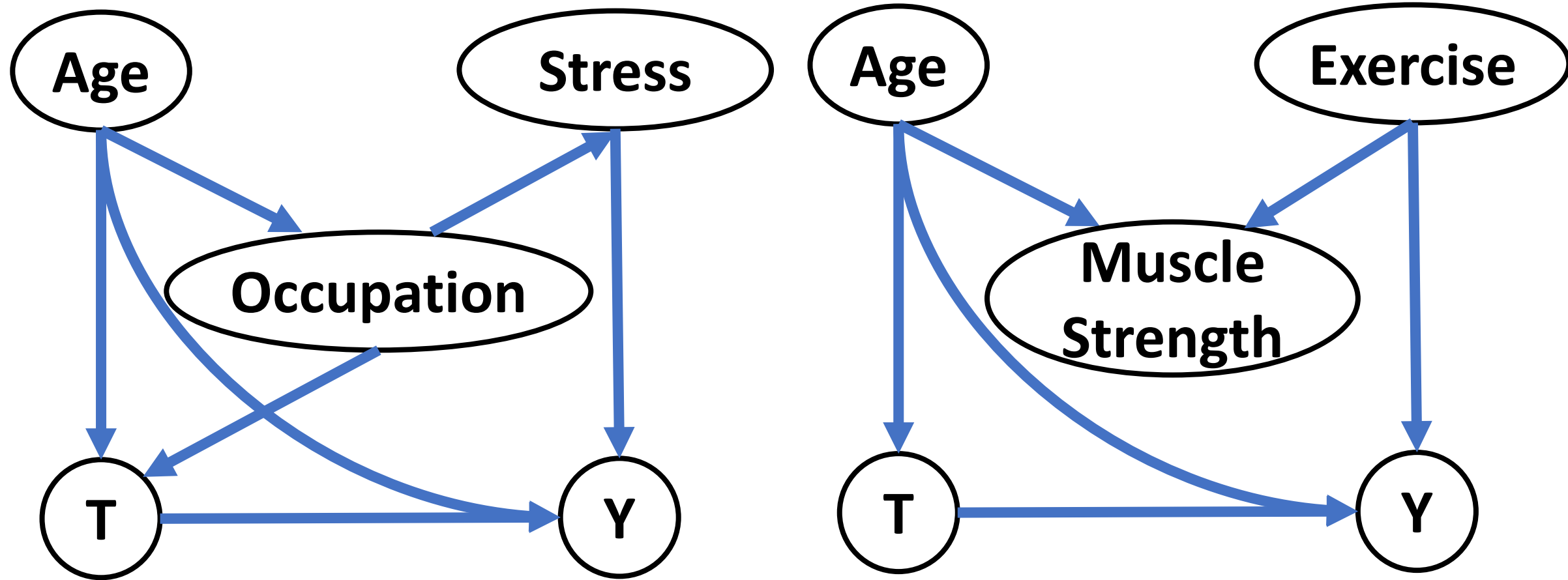
Structural Causal Model: Causal effect is represented by the intervention distribution



$$E^*[Y] = \int_{\phi} y \left[ \frac{P^*(T|Age, Occ.)}{P(T|Age, Occ.)} \right] P(\phi)$$

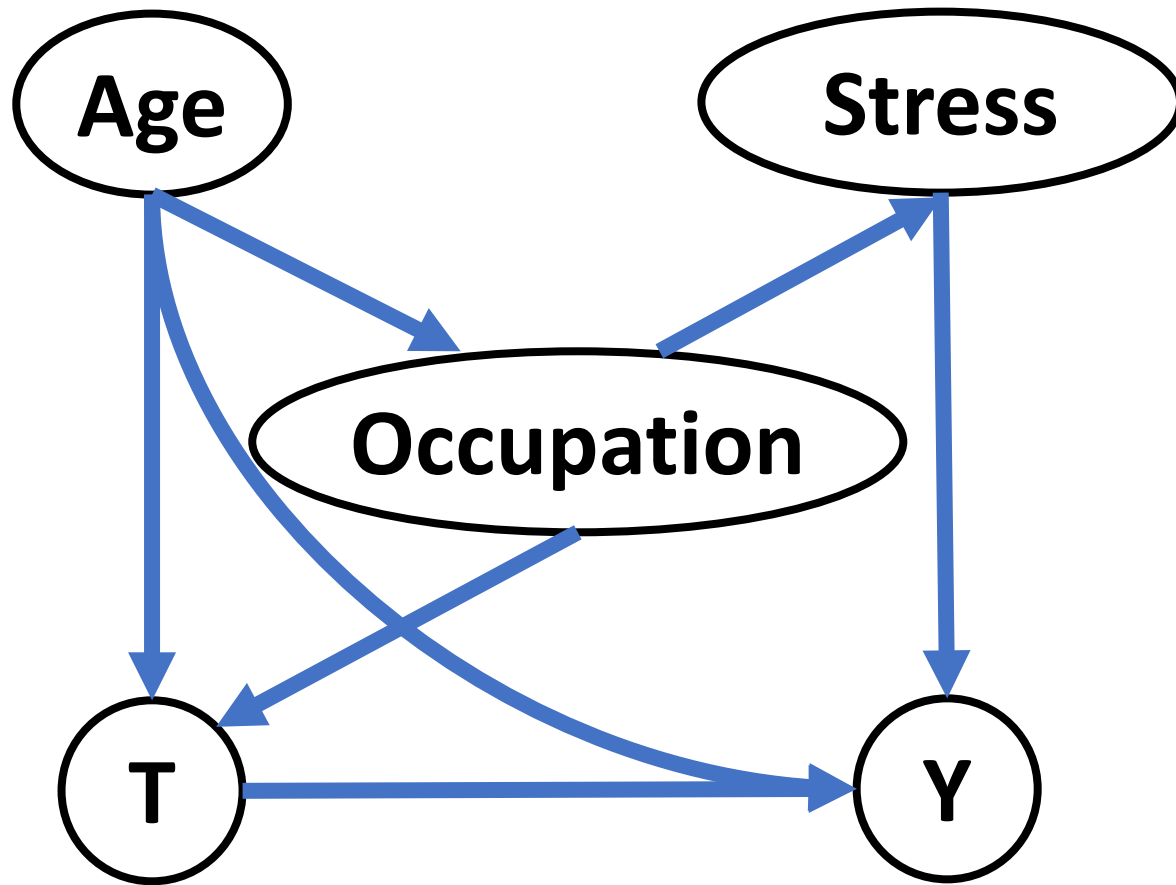
**Causal Effect:**  $E^*[Y|T = 1] - E^*[Y|T = 0]$

Structural Causal Model makes assumptions explicit



The graph encodes all causal assumptions.

Important: Assumptions are the edges that are *missing*



**Assumption 1:** Occupation does affect outcome Y.

**Assumption 2:** Age does affect stress.

**Assumption 3:** Stress does not affect Occupation.

**Assumption 4:** Treatment does not affect stress.

*..and so on.*

**Condition for validity:** The graph reflects all relevant causal processes.

Important: SCM and Potential Outcome frameworks are equivalent

**Potential Outcomes**

$$E[Y_{T=1}] - E[Y_{T=0}]$$

**Structural Causal Model**

$$E^*[Y|T = 1] - E^*[Y|T = 0]$$

If we denote  $E[Y_T] \leftarrow E^*[Y|T]$ , then the formulations are equivalent.

More formally, a theorem in one framework is a theorem in another.

# Key Benefit (1) of SCM: Provides a language for expressing counterfactuals

*If a person was given treatment, what is the probability that he would be cured if he was not given treatment?*

$$P(Y = 1|T = 1, T = 0)$$

**Non-sensical.**

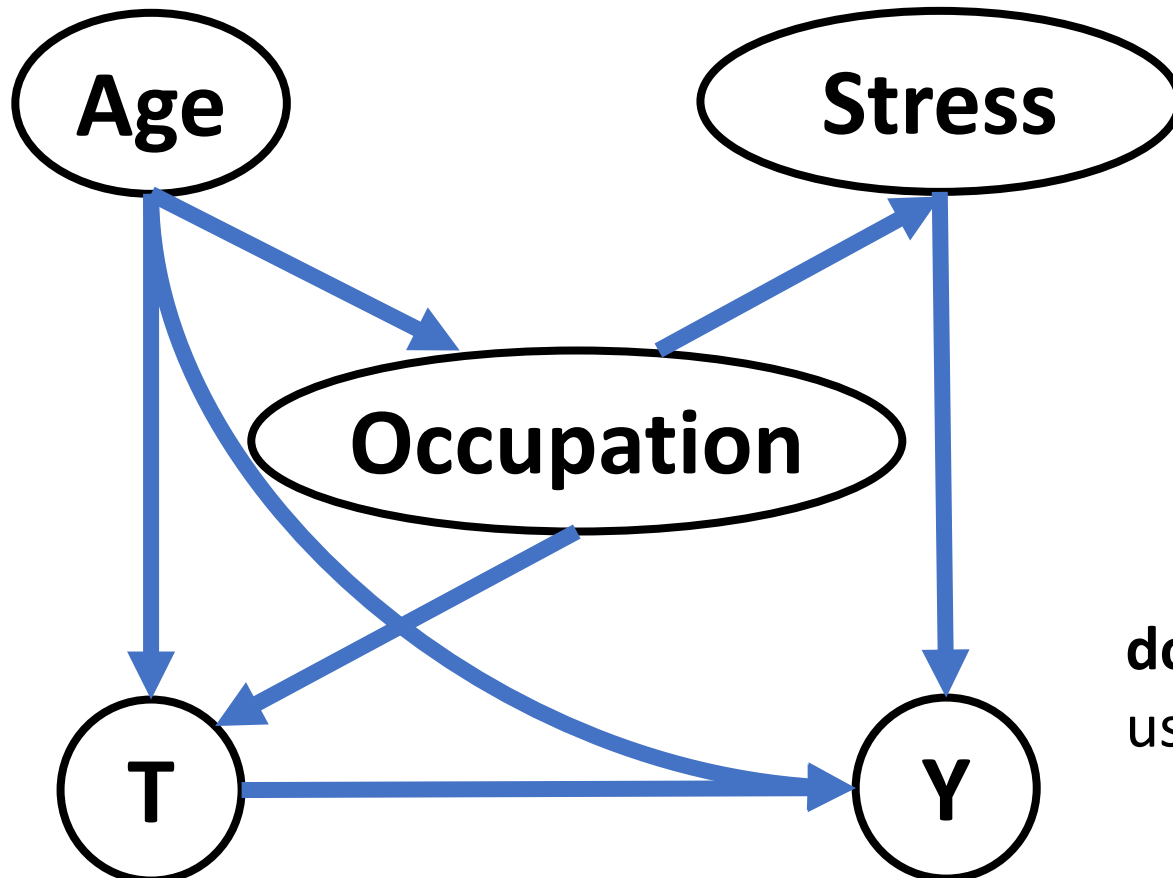
*Can write it as:*

$$P(Y_{T=0} = 1|T = 1), \text{ or } \\ P(Y = 1|T = 1, do(T = 0))$$

$P(Y|do(T))$  avoids confusion with  $P(Y|T)$

# Key Benefit 2 of SCM: Provides a mechanistic way of identifying causal effect

**do-calculus:** A rule-based calculus that can help identify any counterfactual quantity.



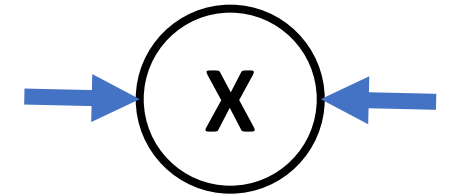
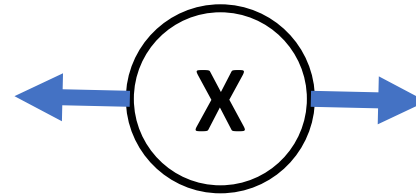
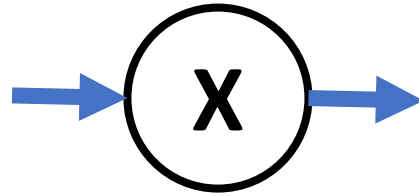
E.g.,  
 $P(Y|do(T))$   
 $= \dots do\text{-calculus rules} \dots$

$$= \sum_{Age, Stress} P(Y|T, Age, Stress) P(Age, Stress)$$

**do-calculus is complete:** If we cannot identify using do-calculus, causal effect is unidentifiable.

# Advanced Topic: Back-door criterion

Three kinds of node-edges



**Path is “blocked”**

If conditioned on X

If conditioned on X

If not conditioned on X

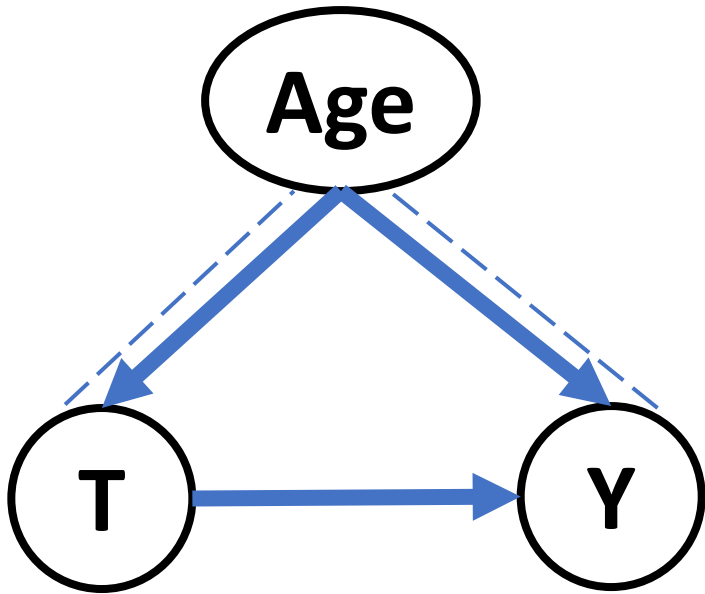
**“Back-door” path:** Any undirected path that starts with  and ends with 

**Back-door criterion:** If conditioning on X blocks all back-door paths between treatment T and outcome Y, then

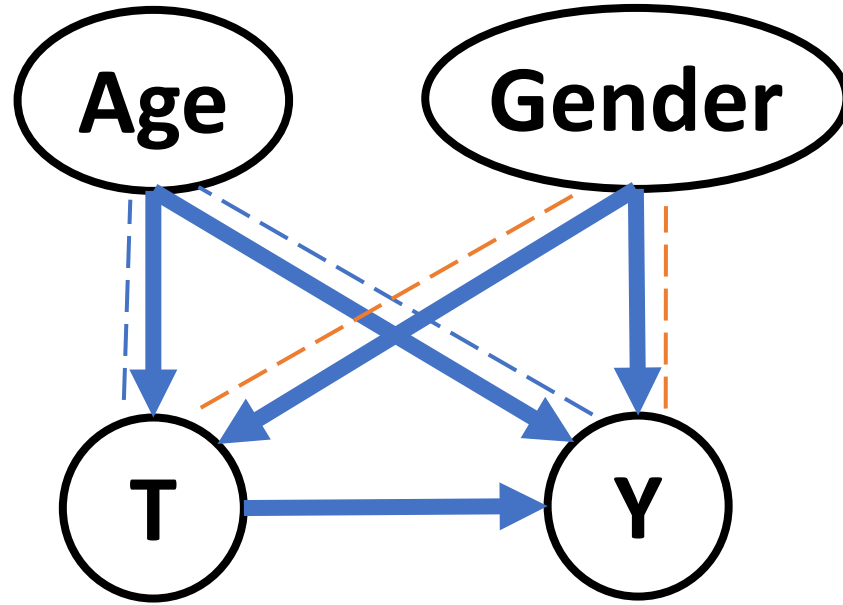
$$P(Y|do(T)) = \sum_x P(Y|T, X = x)P(X = x)$$



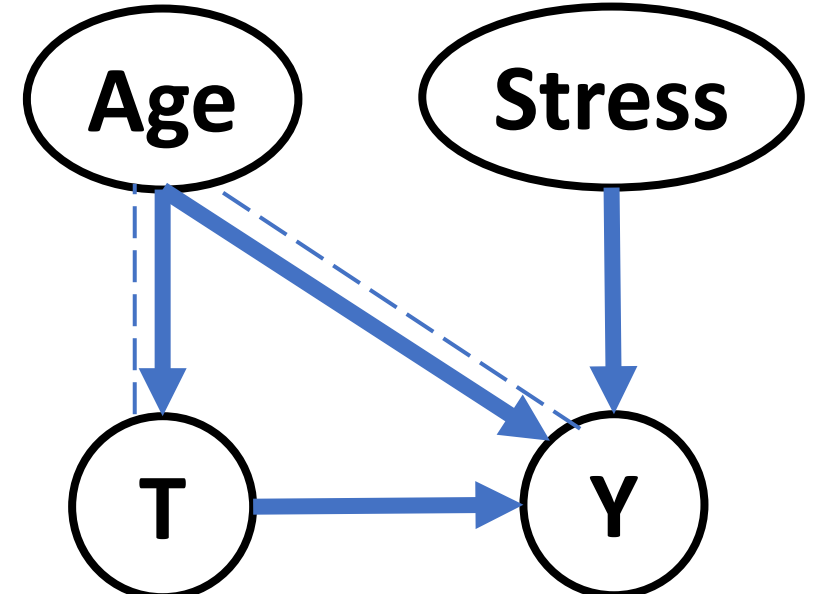
Let us return to our examples



$X = \{Age\}$

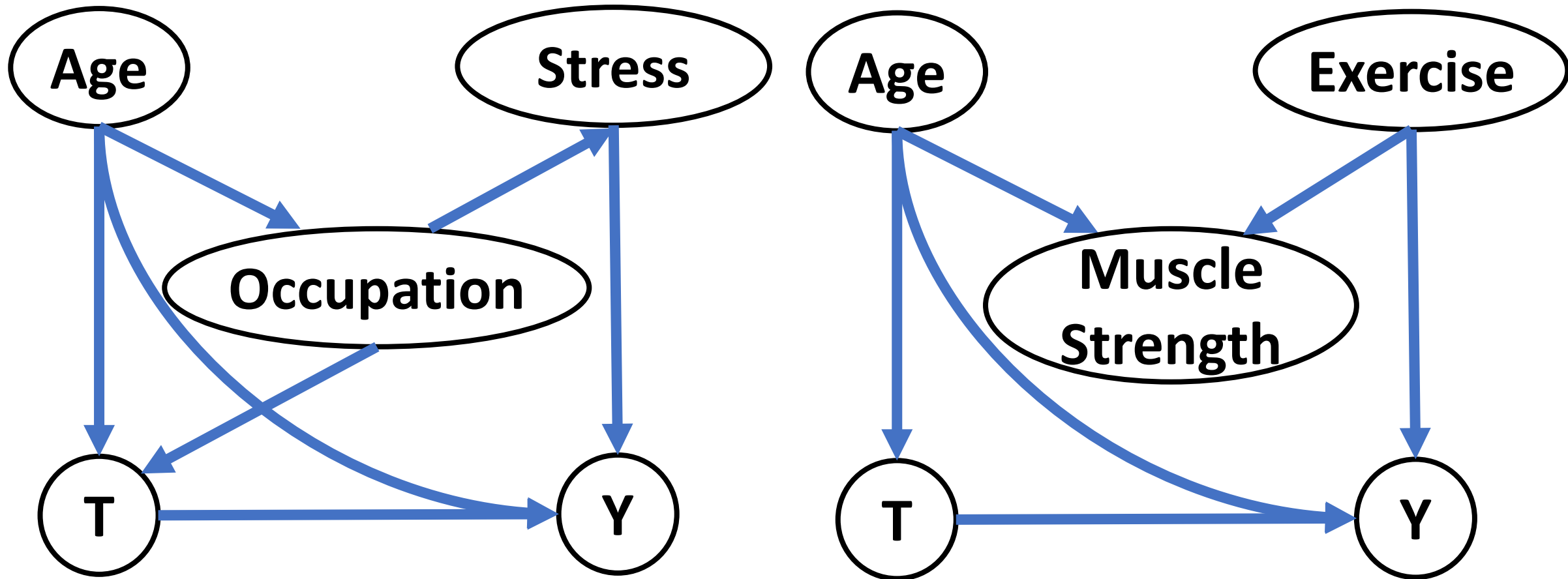


$X = \{Age, Gender\}$



$X = \{Age\}$

Back-door criterion provides a precise way to find variables to condition to



$X = \{Age, Stress\}$   
 $X = \{Age, Occupation\}$

$X = \{Age, Exercise\}$   
 $X \neq \{Age, MuscleStrength\}$

# Both frameworks have merits

**Use structural causal model and do-calculus for**

modeling the problem

making assumptions explicit

identifying the causal effect

**Use potential outcomes-based methods for**

estimating the causal effect

# Recap: Structural Causal Models

- Allow us to make causal assumptions explicit
  - Assumptions are the missing edges!
- Provide language for expressing counterfactuals
- Well-defined mechanisms for reasoning about causal relationships
  - E.g., Backdoor criterion

# Recap: Section 1 - Introduction

- **Causality** is important for decision-making and study of effects
- **Potential Outcomes Framework** gives practical method for estimating causal effects
  - Translates causal inference into counterfactual estimation
- **Unobserved confounds** are a critical challenge
- **Structural Causal Model Framework** gives language for expressing and reasoning about causal relationships

PART I. Introduction to Counterfactual Reasoning

**PART II. Methods for Causal Inference**

**PART III. Large-scale and Network Data**

**PART IV. Broader Landscape**

PART II.  
Methods for Causal  
Inference

**PART II.  
Methods  
for Causal  
Inference**

**Observational Studies**

**Natural Experiments**

**Refutations**



# Review: Treatment, Outcome and Confound

**Goal:** Estimate effect of a treatment  $T$  on an outcome  $Y$

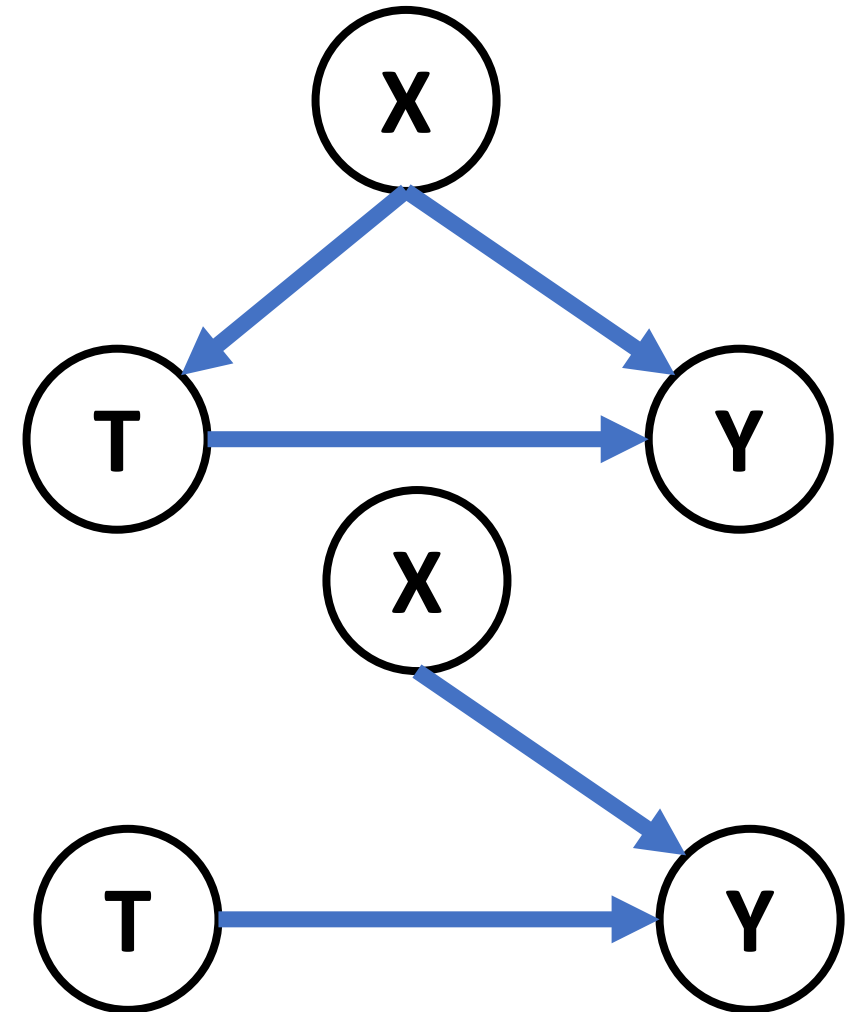
But, confound  $X$  influences both  $T$  and  $Y$

To estimate  $T \rightarrow Y$ , break the dependence  $X \rightarrow T$  (that is,  $T \perp\!\!\!\perp X$ )

- $Y \perp\!\!\!\perp X$  also works, but much less practical.

