

# Modeling Social Data: Causality and Experiments

Guest lecturers: Andrew Mao, Amit Sharma

# Intro to Causal Inference

Amit #1

# Prediction

Make a forecast, leaving the world as it is  
(seeing my neighbor with an umbrella might predict rain)

vs.

# Causation

Anticipate what will happen when you make a change in the world  
(but handing my neighbor an umbrella doesn't cause rain)

# “Causes of effects”

---

It's tempting to ask “what caused Y”, e.g.

- What makes an email spam?
- What caused my kid to get sick?
- Why did the stock market drop?

This is “reverse causal inference”, and is generally quite hard

# “Effects of causes”

---

Alternatively, we can ask “what happens if we do X?”, e.g.

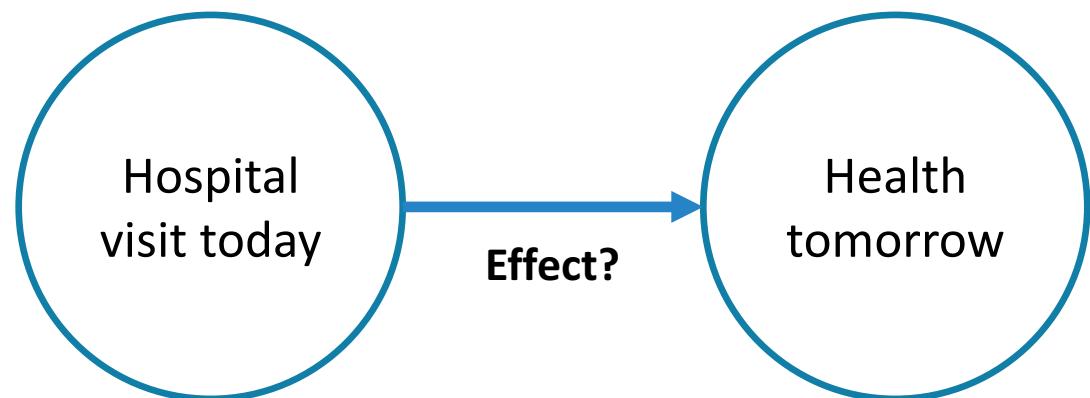
- How does education impact future earnings?
- What is the effect of advertising on sales?
- How does hospitalization affect health?

This is “forward causal inference”: still hard, but less contentious!

# Example: Hospitalization on health

---

What's wrong with estimating this model from observational data?