**Randomized experiments** actively assign treatment  $T$  independent of any confound  $X$

Thus, by construction:  $T \perp\!\!\!\perp X$



# Review: Treatment, Outcome and Confound

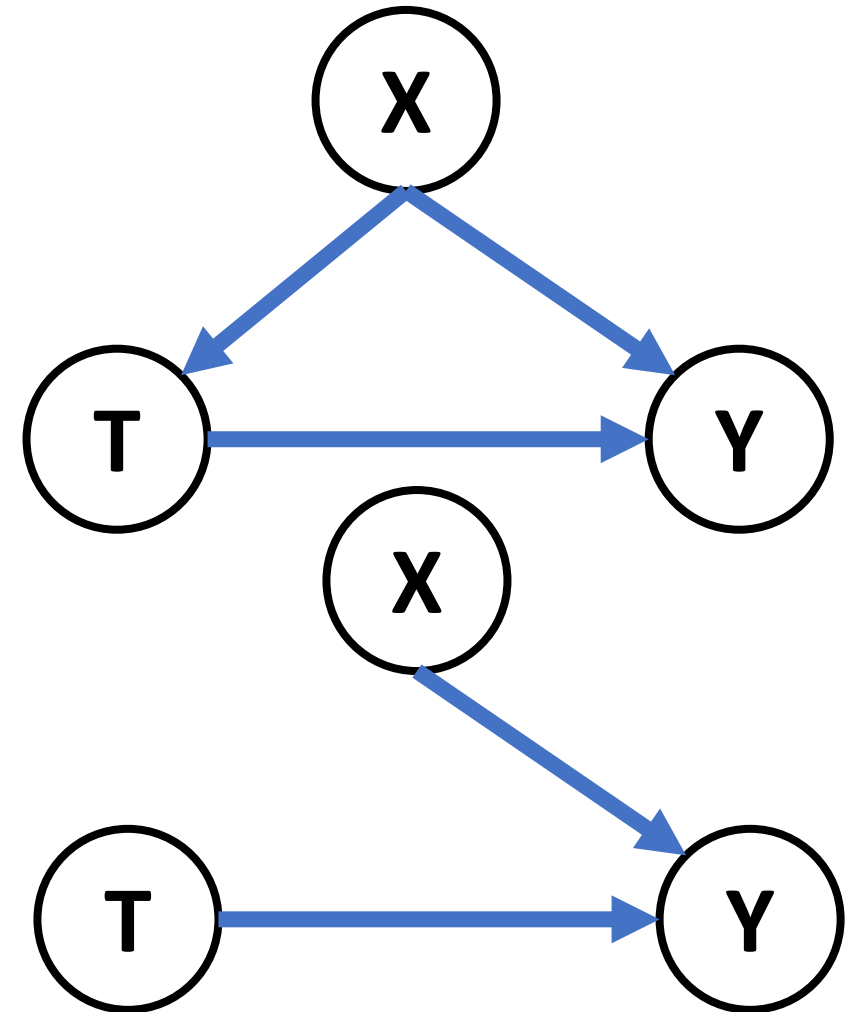
Goal: Estimate effect of a treatment  $T$  on an outcome  $Y$

But, confound  $X$  influences both  $T$  and  $Y$

To estimate  $T \rightarrow Y$ , break the dependence  $X \rightarrow T$  (that is,  $T \perp\!\!\!\perp X$ )

**Randomized experiments** actively assign treatment  $T$  independent of any confound  $X$

Thus, by construction:  $T \perp\!\!\!\perp X$



# Running example

## Review: Exercise, Cholesterol, and Age

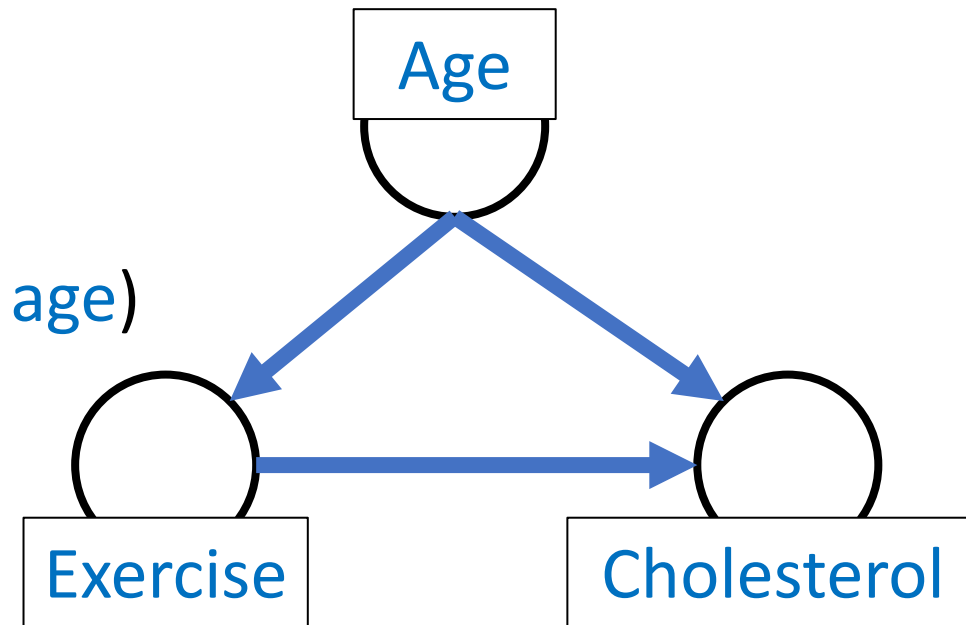
Goal: Estimate effect of **exercise** on **cholesterol**

But, **one's age** influences both **exercise** and **cholesterol**

To estimate **exercise**→**cholesterol**, break the dependence **age**→**exercise** (that is, **exercise**  $\perp$  **age**)

**Randomized experiments** actively assign **exercise** independent of any **age**

Thus, by construction: **exercise**  $\perp$  **age**



# Running example

## Review: Exercise, Cholesterol, and Age

Goal: Estimate effect of **exercise** on **cholesterol**

But, **one's age** influences both **exercise** and **cholesterol**

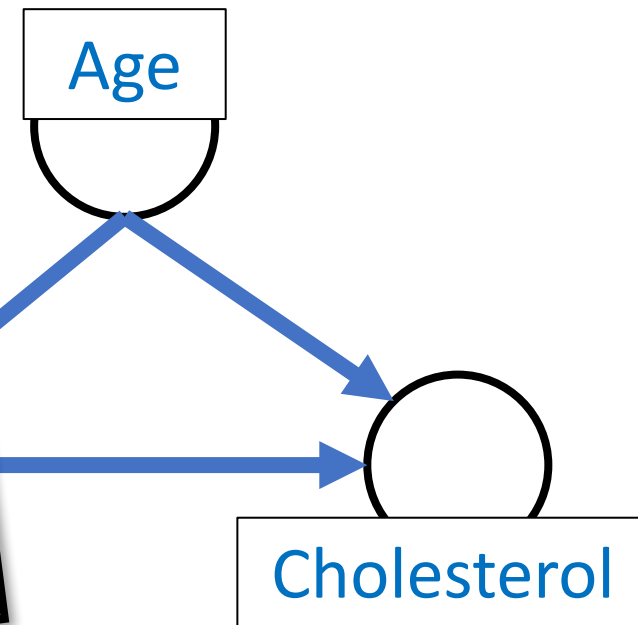
To estimate **exercise**→**cholesterol**, break the dependence **age**→**exercise** (that is, **exercise**  $\perp$  **age**)

Randomized experiment

**exercise**

Thus,

But, what if we cannot actively intervene?



# Part II.A. Observational Studies

*“Simulating  
randomized  
experiments”*

Conditioning on Key Variables

Matching and Stratification

Weighting

Regression

Doubly Robust

Synthetic Controls

# Part II.A. Observational Studies

*“Simulating  
randomized  
experiments”*

Conditioning on Key Variables

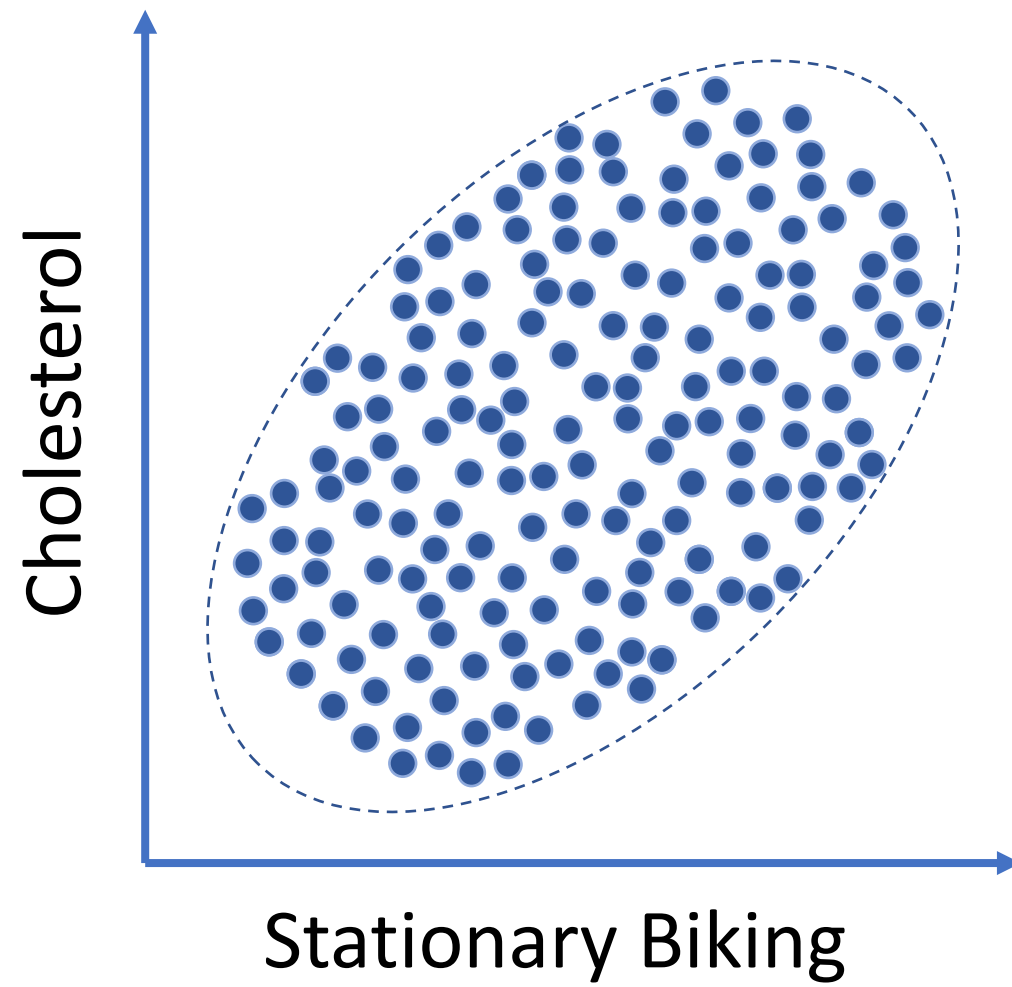
Matching and Stratification

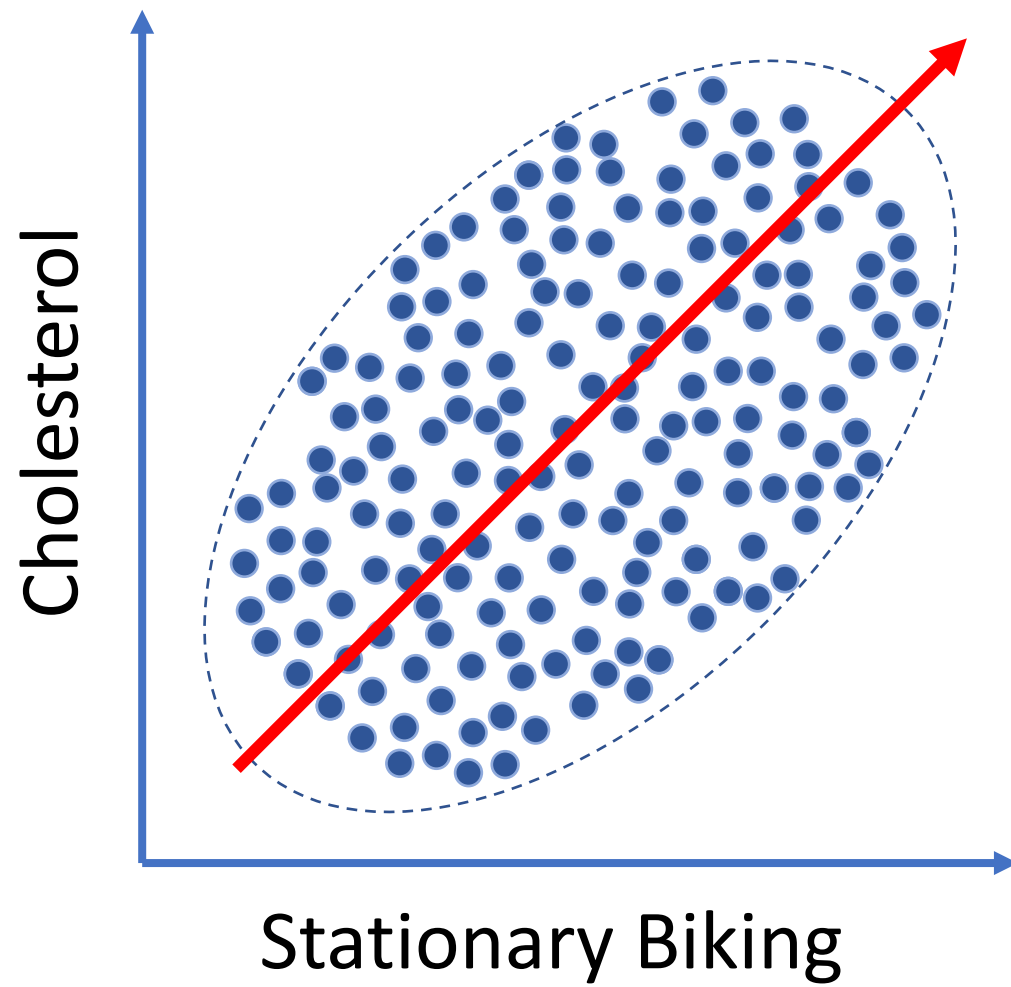
Weighting

Regression

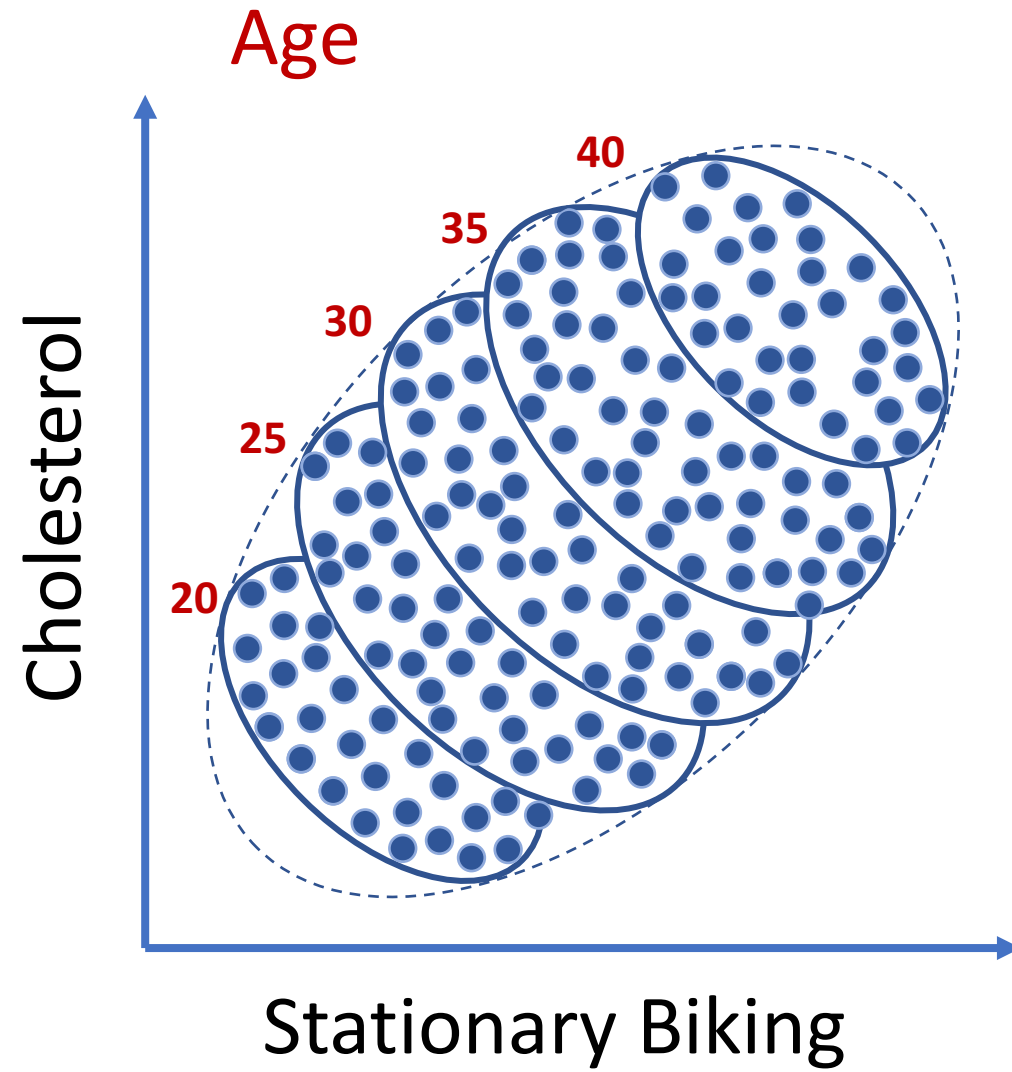
Doubly Robust

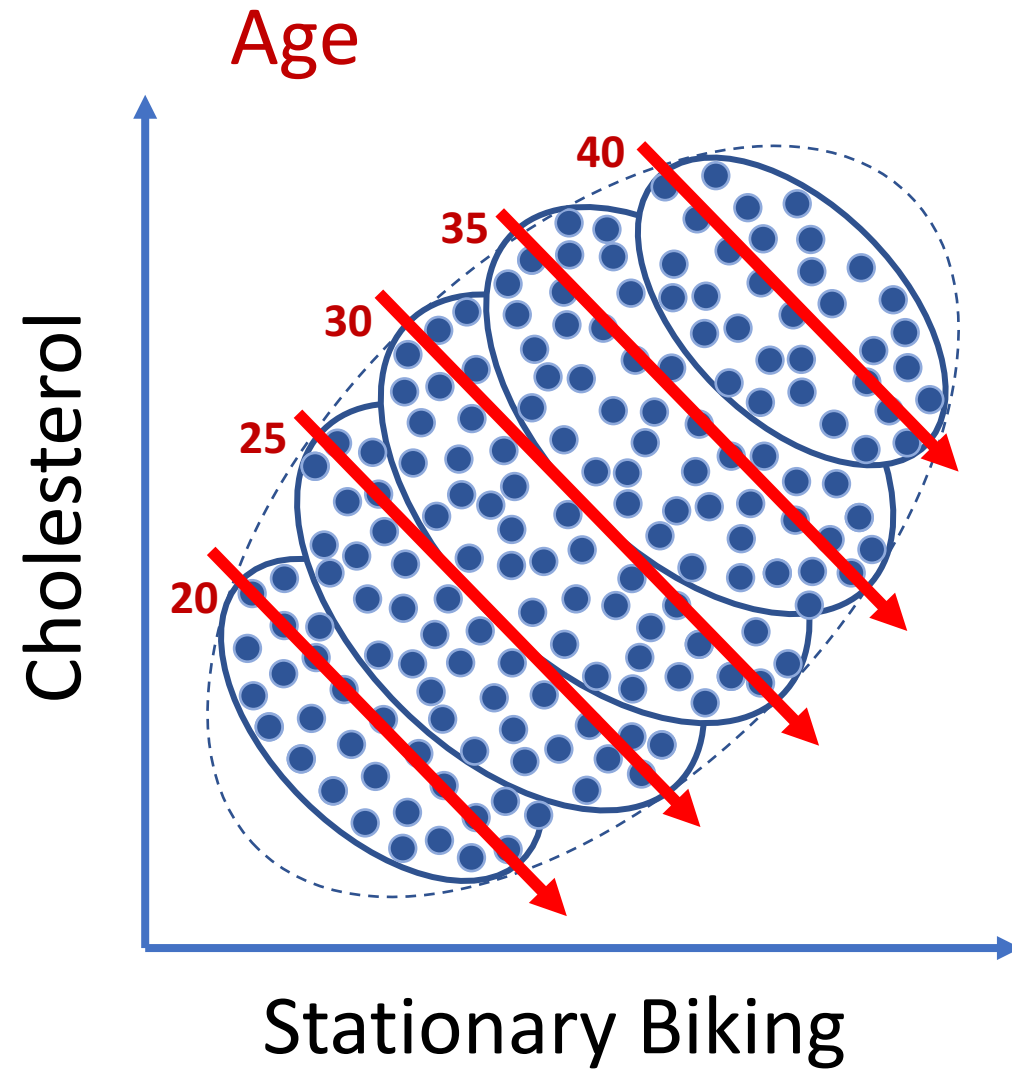
Synthetic Controls











# Recapping what just happened

- At first, more *stationary biking* seems to lead to higher *cholesterol*
- But, we realize that there is a confounder, *age*, that influences both *stationary biking* and *cholesterol*
- We condition on age (by analyzing each age group separately)
- And find stationary biking now seems to lead to lower cholesterol

**Conditioning:**

$$P(\text{Cholesterol} \mid \text{do}(S\_Biking)) = \sum_{age} P(\text{Cholesterol} \mid S\_Biking, age) P(age)$$

# Conditioning

Table 1: Yule-Simpson's Paradox

Population	Survive	Die	Survive Rate
Treatment	20	20	50%
Control	16	24	40%
Male	Survive	Die	Survive Rate
Treatment	18	12	60%
Control	7	3	70%
Female	Survive	Die	Survive Rate
Treatment	2	8	20%
Control	9	21	30%

$$\begin{aligned}\widehat{ACE}_{unadj} &= \widehat{P}(Y = 1 | Z = 1) - \widehat{P}(Y = 1 | Z = 0) \\ &= 0.50 - 0.40 = 0.10 > 0.\end{aligned}$$

$$\begin{aligned}\widehat{ACE}_{adj} &= \{\widehat{P}(Y = 1 | Z = 1, X = 1) - \widehat{P}(Y = 1 | Z = 0, X = 1)\} \widehat{P}(X = 1) \\ &\quad + \{\widehat{P}(Y = 1 | Z = 1, X = 0) - \widehat{P}(Y = 1 | Z = 0, X = 0)\} \widehat{P}(X = 0) \\ &= (0.60 - 0.70) \times 0.5 + (0.20 - 0.30) \times 0.5 \\ &= -0.10 < 0.\end{aligned}$$

male female

# What are the assumptions we made?

- **Assumption:** *age* is the only confounder
  - “*Ignorability*” or “*selection on observables*” assumption
  - How do we know what we must condition on?
- **Assumption:** effect of *stationary biking* doesn’t depend on friends’ exercise
  - Stable Unit Treatment Value (SUTVA) assumption
  - Are there network effects?
- **Assumption:** our observations of exercise/no-exercise cover similar people
  - “*Common support*” or “*Overlap*” assumption
- **Also:** data is not covering all combinations of age and levels of exercise
  - Will our lessons generalize beyond the observed region?

# A1: Ignorability

- Conditional Independence Assumption (CIA)
  - Under random experiments,  $T \perp X$  for both observed and unobserved covariates
  - But conditioning and related techniques can only construct  $T \perp X$  for observed covariates.
- So assume that after conditioning on observed covariates, any unmeasured covariates are irrelevant.

## Ignorability

- Let  $X = \{X_{obs}, X_{unobs}\}$
- Then  $P(Y_T | X_{obs}) = P(Y_T | X_{obs}, T)$  [where  $Y_T = Y|do(T)$ ]

## A2. Stable Unit Treatment Value

The effect of treatment on an individual is independent of whether or not others are treated.

I.e., no spillover or network effects

**SUTVA**

$$P(Y_i | do(T_i, T_j)) = P(Y_i | do(T_i))$$

Example: What is the effect of giving a fax machine to an individual?

- It depends on whether or not

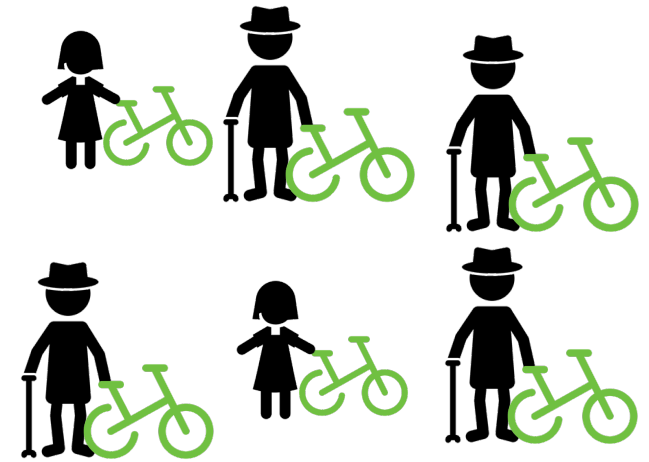
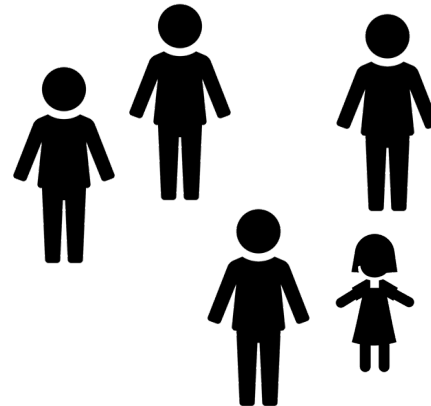
Do people here know / remember what a fax machine is?

# A3. Common support

- The treated and untreated populations have to be similar.
- That is, there should be overlap on observed covariates between treated and untreated individuals.
- Otherwise, cannot estimate counterfactual outcomes.

**Common support**

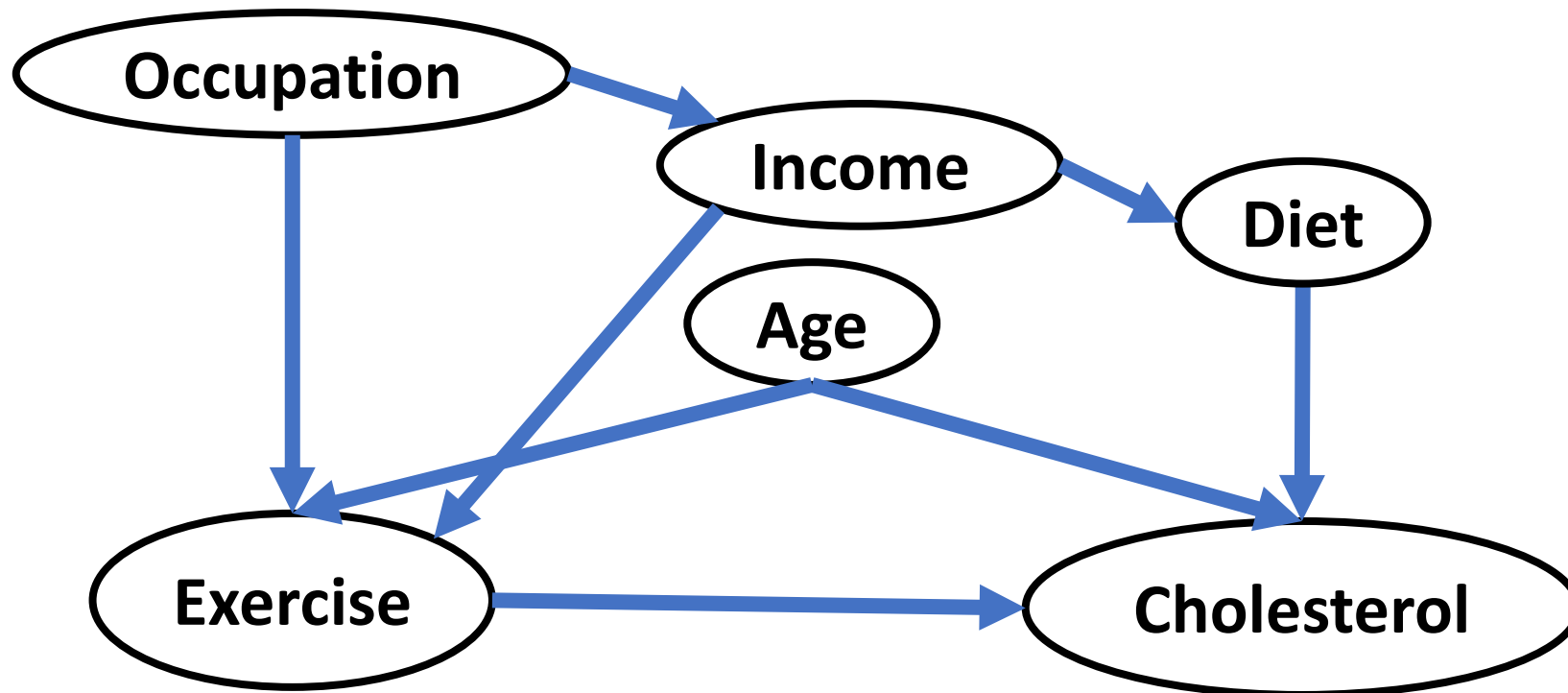
$$0 < P(T = 1 | X = x) < 1$$





# Advanced: How to know we have the right variables? *Backdoor criterion*

1. Use domain knowledge to build a model of the causal graph
2. Condition on enough variables to cover all backdoor paths



**Caveat:** Causal effect only if assumed graphical model is correct

# What we just learned: Simple Conditioning

**Definition** Conditioning calculates treatment effects by identifying groups of individuals with the same covariates, where individuals in one group are treated and in the other group are not.

**Intuition** Conditioning our analysis of  $T \rightarrow Y$  on  $X$  breaks the dependence between confounds  $X$  and the treatment  $T$

**Example** In the cartoon relationship between exercise and cholesterol, age is a confounder, as it influences both levels of exercise and cholesterol. By conditioning analysis on age, we can identify the effect of exercise.

**Keep in mind** How do we know what to condition on?

Grouping becomes harder as dimensionality of  $X$  increases

# Part II.A. Observational Studies

*“Simulating  
randomized  
experiments”*

Conditioning on Key Variables

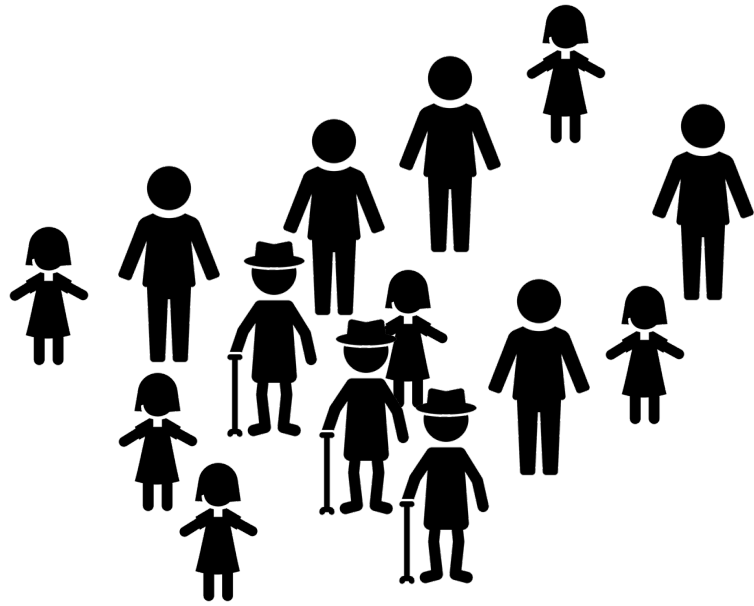
Matching and Stratification

Weighting

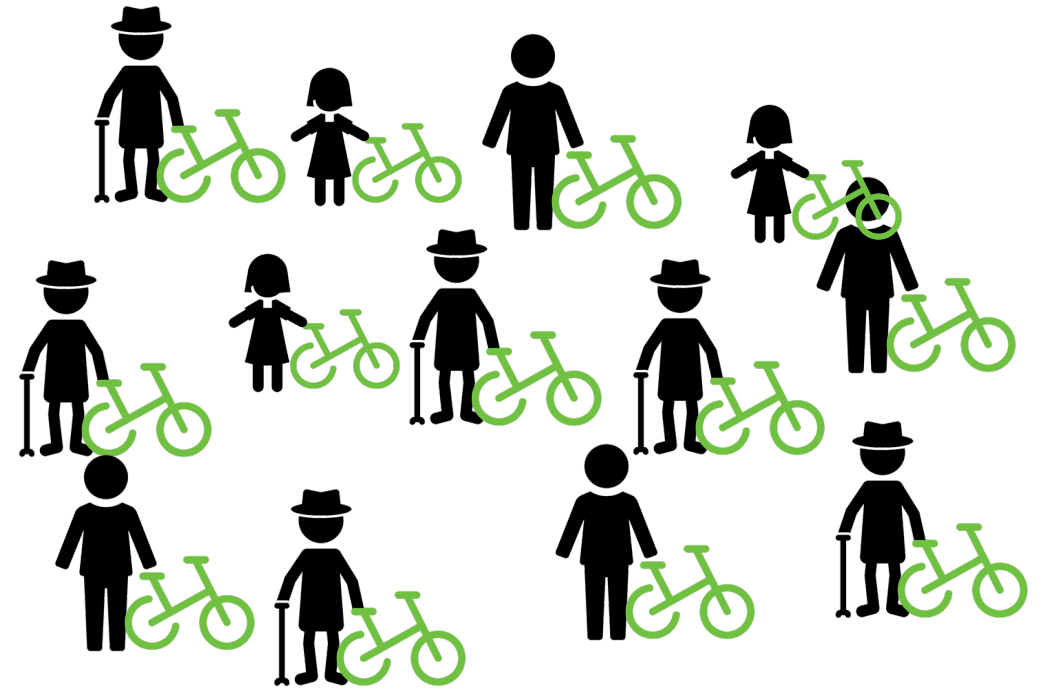
Simple Regression

Doubly Robust

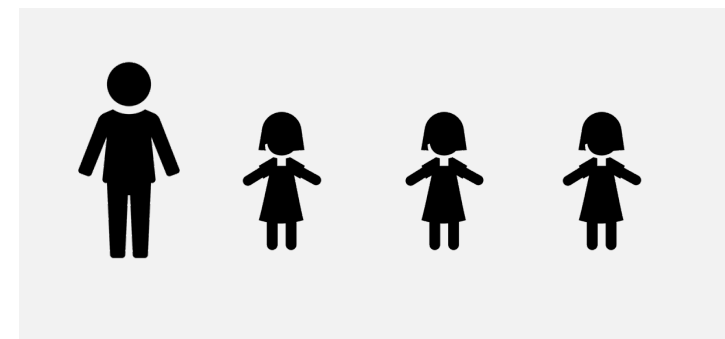
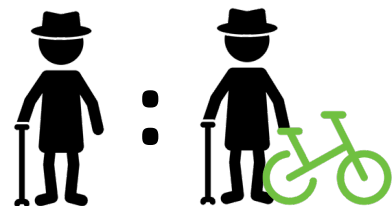
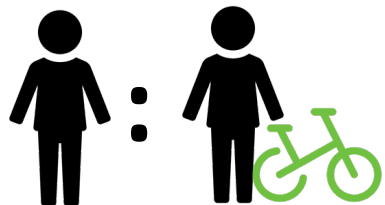
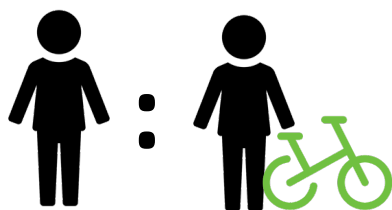
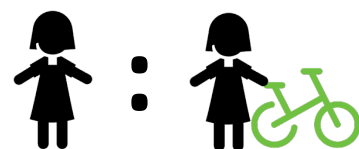
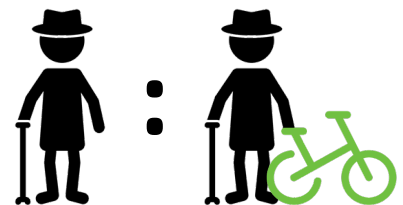
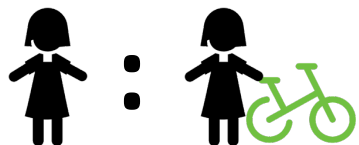
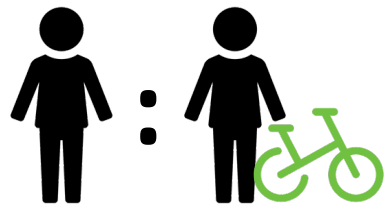
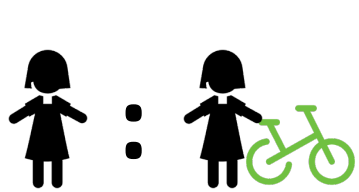
Synthetic Controls



**Avg Cholesterol = 200**



**Avg Cholesterol = 206**



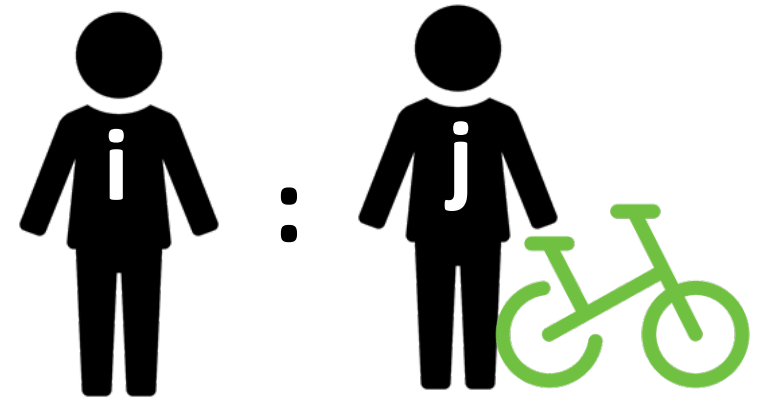
# Matching

Identify pairs of treated and untreated individuals who are very similar or even identical to each other

$$\text{Very similar} ::= \text{Distance}(X_i, X_j) < \epsilon$$

Paired individuals provide the counterfactual estimate for each other.

Average the difference in outcomes within pairs to calculate the *average-treatment-effect on the treated*



# Exact Match

Simple:

$$Distance(\vec{x}_i, \vec{x}_j) = \begin{cases} 0, & \vec{x}_i = \vec{x}_j \\ \infty, & \vec{x}_i \neq \vec{x}_j \end{cases}$$

Use this in low-dimensional settings when overlap is abundant

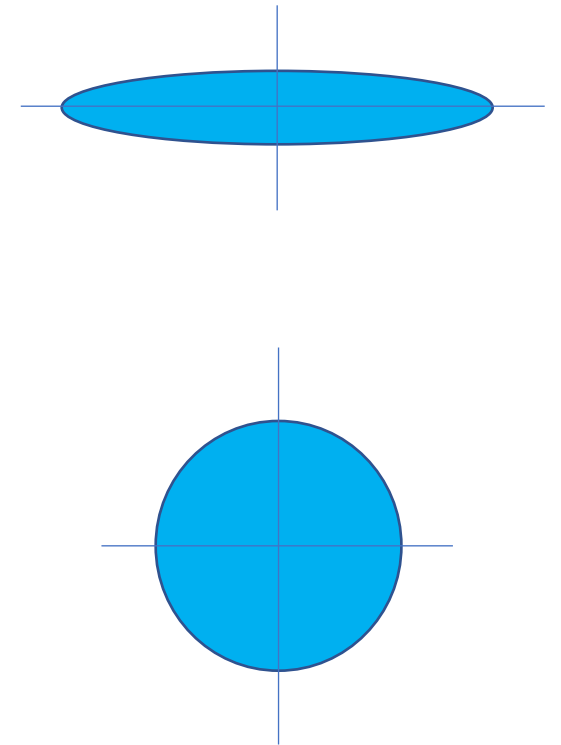
But in most cases, there will be too few exact matches ...

# Mahalanobis Distance

*Mahalanobis distance* accounts for unit differences by normalizing each dimension by the standard deviation.

$$\text{Mahalanobis}(\vec{x}_i, \vec{x}_j) = \sqrt{(\vec{x}_i - \vec{x}_j)^T S^{-1} (\vec{x}_i - \vec{x}_j)}$$

And  $S$  is the covariance matrix.





# Propensity Score

Propensity score is an individual's *propensity to be treated*

$$\hat{e}(X) = P(T = 1|X)$$

- Propensity scores are estimated or modeled, *not observed*.
- Rare exception is if you know likelihood of randomized assignment

**Breaks influence of confound  $X$ ,  
allowing estimate of  $T \rightarrow Y$**

Propensity scores subdivide observational data s.t.  $T \perp\!\!\!\perp X \mid \text{score}$

# How to match with propensity score

1. Train a machine learning model to predict treatment status
  - **Supervised learning:** We are trying to predict a known label (treatment status) based on observed covariates.
  - Conventionally, use a logistical regression model, but SVM, GAMs, are fine
  - But score must be well-calibrated. I.e.,  $(100 * p)\%$  of individuals with score of  $p$  are observed to be treated

2. Distance is the difference between propensity scores

$$Distance(\vec{x}_i, \vec{x}_j) = |\hat{e}(\vec{x}_i) - \hat{e}(\vec{x}_j)|$$

# Propensity score, FAQ

**Q: Wait, why does this work?**

A: Individuals with similar covariates get similar scores, and all individuals mapped to a similar score have similar treatment likelihoods.

**Q: What if my propensity score is not accurate? (i.e., can't tell who is treated)**

A: That's ok. The role of the model is to balance covariates given a score; not to actually identify treated and untreated.

**Q: What if my propensity score is very accurate? (i.e., *can* tell who is treated)**

A: Means we cannot disentangle covariates from treatment status. Any effect we observe could be due either to the treatment or to the correlated covariate.

Consider redefining the treatment or general problem statement. **Don't** dumb down model!

# Propensity score matching python code

```
# learn propensity score model
psmodel = linear_model.LinearRegression()
psmodel.fit(covariates, treatment_status)
data['ps'] = psmodel.predict(covariates)
# find nearest neighbor matches
controlMatcher = NearestNeighbors().fit(untreated['ps'])
distances, matchIndex = controlMatcher.kneighbors(treated['ps'])
# iterate over matched pairs and sum difference in outcomes
for i in range(numtreatedunits):
    treated_outcome = treated.iloc[i][outcome_name].item()
    untreated_outcome = untreated.iloc[matchIndex[i]][outcome_name].item()
    att += treated_outcome - untreated_outcome
# normalize
att /= numtreatedunits
```

# Advanced: Matching

- When matching, should we allow replacement?
  - It's a bias / variance trade-off
- When matching, what if nearest neighbor is far away?
  - Use a caliper threshold to limit acceptable distance
- What if not all treated individuals are matched to untreated?
  - This will bias results. Consider redefining original cohort / population to cleanly exclude treated who won't have matches in untreated population.
- Treatment should be a binary point treatment
  - Advanced variants allow multi-dose, and other treatment regimens

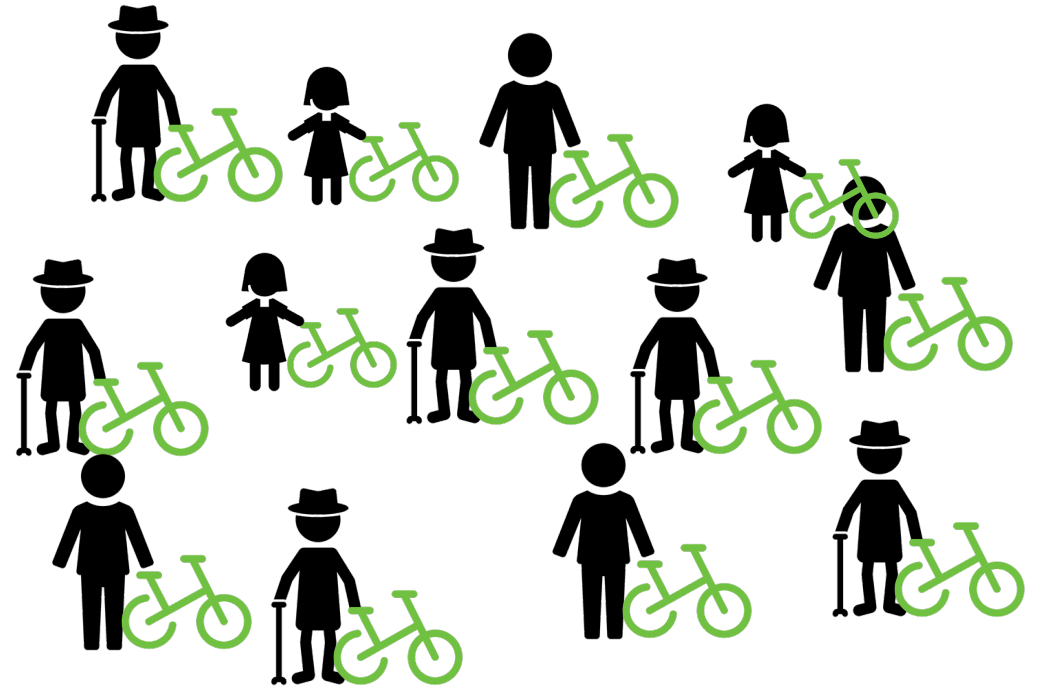
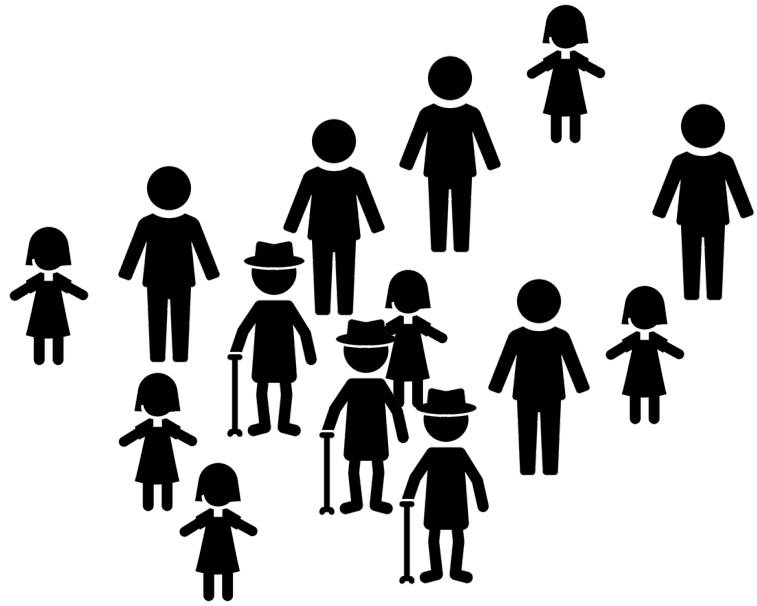
# What we just learned: Matching

**Definition** Matching calculates treatment effects by identifying pairs of similar individuals, where one is treated and the other is not.

**Intuition** The paired individuals stand-in as the counterfactual observations for one another.

**Example** In our cartoon, we create pairs of individuals matched exactly on their age. More generally, we can use Mahalanobis distance or propensity score matching to find similar individuals to be matched.

**Keep in mind** Matching calculates the treatment effect on the treated population. We do not know what might happen if people who would never get treatment are suddenly treated.



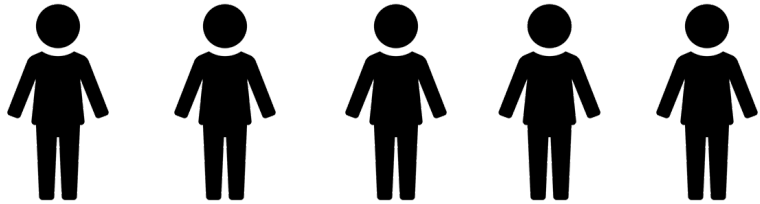
180



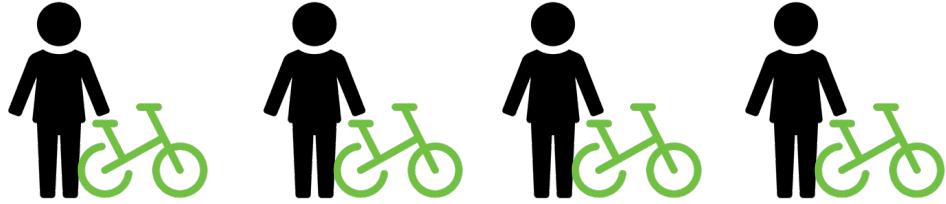
180



200



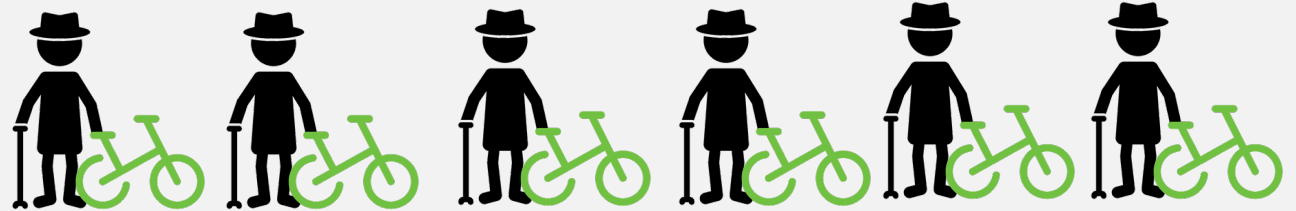
190



240



230





# From Matching to Stratification

- 1: 1 matching generalizes to *many:many* matching.
- Stratification identifies paired *subpopulations* whose covariate distributions are similar.
- There can still be error, if strata are too large.

# How to stratify with propensity score

1. Train a machine learning model to predict treatment status
  - **Supervised learning:** We are trying to predict a known label (treatment status) based on observed covariates.
  - Conventionally, use a logistical regression model, but SVM, GAMs, are fine
  - But score must be well-calibrated. I.e.,  $(100 * p)\%$  of individuals with score of  $p$  are observed to be treated

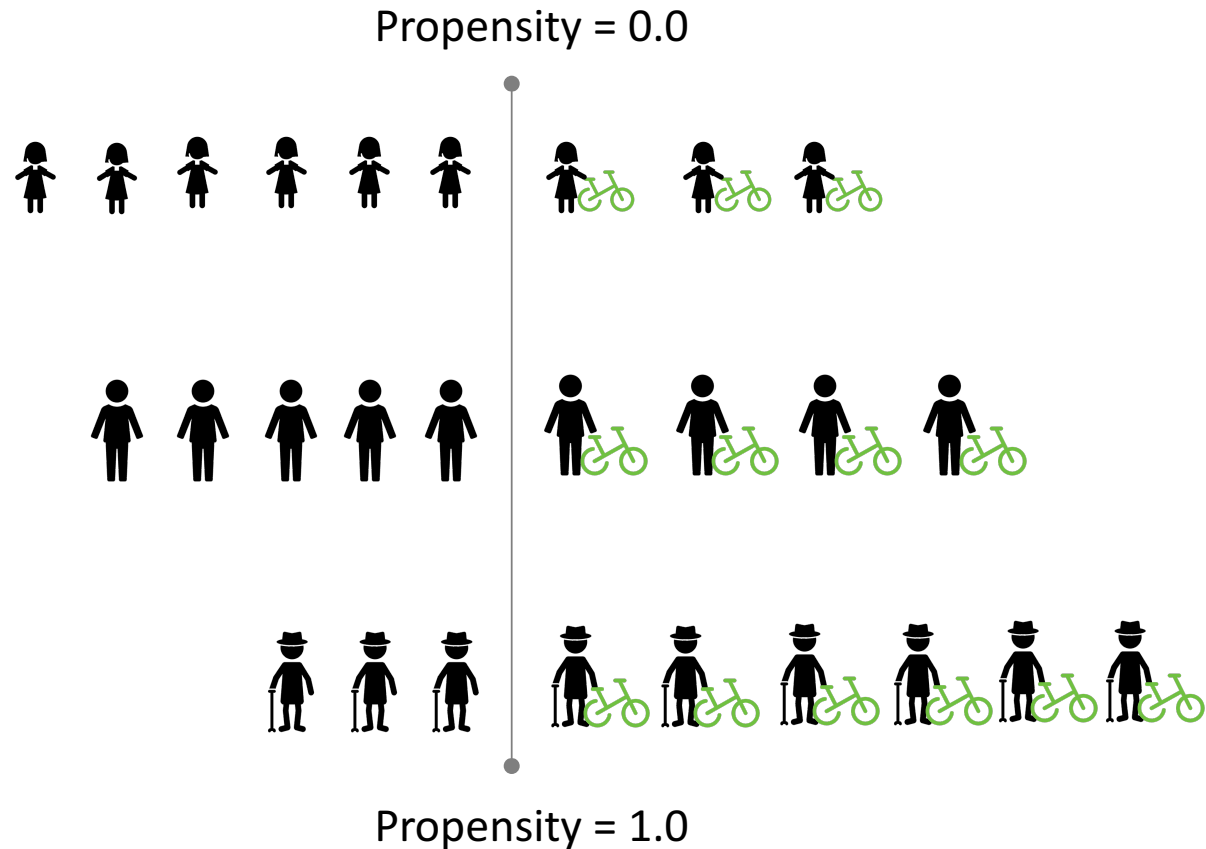
2. Distance is the difference between propensity scores

$$Distance(\vec{x}_i, \vec{x}_j) = |\hat{e}(\vec{x}_i) - \hat{e}(\vec{x}_j)|$$

# Propensity Score Stratification

We can use propensity score to stratify populations

1. Calculate propensity scores per individual as in matching.
2. But instead of matching, stratify based on score.
3. Calculate average treatment effect as weighted average of outcome differences per strata.
4. Weight by number of treated in the population for ATE on treated.



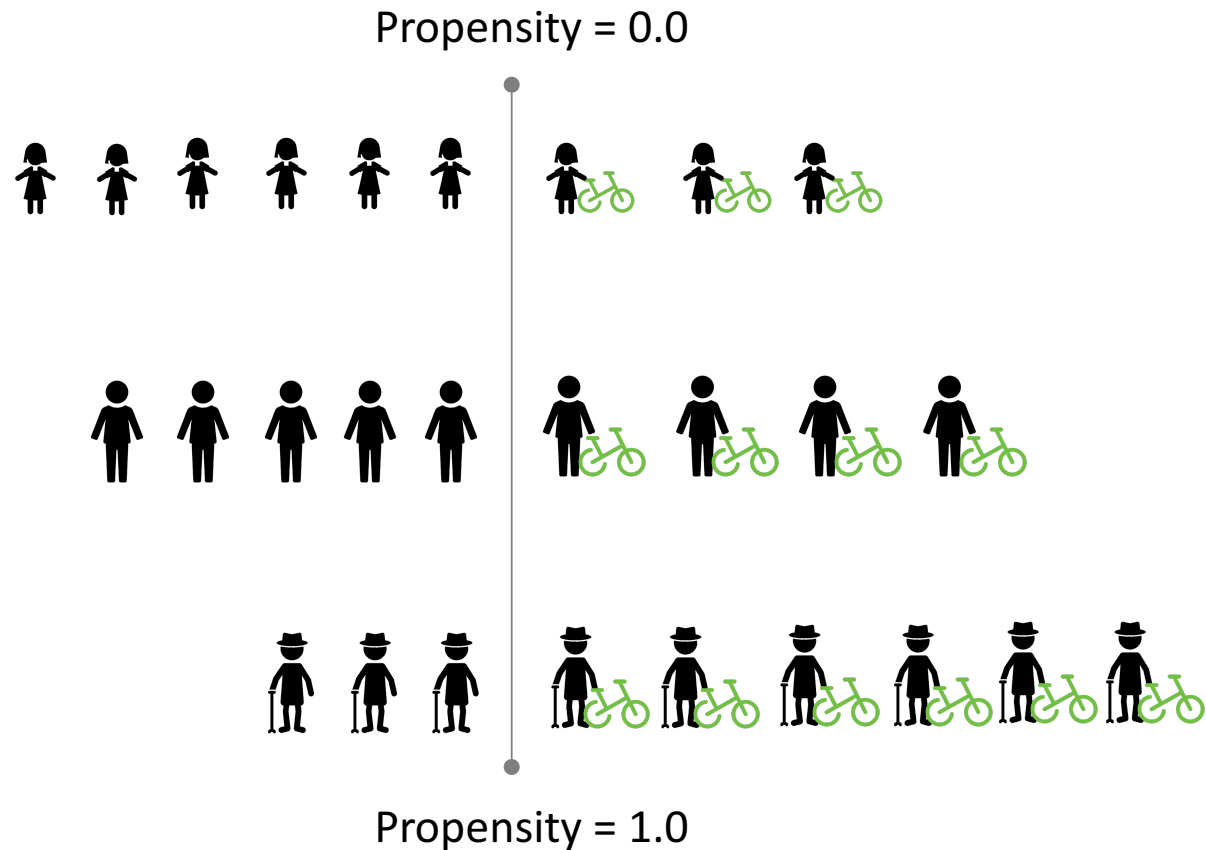
# Propensity Score Stratification

$$ATE = \sum_{s \in \text{strata}} \frac{1}{N_{s,T=1}} (\bar{Y}_{s,T=1} - \bar{Y}_{s,T=0})$$

where,

$\bar{Y}_{s,T}$  is the average outcome at strata  $s$  and treatment status  $T$

And  $N_{s,T=1}$  is the number of treated individuals in strata  $s$



# Propensity score stratification python code

```
# build propensity score model and assign each item a score as earlier..

# create a column 'strata' for each element that marks what strata it belongs to
data['strata'] = ((data['ps'].rank(ascending=True) / numrows) * numStrata).round(0)
data['T_y'] = data['T'] * data['outcome'] # T_y = outcome iff treated
data['Tbar'] = 1 - data['treated'] # Tbar = 1 iff untreated
data['Tbar_y'] = data['Tbar'] * data['outcome'] # Tbar_y = outcome iff untreated
stratified = data.groupby('strata')
# sum weighted outcomes over all strata (weight by treated population)
outcomes = stratified.agg({'T':['sum'],'Tbar':['sum'],'T_y':['sum'],'Tbar_y':['sum']})
# calculate per-strata effect
outcomes['T_y_mean'] = outcomes['T_y_sum'] / outcomes['T']
outcomes['Tbar_y_mean'] = outcomes['Tbar_y_sum'] / outcomes['dbar_sum']
outcomes['effect'] = outcomes['T_y_mean'] - outcomes['Tbar_y_mean']
# weighted sum of effects over all strata
att = (outcomes['effect'] * outcomes['T']).sum() / totaltreatmentpopulation
```

# P.S. Stratification, Practical Considerations

- How many strata do we pick?
  - Scale will depend on data. Want each stratum to have enough data in it.
  - Conventional, small-data literature (e.g., ~100 data points) picked 5.
  - With 10k to 1M or more data points, I pick 100 to 1000 strata.
  - Set strata boundaries to split observed population evenly
  - Aside: why not always pick a small number of strata? It's a bias-variance trade-off...
- What if there aren't enough treated or untreated individuals in some of my stratum to make a meaningful comparison?
  - This often happens near propensity score 0.0 and near 1.0
  - Drop ("Clip") these strata from analysis. Technically, you are now calculating a local-average-treatment-effect.

# What we just learned: Stratification

**Definition** Stratification calculates treatment effects by identifying groups of individuals with similar distributions of covariates, where individuals in one group are treated and in the other group are not.

**Intuition** The difference in average outcome of paired *groups* tells us the effect of the treatment on that subpopulation. Observed confounds are balanced, due to covariate similarity across paired groups.

**Example** In our cartoon example, we stratified based on propensity score into 3 strata. ATE is the weighted sum of differences in avg outcomes in each strata.

**Keep in mind** Make sure there are enough comparable individuals in each strata

# Part II.A. Observational Studies

*“Simulating  
randomized  
experiments”*

Conditioning on Key Variables

Matching and Stratification

Weighting

Simple Regression

Doubly Robust

Synthetic Controls



# Weighting: An alternative to conditioning

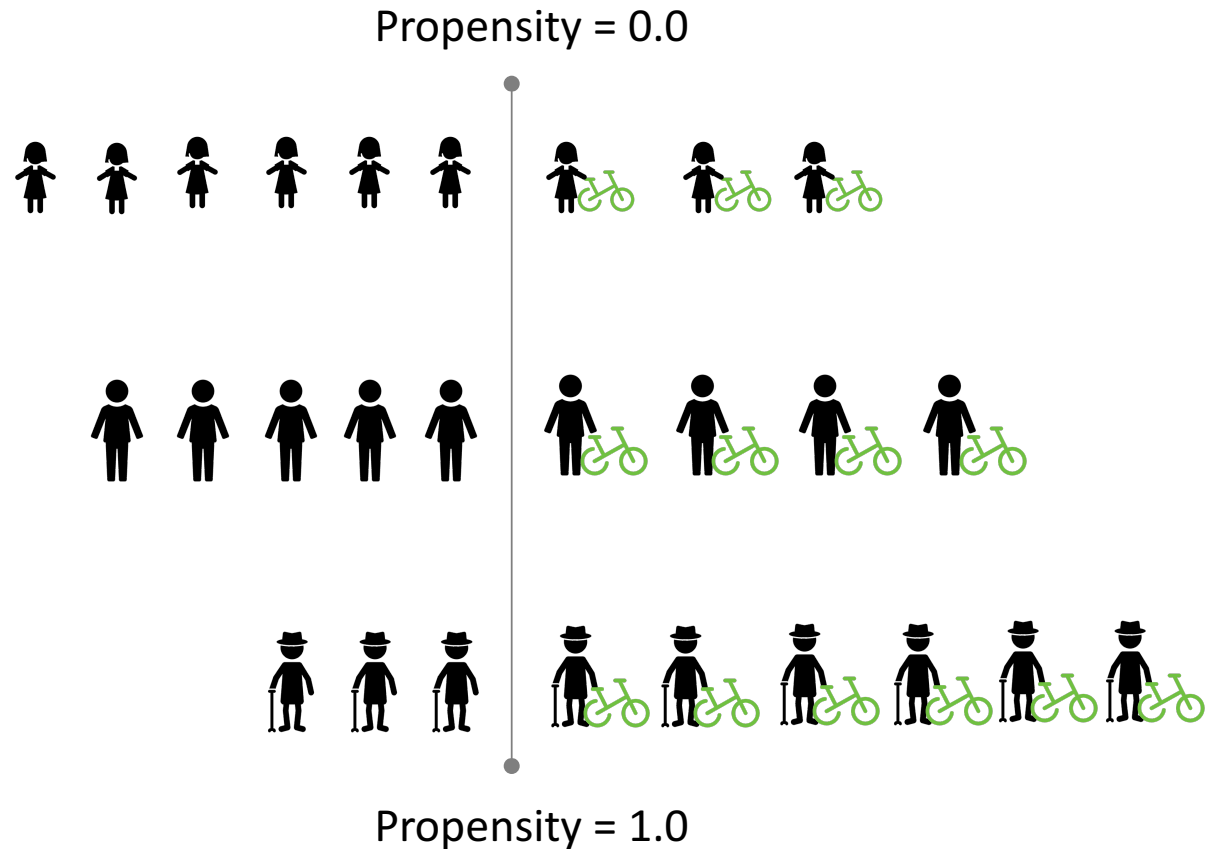
*What if we assign weights to observations to simulate randomized experiment?*

Stratification weights strata results by number of treated

Weighting by treated population  $\sim$  weighting by propensity score.

Generalized weighting: Calculate effect by weighted sum over all individual outcomes

Many weighting methods to generate a balanced dataset



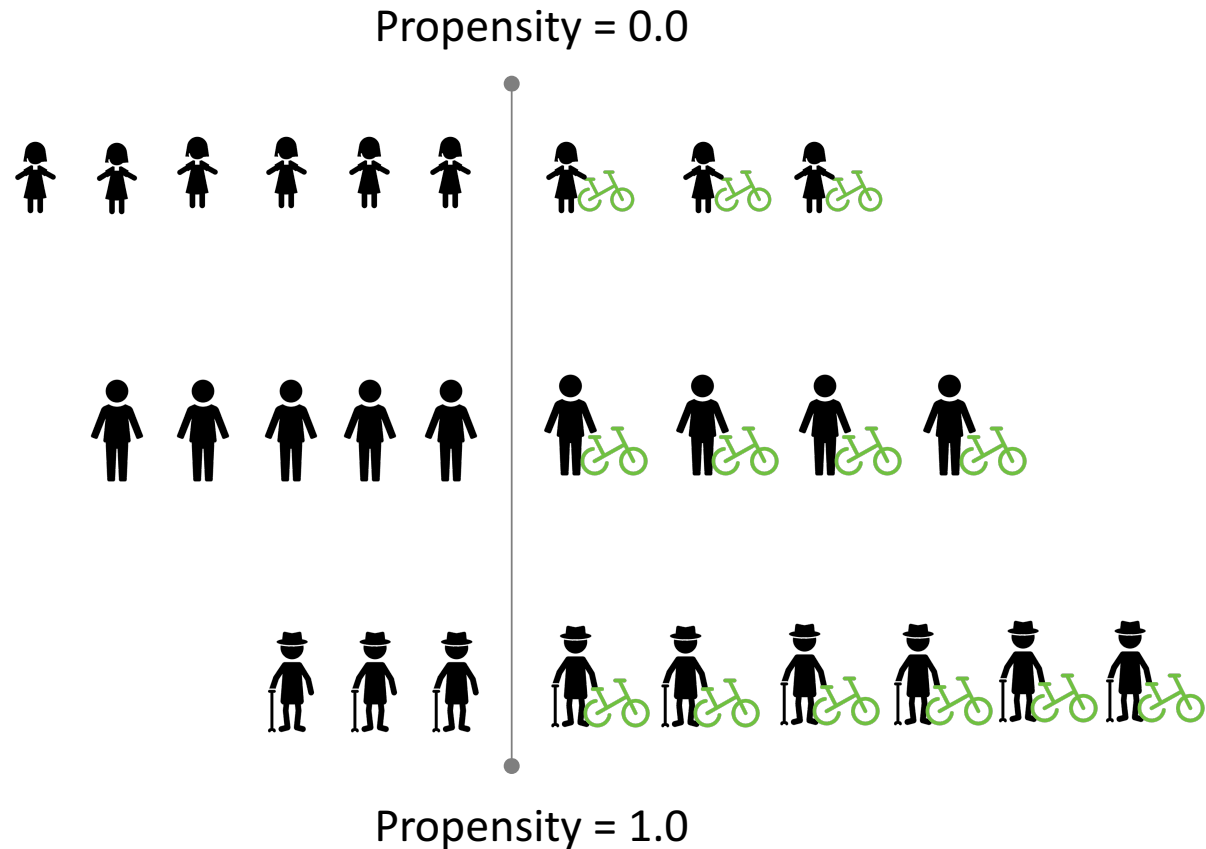
# Weighting

Stratification weights strata results by number of treated

Weighting by treated population  
~ weighting by propensity score.

Generalized weighting: Calculate effect by weighted sum over all individual outcomes

Many weighting methods to generate a balanced dataset



# Weighting

$$ATE = \frac{1}{N_{T=1}} \sum_{i \in \text{treated}} w_i Y_i - \frac{1}{N_{T=0}} \sum_{j \in \text{untreated}} w_j Y_j$$

Inverse Probability of Treatment Weighting (IPTW)

$$w_i = \frac{T}{e} + \frac{1-T}{1-e};$$
$$N_{T=1} = \sum \frac{T}{e}; \quad N_{T=0} = \sum \frac{1-T}{1-e}$$

# Weighting: Caveats and Practical notes

- High variance when  $e$  close to 0 or 1  
A single value can derail the estimate.
- Many heuristics for clipping weights; stabilizing weights; etc.
- Assumes propensity score model is correctly specified (i.e., that  $e$  is correctly estimated for all individuals)
- Variants of weighting: calculate average treatment effect on treated

# What we just learned: Weighting

**Definition** Weighting calculates average treatment effect as the difference between the weighted sum of the treated and untreated populations

**Intuition** Weights on each individual act to balance the distribution of covariates in the treated and untreated groups. (i.e., break the dependence between treatment status and covariates)

**Keep in mind** High variance when propensity scores are very high or very low  
Many variants of weighting schemes

## Part II.A. Observational Studies

*“Simulating  
randomized  
experiments”*

Conditioning on Key Variables

Matching and Stratification

Weighting

Simple Regression

Doubly Robust

Synthetic Controls

# Regression (or supervised learning)

In regression analysis, we build a model of  $Y$  as a function of covariates  $X$  and  $T$ , and interpret coefficients of  $X$  and  $T$  causally:

$$E(Y|X, T) = \alpha_1 X_1 + \alpha_2 X_2 + \cdots \alpha_n X_n + \alpha_T T$$

Example:

$$\textit{Cholesterol} = \alpha_{\textit{age}} \textit{Age} + \alpha_{\textit{exercise}} \textit{Exercise}$$

Model is fit with standard methods (e.g., MLE)

The bigger  $\alpha$  is, the stronger the causal relationship to  $Y$

# Regression warnings

Causal interpretation of regressions requires many assumptions.

Threats to validity include:

- **Model correctness:** e.g., what if we use a linear model and causal relationship is non-linear
- **Multicollinearity:** if covariates are correlated, can't get accurate coefficients
- **Ignorability (Omitted variables):** Omission of confounds will invalidate findings



# What we just learned: Regression

**Definition** Use a regression-based causal analysis, we interpret coefficients as the strength of causal relationship

**Example** *Modeling cholesterol as a function of exercise and age*

**Keep in mind** Analysis must be carefully designed to ensure causal interpretability, avoiding collinearity and including all relevant confounds

Avoid unless you are absolutely sure of what you are doing.

# Part II.A. Observational Studies

*“Simulating  
randomized  
experiments”*

Conditioning on Key Variables

Matching and Stratification

Weighting

Simple Regression

Doubly Robust

Synthetic Controls

# Doubly robust: Best of both worlds?

- Both propensity score weighting and regression models require correctly specified models
  - E.g., if propensity score or regression is modeled as a linear combination, but is non-linear, then it is not correctly specified
- Doubly robust methods combine “best of” propensity score and regression methods
- If either propensity score or regression is correctly specified, then doubly robust is correct.

# DR: Combines 3 components

Learn 3 models:

**1,2:** Models of outcome given treatment and covariates:  $\hat{Y}_{T=0}$  ,  $\hat{Y}_{T=1}$

**3:** Propensity of treatment given covariates:  $\hat{e}$

Combine to calculate doubly robust estimators,  $DR_1$  and  $DR_0$ , for each individual:

$$DR_1 = \begin{cases} \frac{Y}{\hat{e}} - \frac{\hat{Y}_{T=1}(1 - \hat{e})}{\hat{e}}, & T = 1 \\ \hat{Y}_{T=1}, & T = 0 \end{cases}$$

$$DR_0 = \begin{cases} \hat{Y}_{T=0}, & T = 1 \\ \frac{Y}{1 - \hat{e}} - \frac{\hat{Y}_{T=1}\hat{e}}{1 - \hat{e}}, & T = 0 \end{cases}$$

Finally, calculate mean  $\overline{DR_1}$  and  $\overline{DR_0}$  over the whole study population, and take difference as the causal effect of  $T$

# Doubly Robust: Caveat

If either propensity score or regression is correctly specified, then doubly robust is unbiased.

Seems like doubly robust should be strictly better (less biased) than either propensity score weighting or regression

But, if both propensity score or regression are *slightly* incorrect, then doubly robust estimator may become *very* biased

# What we just learned: Doubly Robust

**Intuition** Combine propensity score weighting and regression models to provide unbiased estimate when either propensity score or regression is correctly specified

**Keep in mind** Fundamental assumptions (ignorability, etc) must still hold. If both models are slightly incorrect, doubly robust estimator can be more biased

# Part II.A. Observational Studies

*“Simulating  
randomized  
experiments”*

Conditioning on Key Variables

Matching and Stratification

Weighting

Simple Regression

Doubly Robust

Synthetic Controls

# Synthetic control

All previous methods require that we observe both *treated* and *untreated* individuals

What if we are analyzing a scenario where everyone is treated?

*E.g.*, effect of a large marketing campaign, or a global policy change?

Pre/Post comparison is option, but not robust to dynamics, seasonality, ...

Alternative: Build *synthetic controls* that estimate what  $\bar{Y}_{T=0}$  would have been for a population were it not for treatment



# Synthetic controls: Intuition

1. *Decide what the treatment will be*

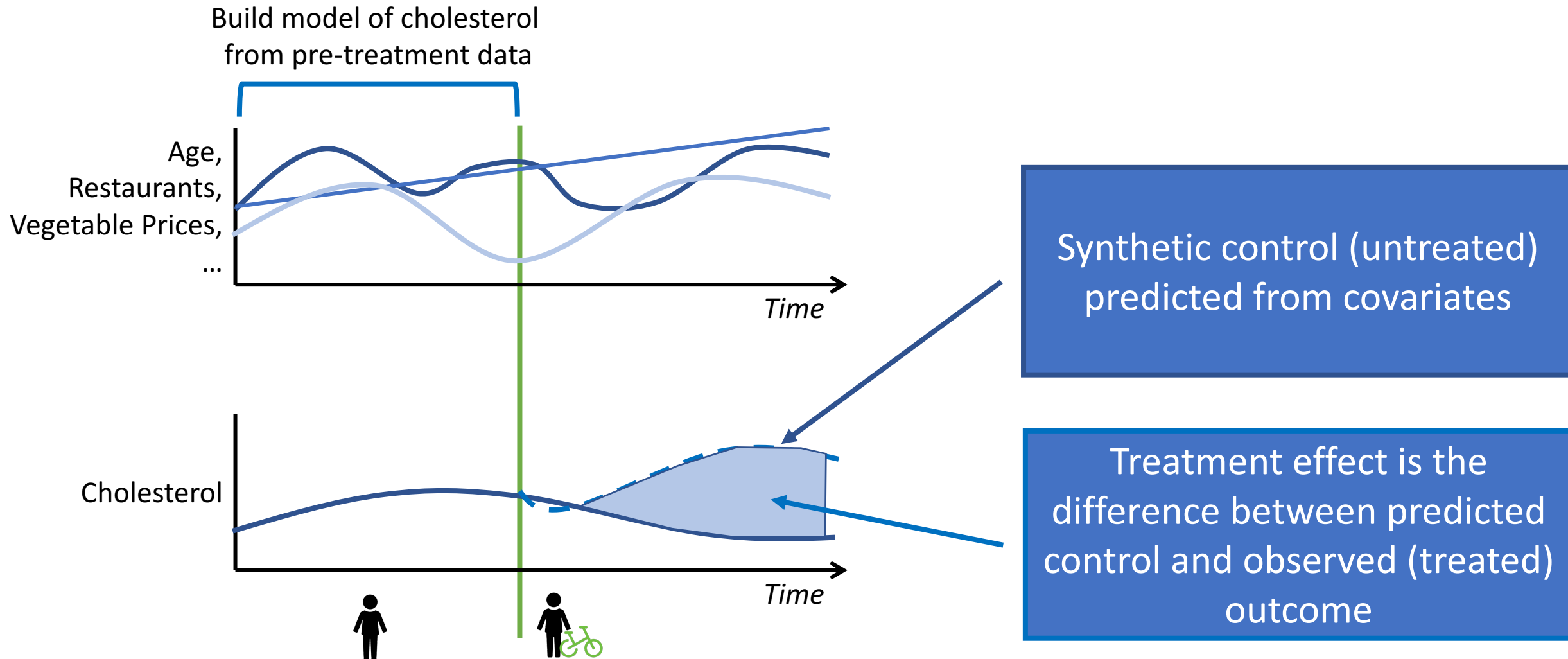
2. *Pre-treatment stage: Observe the world for a while*

- Record the outcome we care
- Record covariates that can help us predict our observed outcome, but will not be effected by the treatment. Use domain-knowledge / theory to identify these covariates.
- Learn a model that predicts outcome based on covariates.

3. *Post-treatment stage:*

- Keep recording outcome. This is now the treated outcome.
- Predict untreated outcome using learned model and current covariates
- ATE = Difference between observed outcome and prediction of untreated outcome

# Example: policy change to encourage exercise



# What we just learned: Synthetic Controls

**Definition** Calculate treatment effect by comparing observed outcomes of treated population with synthetic (predicted) outcomes of an untreated population

**Intuition** If we can measure covariates that are unaffected by the treatment and predictive of untreated outcomes, then we can build a synthetic control

**Example** Predicting effect of global policy change to encourage exercise on population-wide cholesterol

**Keep in mind** Ignorability assumption must still hold;  
Relatedly, be concerned about generalizability/robustness of learned outcome model

**PART II.  
Methods  
for Causal  
Inference**

Observational Studies

Natural Experiments

Refutations

Part II.B.  
Natural  
Experiments

Simple natural  
experiment

Instrumental Variables

Regression  
Discontinuities

# Natural experiments: What can we do without ignorability?

Rather than assume ignorability over the entire dataset, find data subsets that approximate an experiment.

“Natural” → as if Nature *conducted an experiment* for you

**Common sources:** Prior A/B tests, Lottery, any randomized policy, an external shock to the treatment.

Allows common causes of T and Y, as long as the source is not affected by them.

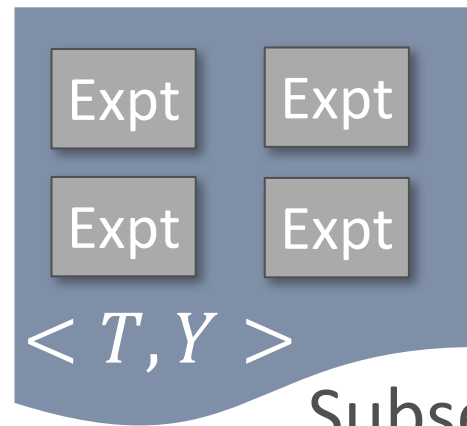
# Finding a natural experiment



Full dataset

$$y = f(t, x)$$

$$t = g(x)$$



Subsets of the data

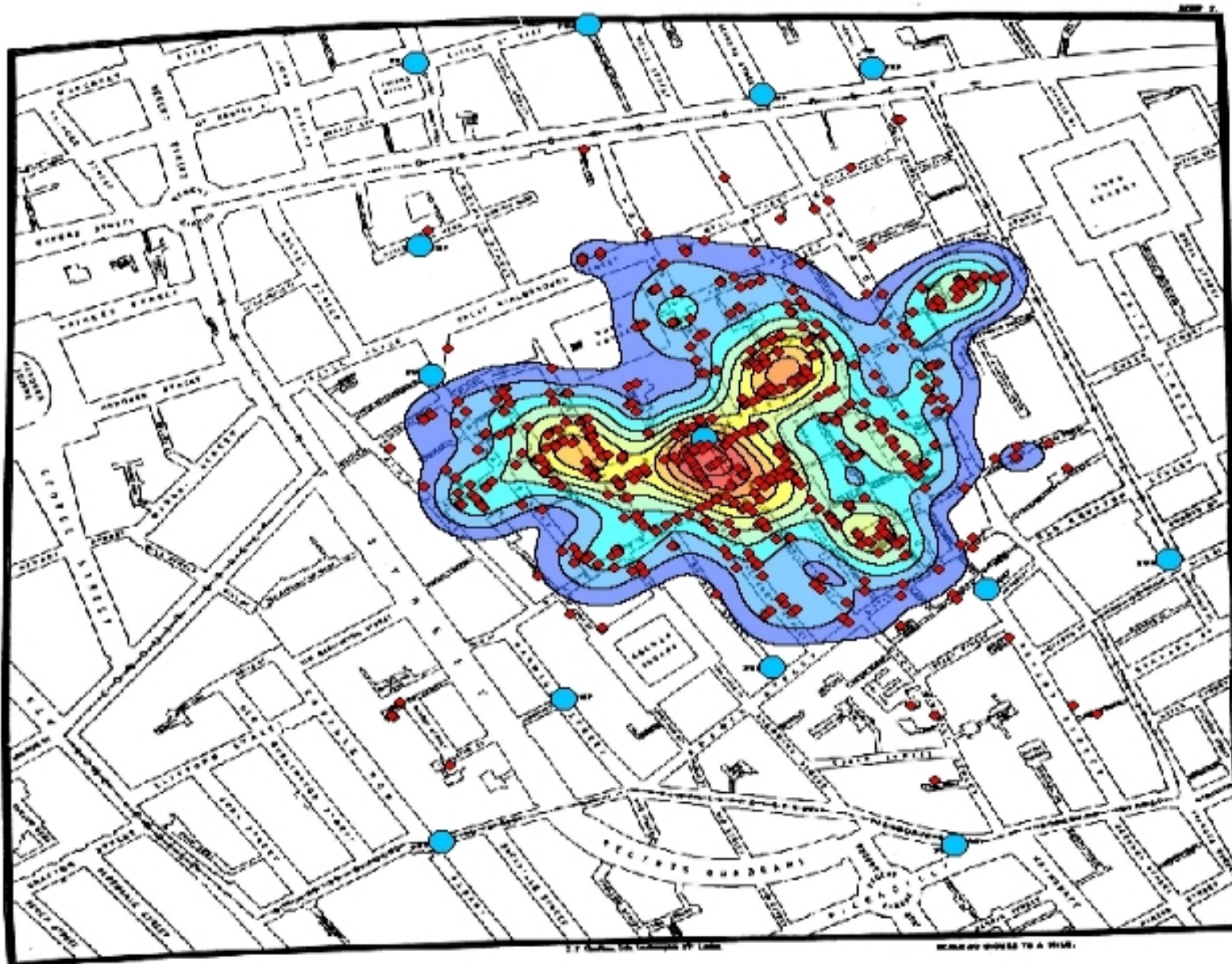
$$y = f(t, u)$$

$$t = g(r)$$

$r$ : randomized

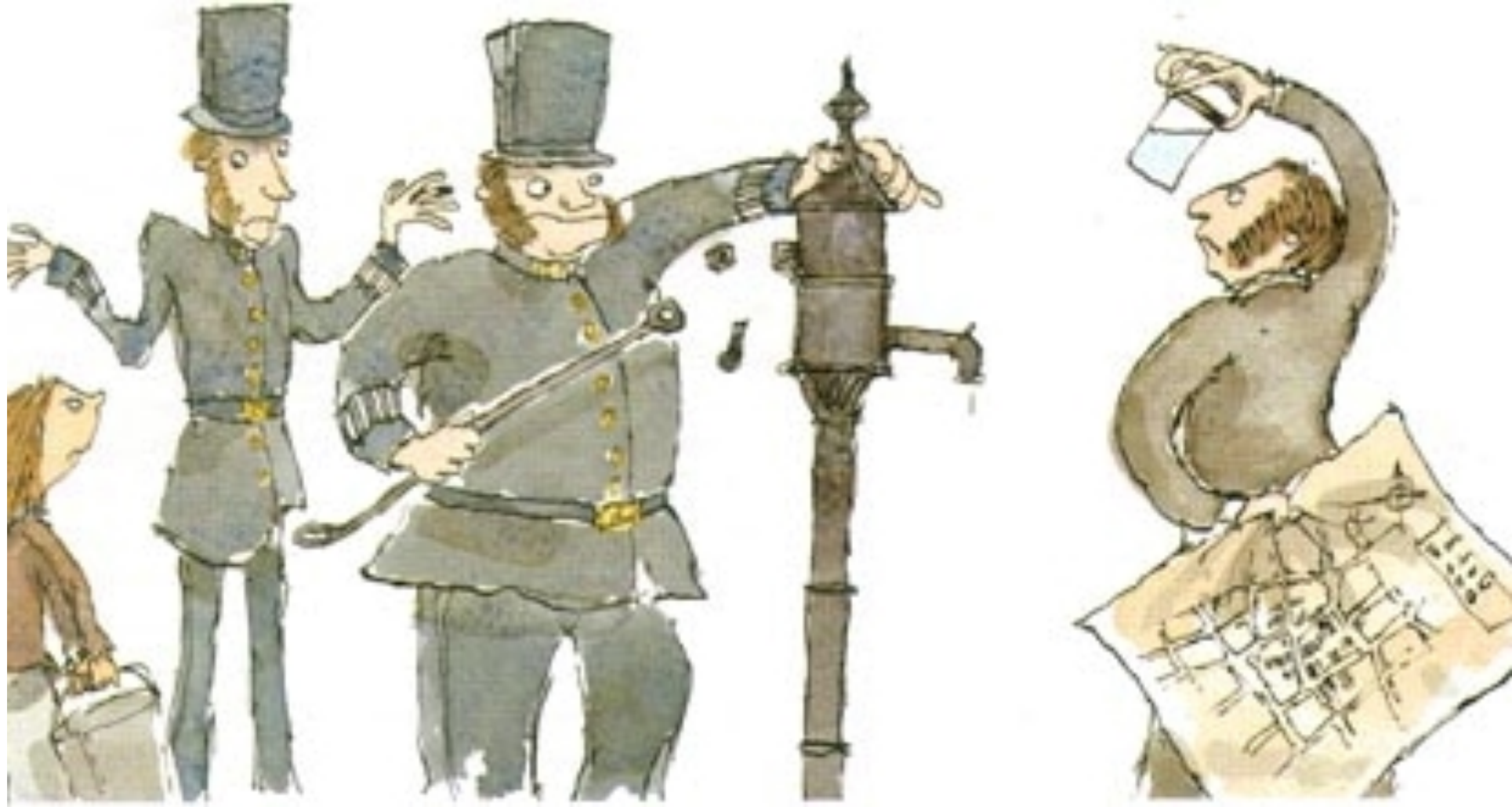
How to find such experiments?

**Example:** Cholera cause estimation in 1850s.

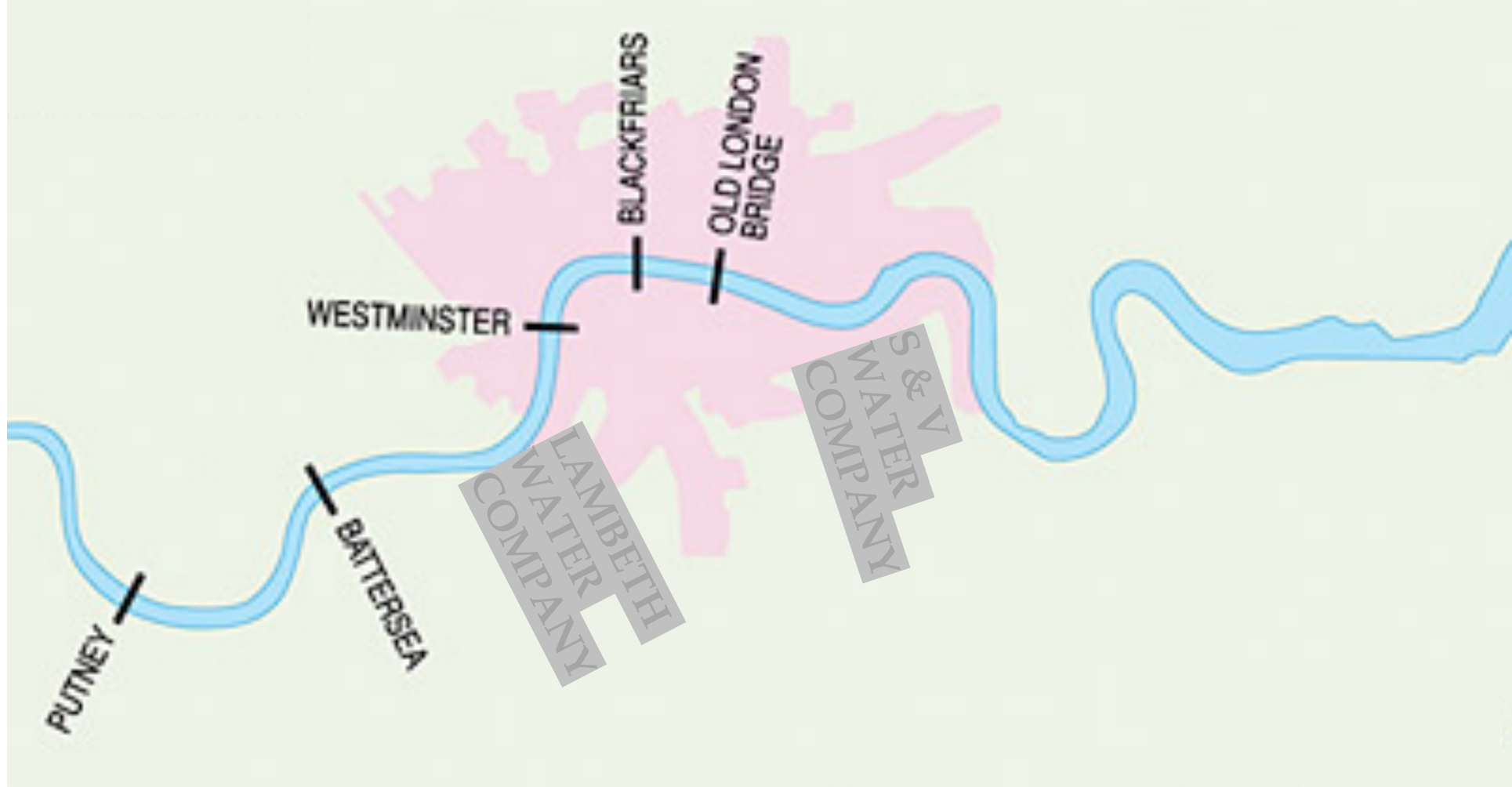


1854: London was having a devastating cholera outbreak

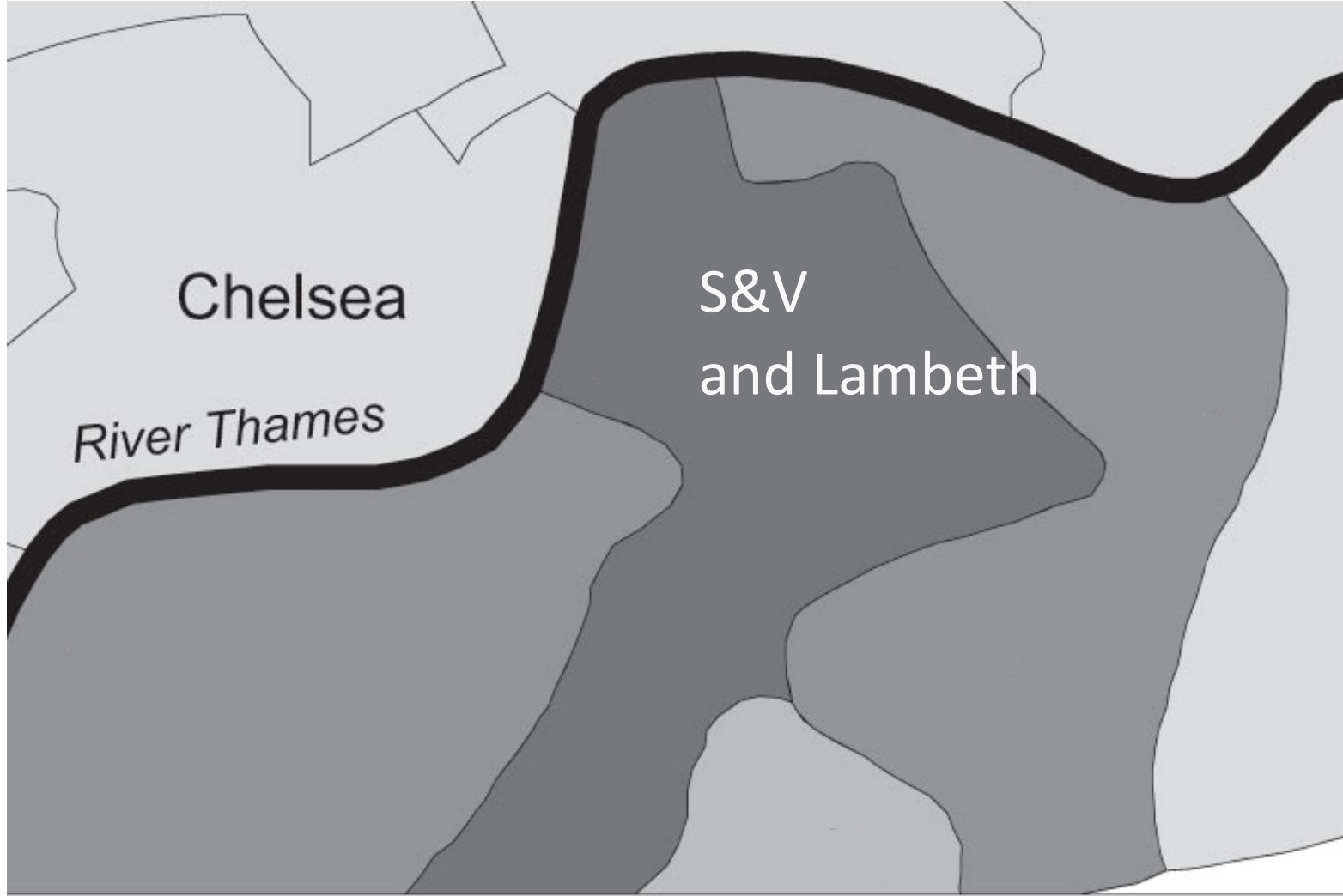




Enter John Snow. He found higher cholera deaths near a water pump, but could be just correlational.



**New Idea:** Two major water companies for London:  
one upstream and one downstream.  
Customers of each company distributed throughout city



No difference in neighborhood, still an 8-fold increase in cholera with the downstream company.

# “Natural” experiments: exploit variation in observed data

Can exploit naturally occurring **as-if random** variation in data.

Since data is not actively randomized, as-if-random remains an assumption.

Also need **exclusion**: the source of variation should not affect the outcome directly, only the treatment.

# What we just learned: Simple natural experiment

**Definition** Exploit “as-if random” assignment of treatments to measure outcome.

**Intuition** When assignment of treatment is unrelated to the measured outcome and their common causes, we can treat it as if it is a randomized experiment to estimate treatment effect.

**Example** What water company do you buy from?

**Keep in mind** As-if random assignments of treatments are hard to find. Estimates very sensitive to violation of exclusion assumption.

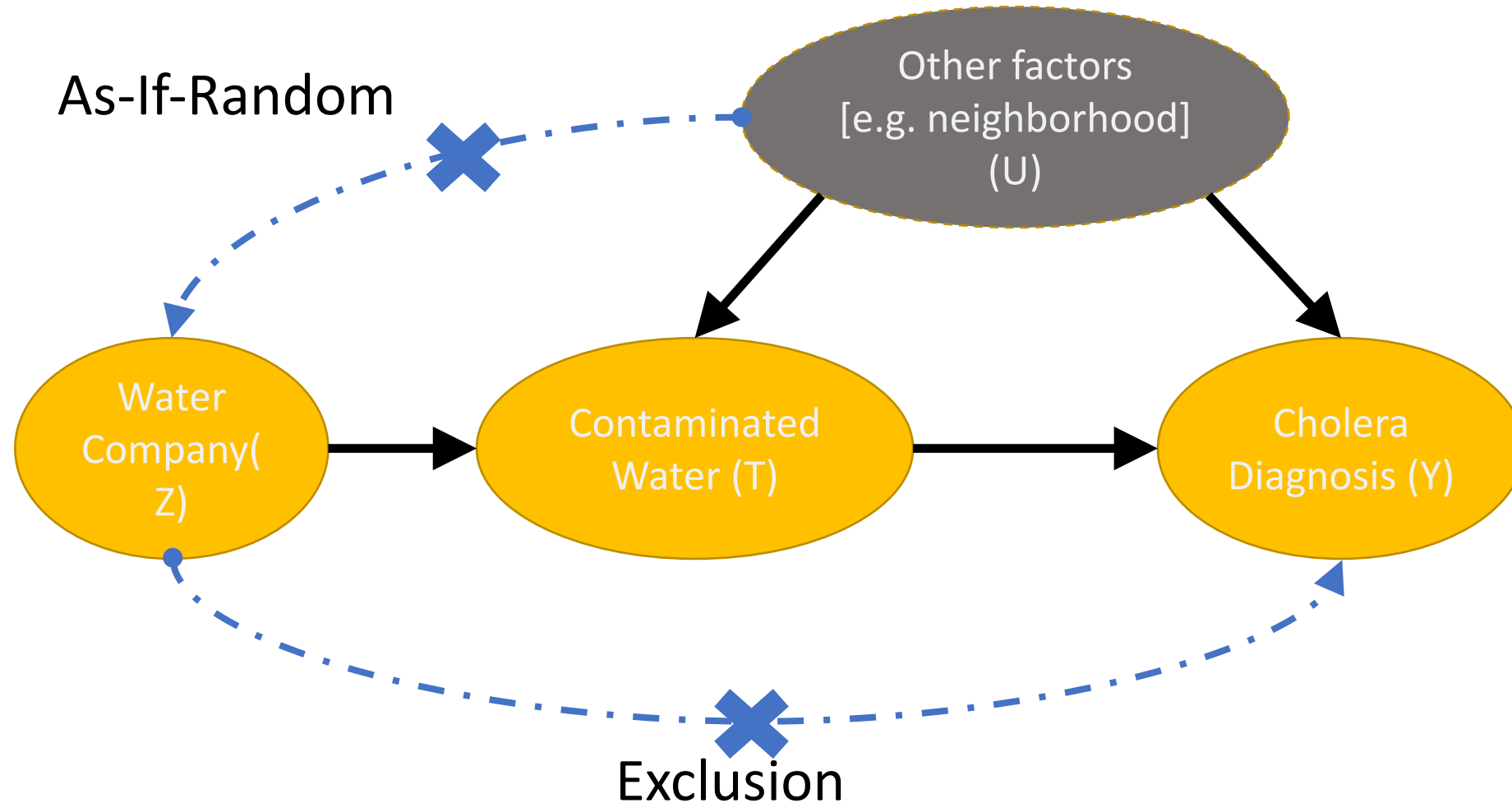
**Part II.B.  
Natural  
Experiments**

As-if Random

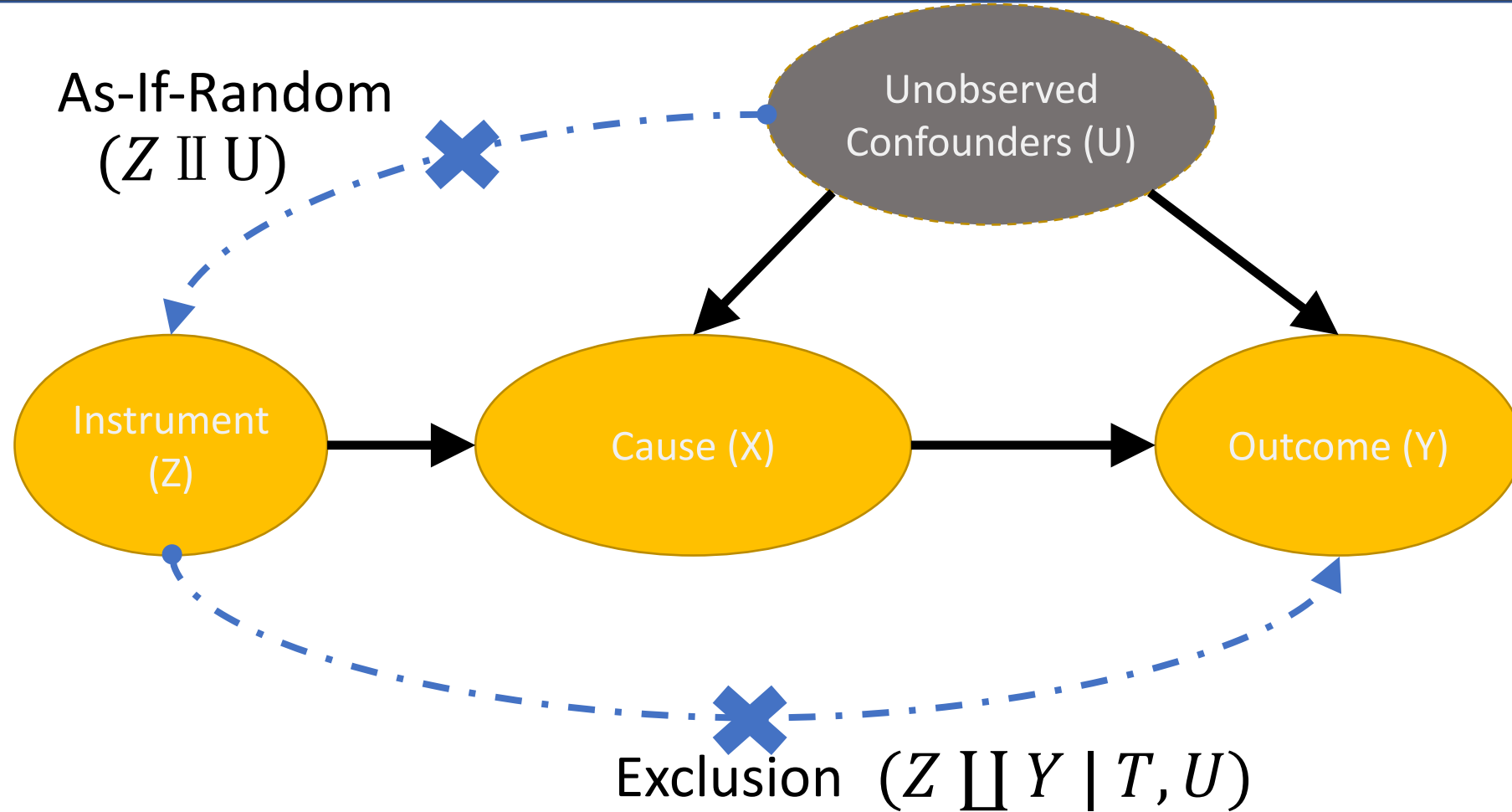
Instrumental Variables

Regression  
Discontinuities

Prior setup can be generalized as search for an “instrumental variable”



# Prior setup can be generalized as search for an “instrumental variable”





# Intuition: Can use this variation to compute causal effect

An increase in Z can lead to a change in Y *only through* X.

So change in Y is a product of change in Z->X and X->Y arrows.

Compare the extent by which random assignment affects X versus Y.

$$\text{Causal effect (X->Y)} = \frac{Y_{Z=1} - Y_{Z=0}}{X_{Z=1} - X_{Z=0}}$$

# A generalized natural experiment: Instrumental Variables

Can look at *as-if random* variations due to external events.

E.g.,

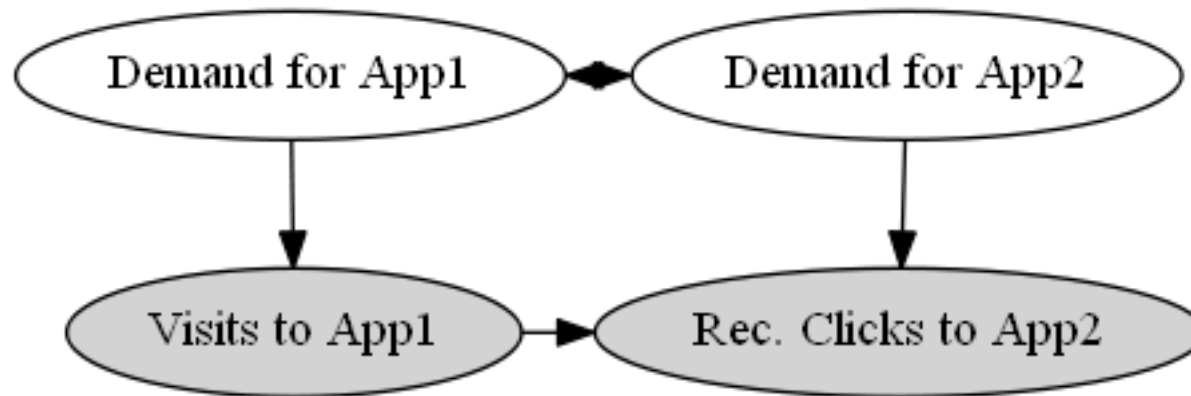
Experimental: Encouraging randomly selected users of an app to exercise.

Observational: Looking at a past A/B test intervention that increased chances of exercise.

*Example:* What is the effect of recommendations on an app store?

*Instrumental Variable:* External sources that drive sudden, large traffic to an app.

# Example: Effect of store recommendations

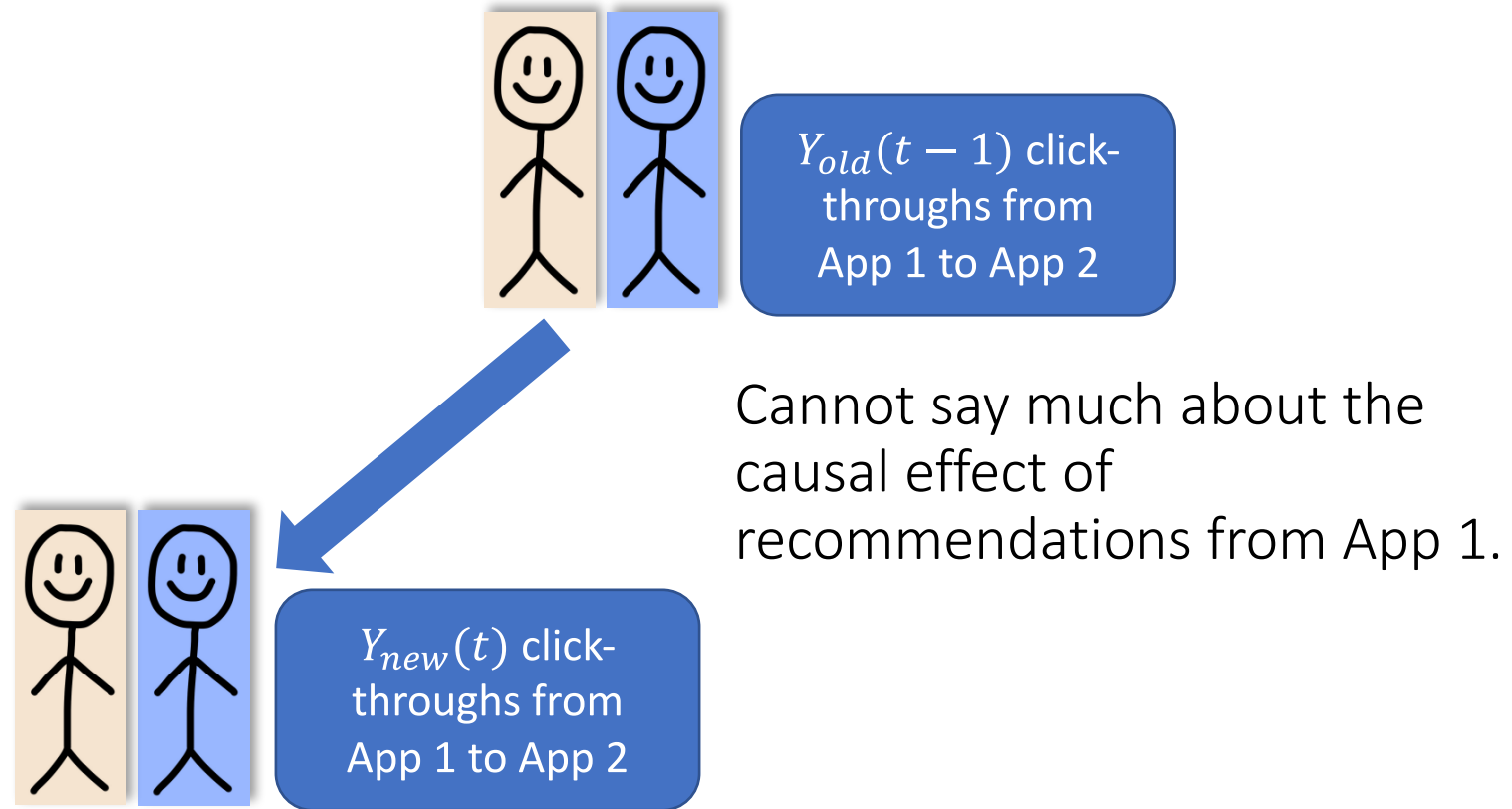


How many new visits are *caused* by the recommender system?

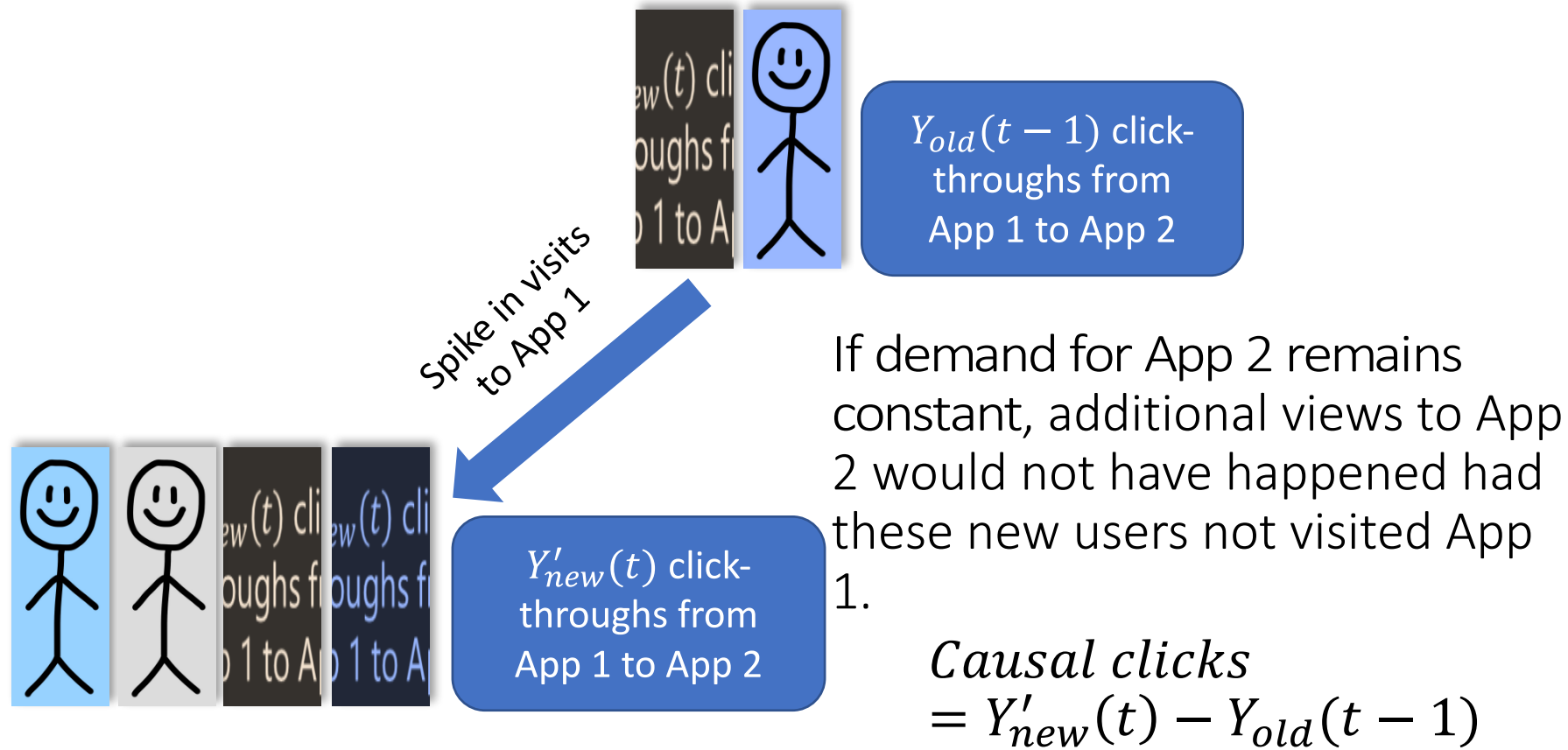
Demand for App 1 is correlated with demand for App 2.

⇒ Users would most likely have visited App 2 even without recommendations.

# Traffic on normal days to App 1



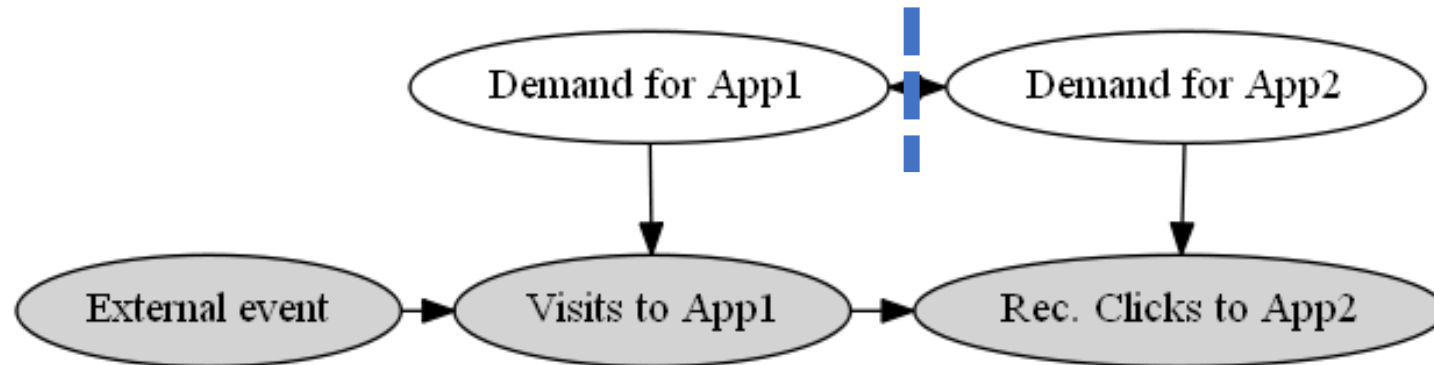
# External shock brings as-if random users to App1



# Exploiting sudden variation in traffic to App 1

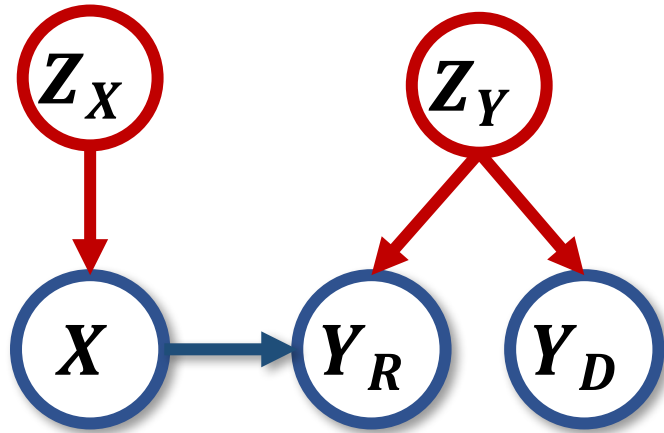
To compute Causal CTR of Visits to App1 on Visits to App2:

- Compare observed effect of external event separately on Visits to App1, and on Rec. Clicks to App2.
- Causal click-through rate =  $\frac{\Delta(\text{Rec. Click-throughs from App1 to App2})}{\Delta(\text{Visits to App1})}$



# Automatically Identifying Natural Experiments

Split-Door Criterion



- Finds 7,000 natural experiments, instead of 133
- Result: Across 10 product categories, half of recommendation clicks would have happened anyway



$\langle T, Y \rangle$

Examples of Instrumental  
Variables



Lottery



Weather



Shocks



Discontinuities

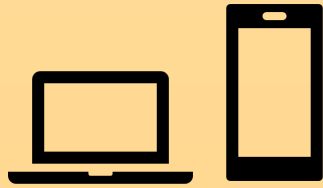


Hard-to-find  
variations



But there are so many natural variations.

$\langle T, Y \rangle$



Change in access of digital services



Change in train stops in a city



Change in medicines at a hospital

...



Lottery



Weather



Shocks



Discontinuities



Hard-to-find variations

# What we just learned: Instrumental Variables

**Definition** Instrumental variables (IV) introduce “as-if random” noise into treatment assignment, and are used to estimate treatment effect

**Intuition** Because IVs are not influenced by confounds, IVs’ indirect effect on outcome  $Y$  is independent of confounds too.  
Because IVs do not directly influence outcome, their effect must be due to the effect of the treatment.

**Examples** Encouraging people to exercise at random.  
Sudden increase in page visits to a product.

**Keep in Mind** Causal Estimate may not generalize to full population.  
Estimate very sensitive to the violations of IV assumptions.

Part II.B.  
Natural  
Experiments

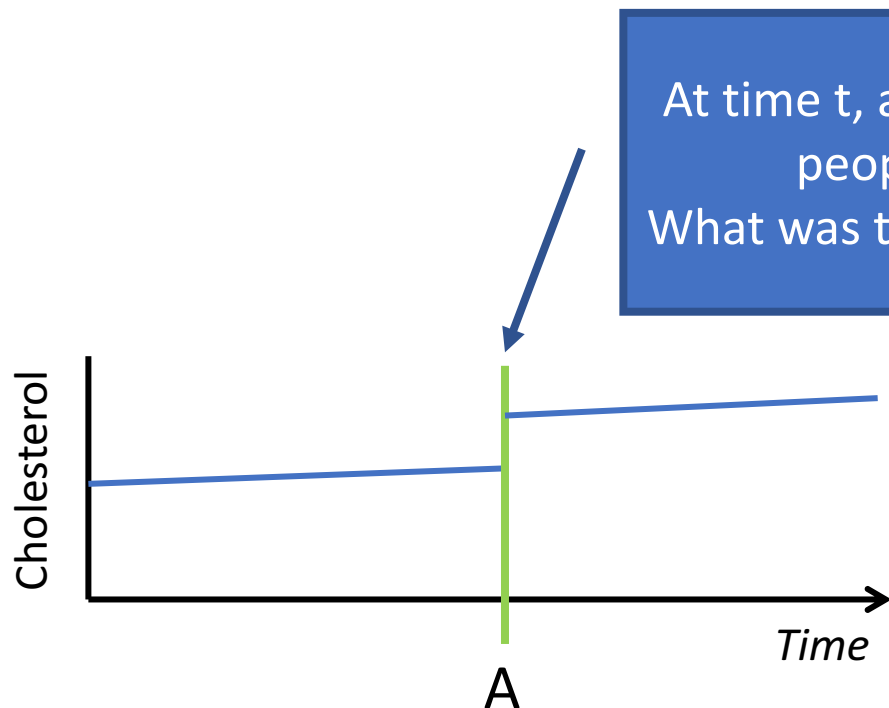
As-if Random

Instrumental Variables

Regression  
Discontinuities

# Regression discontinuities: Look for arbitrary changes to treatment

Instead of an IV changing the distribution of treatment over individuals, an arbitrary change decides the treatment deterministically.

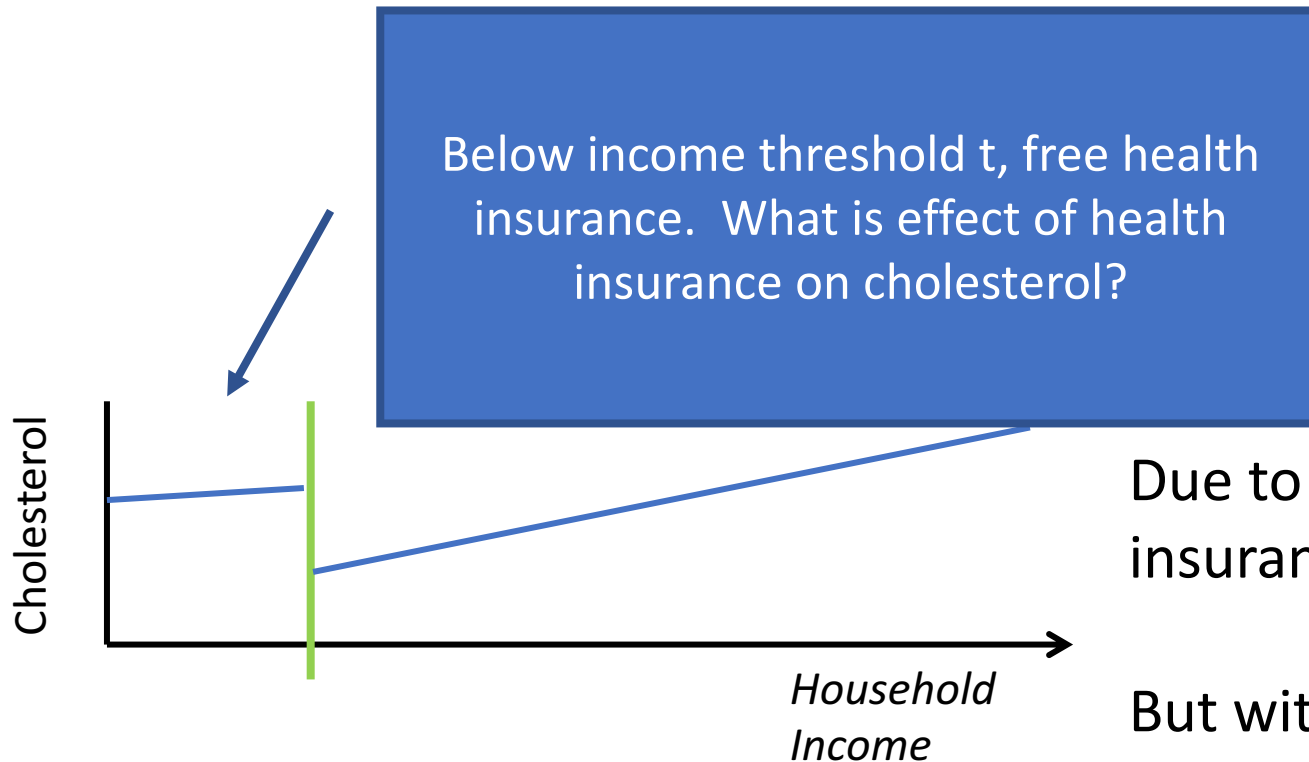


At time  $t$ , a cholesterol drug A is banned, and people switch to another drug B.  
What was the relative effect of drug A over B ?

Due to selection effects, people taking drug A are different from those taking drug B.

But within  $[t-1, t+1]$  duration, patients of A and B can be assumed to be similar.

# Regression discontinuities



Due to selection effects, people with health insurance different from those without.

But within  $[t-1, t+1]$  income, people with or without health insurance are similar.

# Regression discontinuities also depend on as-if-random and exclusion

**As-if-random:** People near the threshold are similar to each other, as if Nature randomized them on either side of the threshold.

**Exclusion:** Merely being on one side of the threshold does not affect the outcome.

**Very common:** Many decisions in organizations, arbitrary decisions in software are examples.

Can be thought of as a special case of an instrumental variable.

# Example: Effect of Store recommendations

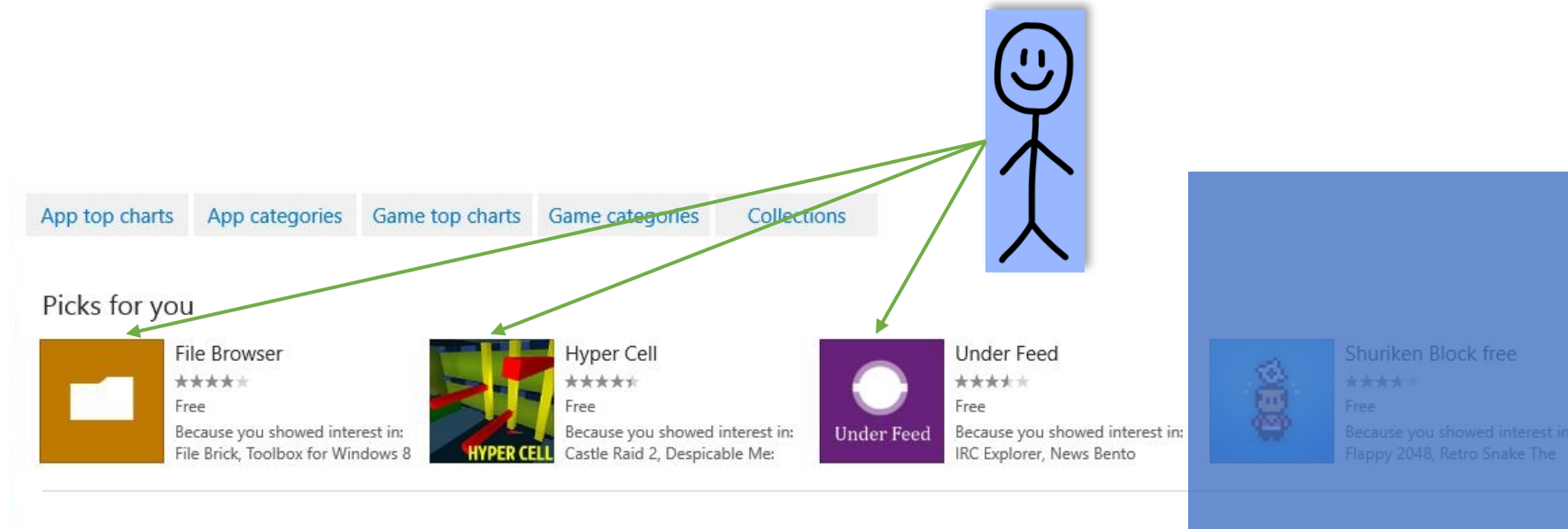
Suppose instead of comparing recommendation algorithms, we want to estimate the causal effect of showing *any* algorithmic recommendation.

Can be used to benchmark how much revenue a recommendation system brings, and allocate resources accordingly.

(and perhaps help analyze the tradeoff with users' privacy)

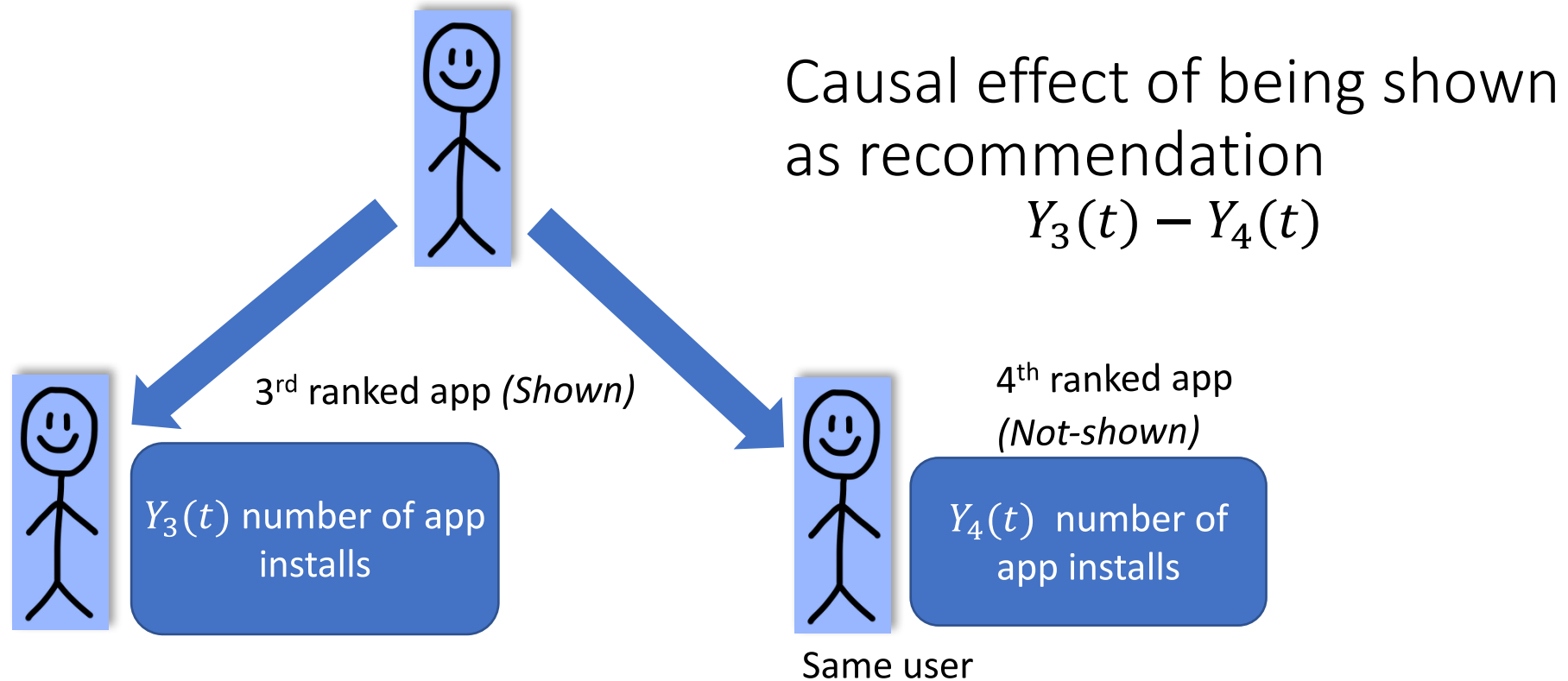


# Exploiting arbitrary cutoffs to recommendations



Only 3 recommendations shown to user.

# Assumption: Closely-ranked not-shown apps are as relevant as shown apps



# Algorithm: Regression discontinuity

For any top-k recommendation list:

- Using logs, identify apps that were similarly ranked but could not make it to the top-k shown apps.
- Measure difference in app installs between **shown and not-shown apps** for each user.

# What we just learned: Regression Discontinuities

**Definition** Regression discontinuities identify arbitrary boundaries between treated and untreated populations, measure treatment effect as difference in outcomes at the boundary

**Intuition** Regression discontinuities approximate randomized experiments as long as no substantial differences between people just on one side or the other. That is, at the boundary,  $T \perp X, U$

**Example** Policy decisions based on income or time; exogenous shocks; and are all common sources of regression discontinuities

**Keep in mind** Only estimates treatment effect at the boundary. Effect may vary elsewhere!

**PART II.  
Methods  
for Causal  
Inference**

Observational Studies

Natural Experiments

Refutations

# Causal inference is only possible with assumptions

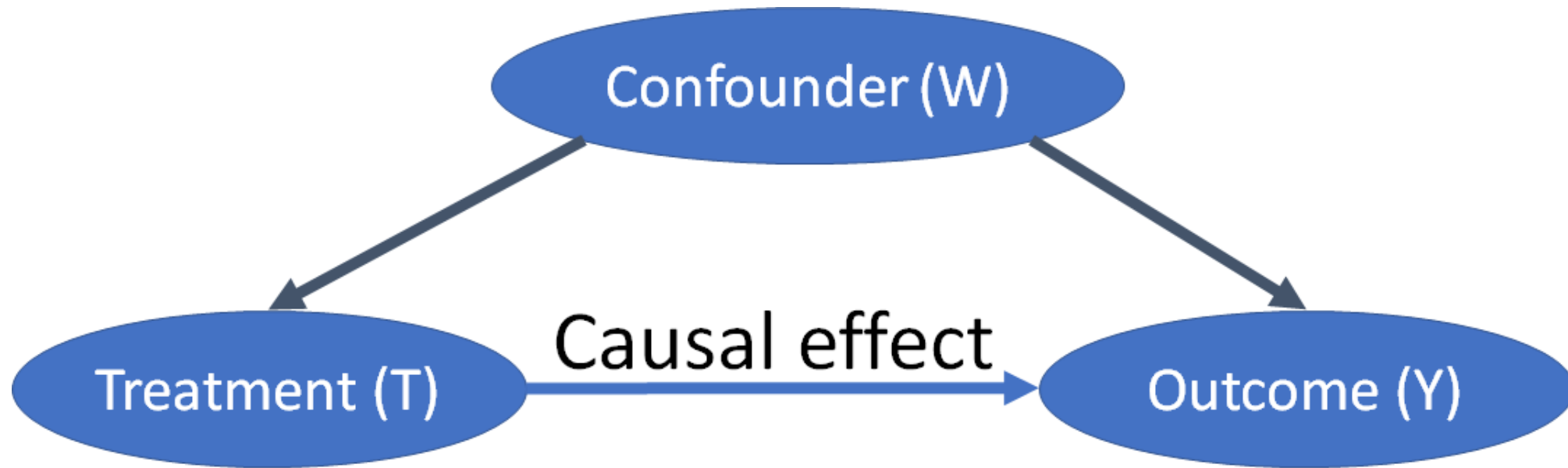
“Causal” part does not come from the data.

It comes from your assumptions that lead to ***identification***.

The data is simply used for statistical ***estimation***.

Critical to verify your assumptions. But how?

# (Step 1): Making explicit the difference between identification and estimation



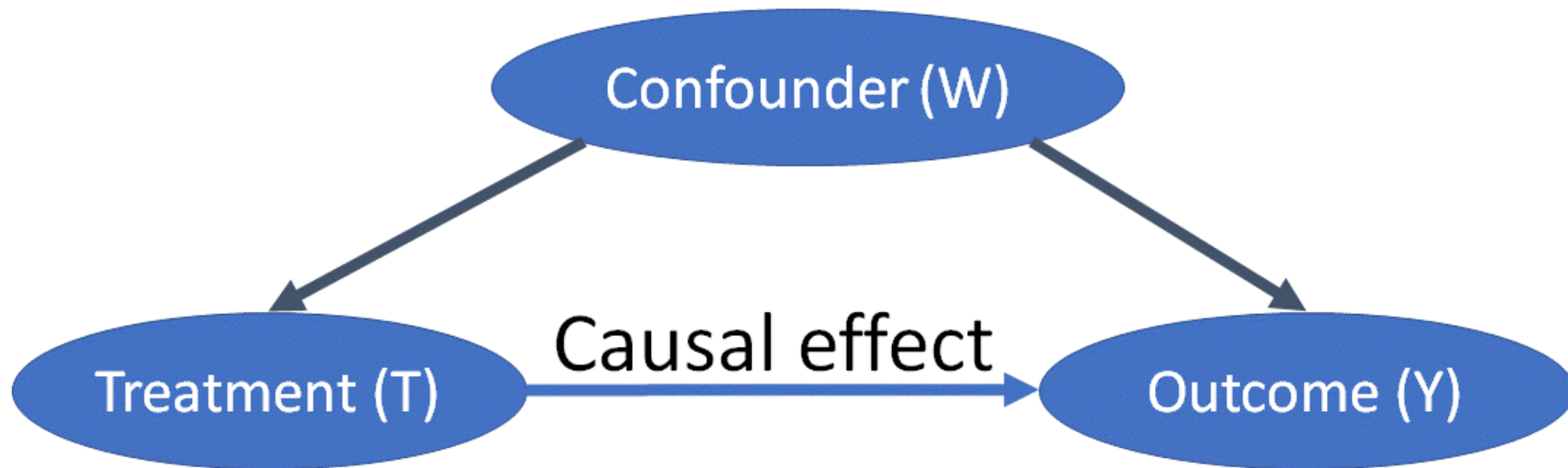
**Identification:** Causal effect  $\rightarrow$  Observed effect conditioned on  $W$ ,  $E[Y|T, W]$

**Estimation:**  $E[Y|T, W] \rightarrow$  Propensity Score Stratification

**Why do observational studies fail?** Most likely due to errors in identification.

--Estimation is a statistical problem, relatively easy.

(Step 2): Explicitly represent your identifying and estimating assumptions.

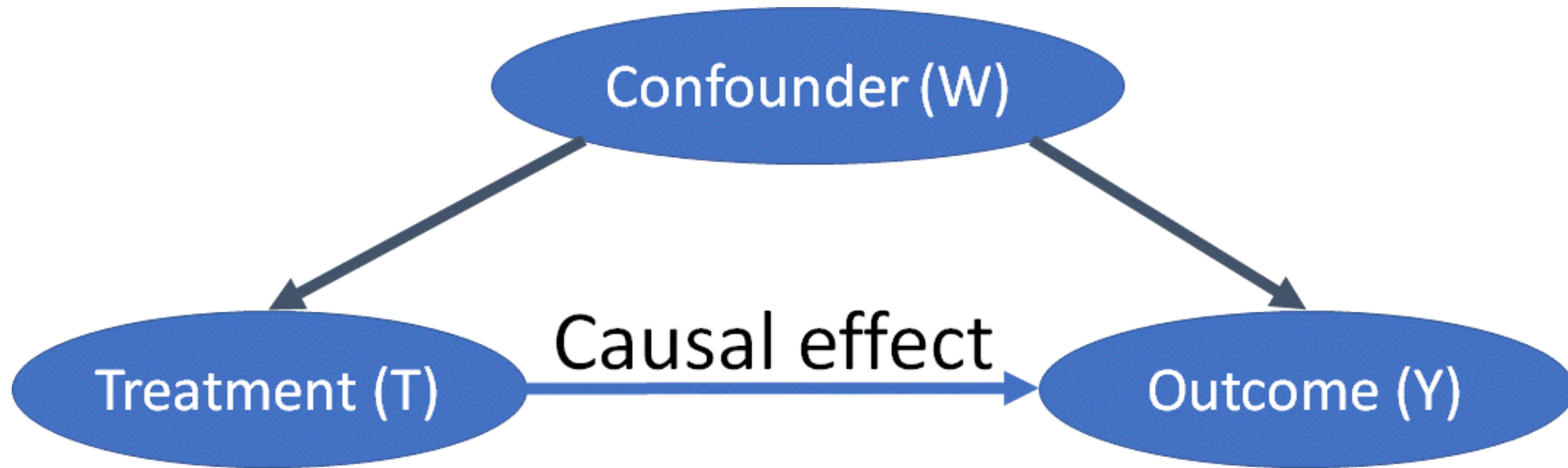


**Identifying assumption:** All the arrows missing in the causal graphical model. E.g. No other common cause exists -> Untestable in general.

**Estimating assumption:** Overlap between treated and untreated population. Can be solved by collecting more data.



(Step 3): Refute your assumptions, and analyze your estimate's sensitivity to violations



**Identifying assumption:** All the arrows missing in the causal graphical model. E.g. No other common cause exists -> Untestable in general.

- *What happens* when another common cause exists?
- *What happens* when treatment is placebo?

# To make these steps easy, we created DoWhy: a python library for causal inference

DoWhy focuses attention on the **assumptions** required for causal inference.

Provides estimation methods such as matching and IV so that you can focus on the identifying assumptions.

- Models assumptions explicitly using causal graphical model.
- Provides an easy way to test them (if possible) or analyze sensitivity to violations.

Unifies all methods to yield **four verbs** for causal inference:

- Model
- Identify
- Estimate
- Refute

# DoWhy: Sample causal inference analysis in 4 lines

```
from dowhy.do_why import CausalModel

# Create a causal model from the data and given graph.
model=CausalModel(
    data = df,
    treatment=data["treatment_name"],
    outcome=data["outcome_name"],
    graph=data["dot_graph"],
)

# Identify causal effect and return target estimands
identified_estimand = model.identify_effect()

# Estimate the target estimand using a statistical method.
estimate = model.estimate_effect(identified_estimand,
    method_name="backdoor.propensity_score_matching")

# Refute the obtained estimate using multiple robustness checks.
refute_results=model.refute_estimate(identified_estimand, estimate,
    method_names=["random_common_cause", "placebo_treatment_refuter",
        "data_subset_refuter"])
```

# Refutation 1: Add random variables to your model

Can add randomly drawn covariates into data

Rerun your analysis.

Does the causal estimate change? (*Hint: it shouldn't*)

# Refutation check 2: Replace treatment by a placebo (A/A test)

Randomize or permute the treatment.

Rerun your analysis.

Does the causal estimate change? (*Hint: it should become 0*)

# Refutation Check 3: Divide data into subsets (cross-validation)

Create subsets of your data.

Rerun your analysis.

Does the causal estimate vary across subsets?  
*(Hint: it shouldn't vary significantly)*

# Refutation Check 4: Test Balance of Covariates

Many methods (e.g., matching, stratification, weighting, regression discontinuity) depend on balancing of covariates

Can test this.

# When refutations are not possible? Sensitivity Analysis to violations of assumptions

**Question:** *How sensitive is your estimate to minor violations of assumptions?*

*E.g. How big should the effect of a confounder be so that your estimate reverses in direction?*

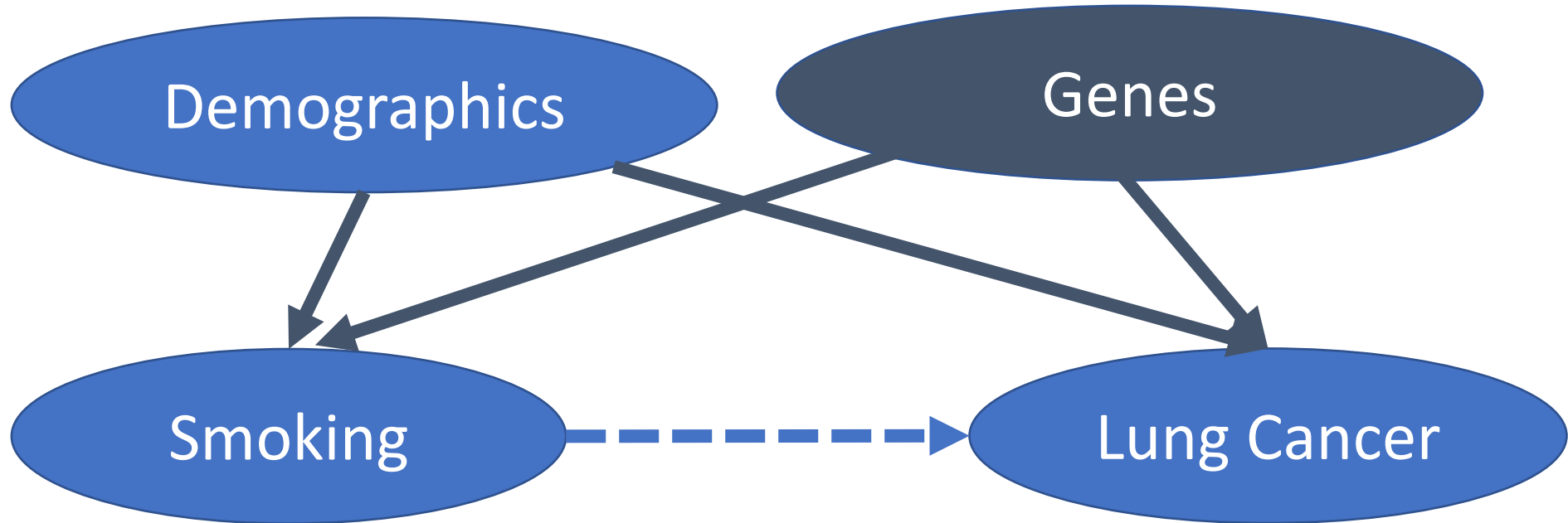
Use simulation to add effect of unknown confounders.

Domain knowledge helps to guide reasonable values of the simulation.

Make comparisons to other known estimates.



# Example: Does smoking cause lung cancer?



Cornwell (1959) showed that the effect of Genes had to be 8 times any known confounder for the effect to go to zero.

# Observational causal inference: Best practices

Always follow the four steps: *Model, Identify, Estimate, Refute*.

Refute is the most important step.

Aim for simplicity.

If your analysis is too complicated, it is most likely wrong.

Try at least two methods with different assumptions.

Higher confidence in estimate if both methods agree.

# Try out DoWhy to see best practices in action

## DoWhy: A Python Library for Causal Inference

**Principled:** Converts prior knowledge to a formal causal graph

**Simple:** Automated analysis of many assumptions, one line of code for powerful causal inference algorithms

**Robust:** Battery of tests to refute obtained estimates

**Modest:** No estimate if the data is insufficient

- **Input:** Observational data, Causal graph
- **Output:** Causal effect between desired variables, “What-if” analysis

Code: <https://github.com/Microsoft/dowhy>

Docs: <http://causalinference.gitlab.io/dowhy>

PART I. Introduction to Counterfactual Reasoning

PART II. Methods for Causal Inference

**PART III. Large-scale and Network Data**

**PART IV. Broader Landscape**

PART I. Introduction to Counterfactual Reasoning

PART II. Methods for Causal Inference

PART III. Large-scale and Network Data

**PART IV. Broader Landscape**

## PART IV.

High-level awareness  
of broader  
landscape in causal  
reasoning

# Outline

- Discovery of causal relationships from data
- Heterogeneous treatment effects
- Machine learning, representations and causal inference
- Reinforcement learning and causal inference
- “Automated” causal inference

Causal discovery



# Effects of causes and causes of effects

- We discussed causal inference: effects of causes
- But a complementary question is causal discovery
  - [Local] Causes of effects
  - [Global] Mapping out causal mechanisms
- In general, a harder problem.
- See Causation [Spirtes (2000)] and Elements of Causal Inference (Scholkopf et al. 2017).

Heterogenous treatment effects

# Average causal effect does not capture individual-level variations

- Stratification is one of the simplest methods for heterogeneous treatment by strata
- Typical strata are demographics.
- Need more data to statistically detect differences
  
- For high-dimensions, can use machine learning methods like random forests [Athey and Wager, 2015]

# Machine learning and causal inference

# Causal inference as a (counterfactual) prediction problem

**Causal inference  $\Leftrightarrow$  robust prediction**

(Supervised) ML

Predicted value under the training distribution

$P(X, y)$ .

$P(X, y): y = k(X) + \epsilon$

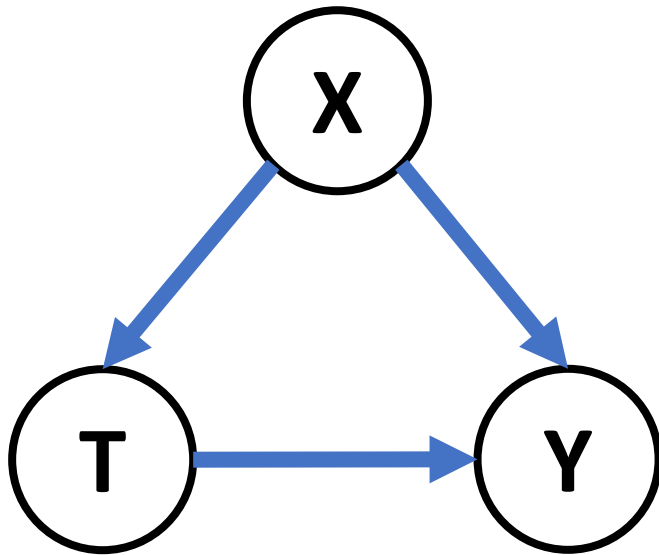
Causal inference

Predicted value under the counterfactual distribution

$P'(X, y)$ .

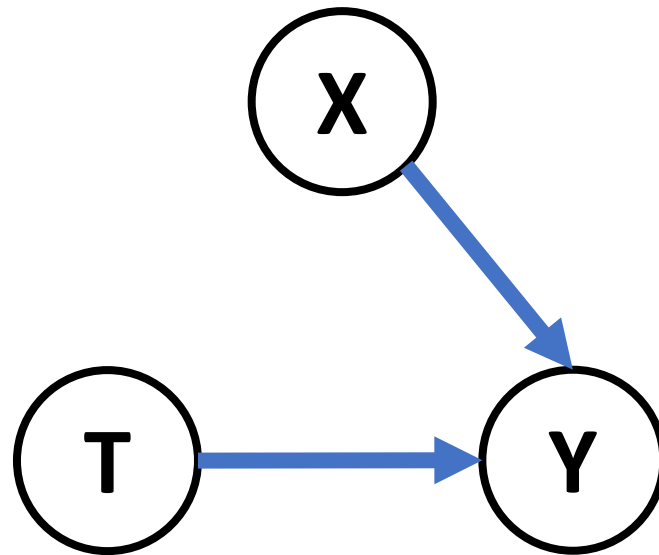
$P'(X, y): y = ?$

# Causal inference: A special kind of domain adaptation



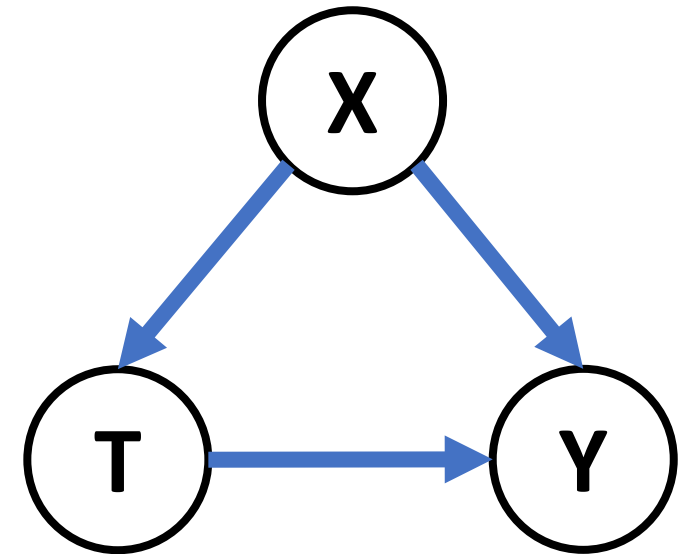
$P(Y,T,X)$

Observed data



$P^*(Y,T,X)$

Randomized  
experiment



$P^{**}(Y,T,X)$

Another domain

# Predicting the counterfactual $\Leftrightarrow$ Causal Inference

Predicting Individual treatment effects can be considered as domain adaptation  
--Use regularization and transformation of input features [Johansson 2016]

Generalizing prediction to new domains

-- Selection bias or covariate shift [Barenboim and Pearl 2013]

-- If predictive model generalizes to new domains, can be considered “causal”  
[Peters et al. 2015]

# Causal inference and machine learning

## **Machine learning**

Use causal inference methods for robust, generalizable prediction.

## Causal inference

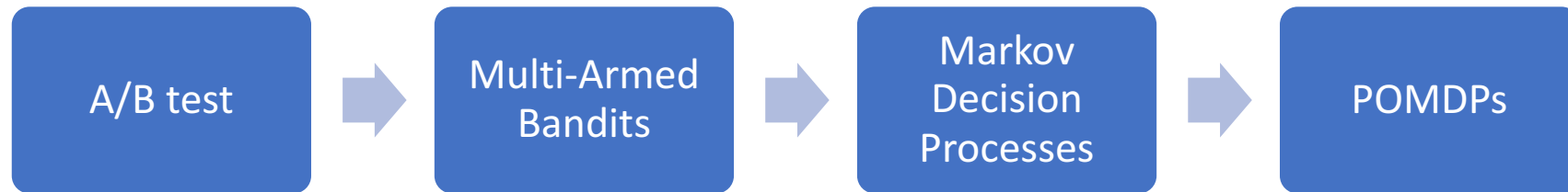
Use ML algorithms to better model the non-linear effect of confounders, or find low-dimensional representations.

In general, be wary of methods that have not been empirically tested, especially ones that you do not understand.



# Reinforcement learning and causal inference

# Generalizing a randomized experiment

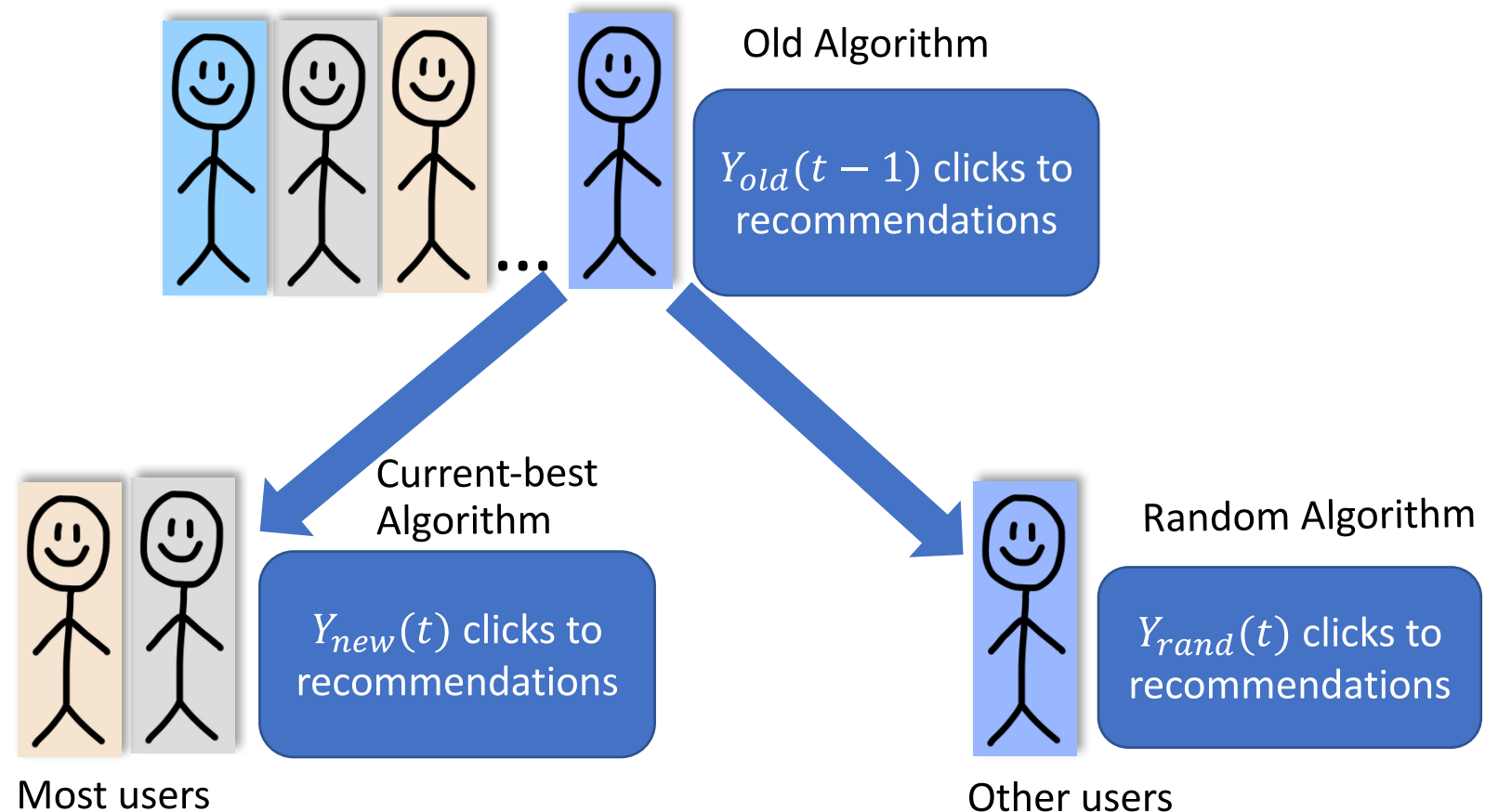


# Efficient randomized experiment: Multi-armed bandits

## Two goals:

1. Show the best known algorithm to most users.
2. Keep randomizing to update knowledge about competing algorithms.

**“Explore and Exploit”  
strategy**



# Algorithm: $\epsilon$ -greedy multi-armed bandits

Repeat:

**(Explore)** With low probability  $\epsilon$ , choose an output item randomly.

**(Exploit)** Otherwise, show the current-best algorithm.

Use CTR results for Random output items to train new algorithms offline.

# Practical Example: Contextual bandits on Yahoo! News

**Actions:** Different news articles to display

A/B tests using all articles inefficient.

Randomize the articles shown using  $\epsilon$ -greedy policy.

Better: Use context of visit (user, browser, time, etc.) to have different current-best algorithms for different contexts.



The screenshot shows a 'Featured' news section with navigation tabs for 'Entertainment', 'Sports', and 'Life'. The main article is titled 'McNair's final hours revealed' with a large 'STORY' overlay. Below it are four smaller article thumbnails labeled F1, F2, F3, and F4. At the bottom right, there is a link to 'More: Featured | Buzz'.

**Featured** | Entertainment | Sports | Life

**McNair's final hours revealed**  
STORY  
Police release 50 text messages that depict the late NFL player's alleged killer as losing control. » **Details**

- UConn murder victim mourned

Find Steve McNair murder case

**F1** McNair's final hours revealed

**F2** Cindy Crawford stays fierce in black mini

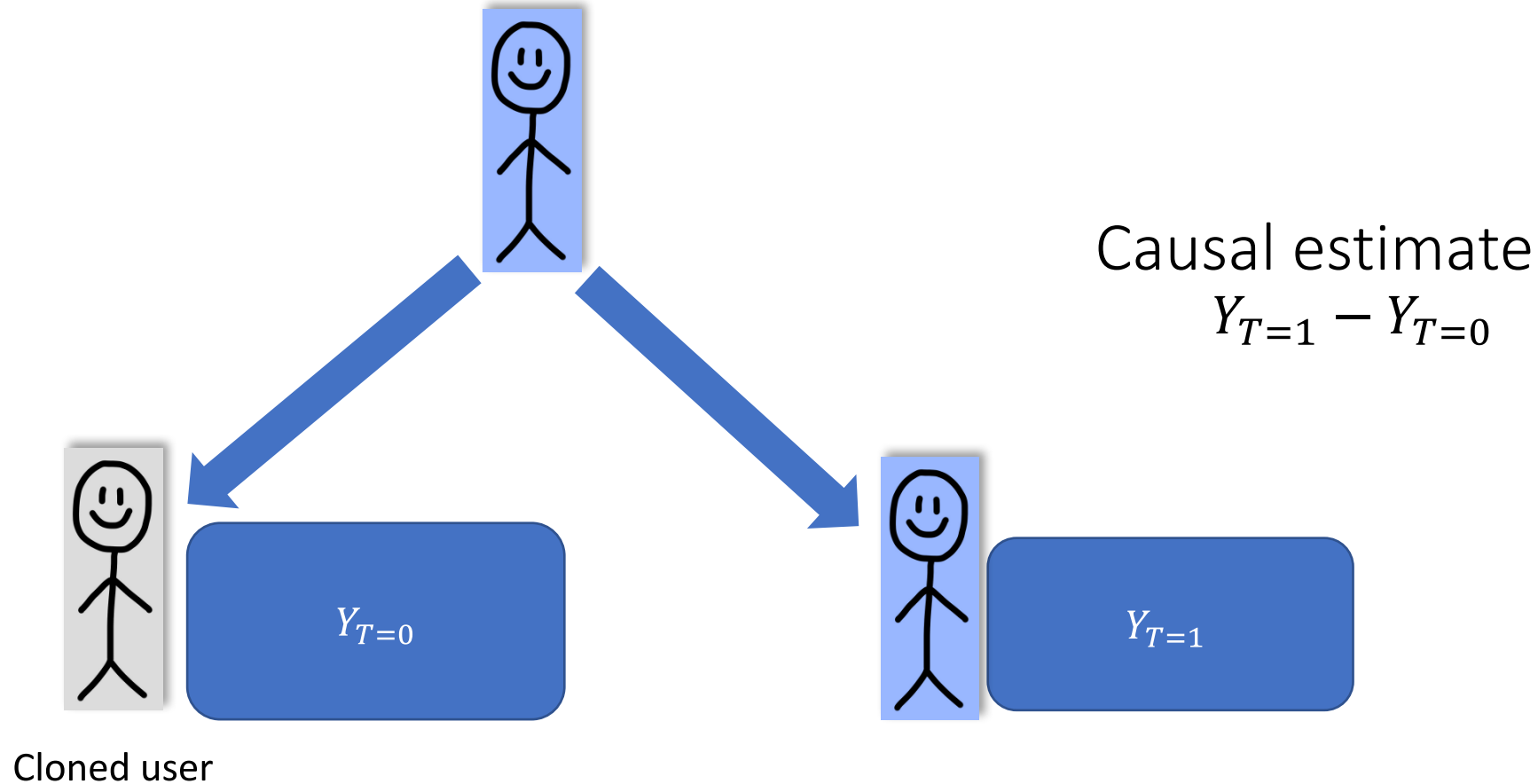
**F3** Watch for dozens of 'shooting stars' tonight

**F4** At team's big moment, star player isn't around

» More: **Featured** | **Buzz**

Many of these techniques can be combined

Remember, we are always looking for the ideal experiment with multiple worlds



# Example: Randomization + Instrumental Variable

**Treatment example:** You cannot randomize who exercises, but maybe can provide incentives to join the gym.

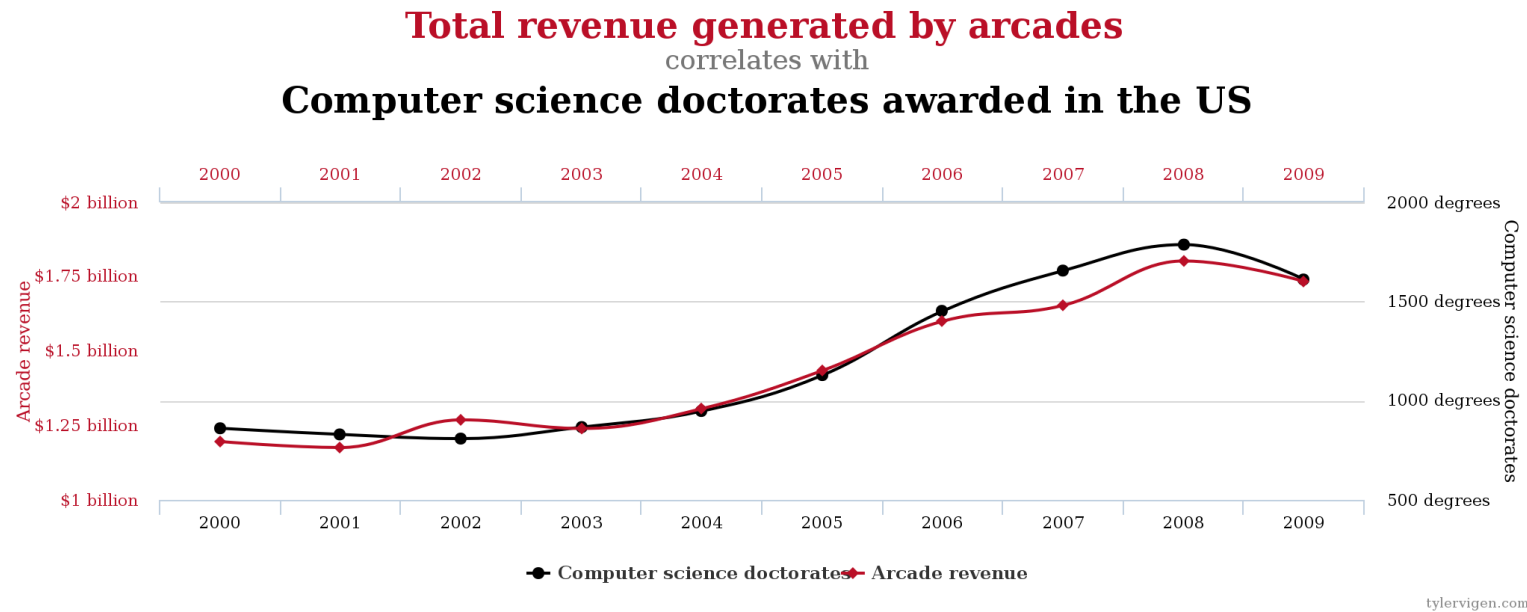
**Algorithm example:** You cannot remove recommendations at random, but could advertise a focal product to a random subset of people on the homepage.



# Conclusions

# Causal inference is tricky

Correlations are seldom enough. And sometimes horribly misleading.



Always be skeptical of causal claims from observational any data.  
More data does not automatically lead to better causal estimates.

# Causal inference: Best practices

**Always follow the four steps: *Model, Identify, Estimate, Refute.***

--Refute is the most important step.

**Aim for simplicity.**

--If your analysis is too complicated, it is most likely wrong.

**Try at least two methods with different assumptions.**

--Higher confidence in estimate if both methods agree.

# Thank you!

Emre Kiciman, Amit Sharma (Microsoft)

@emrek, @amt\_shrma

Tutorial and other resources will be posted at:

<http://causalinference.gitlab.io>

DoWhy library can be accessed at

<http://causalinference.gitlab.io/dowhy>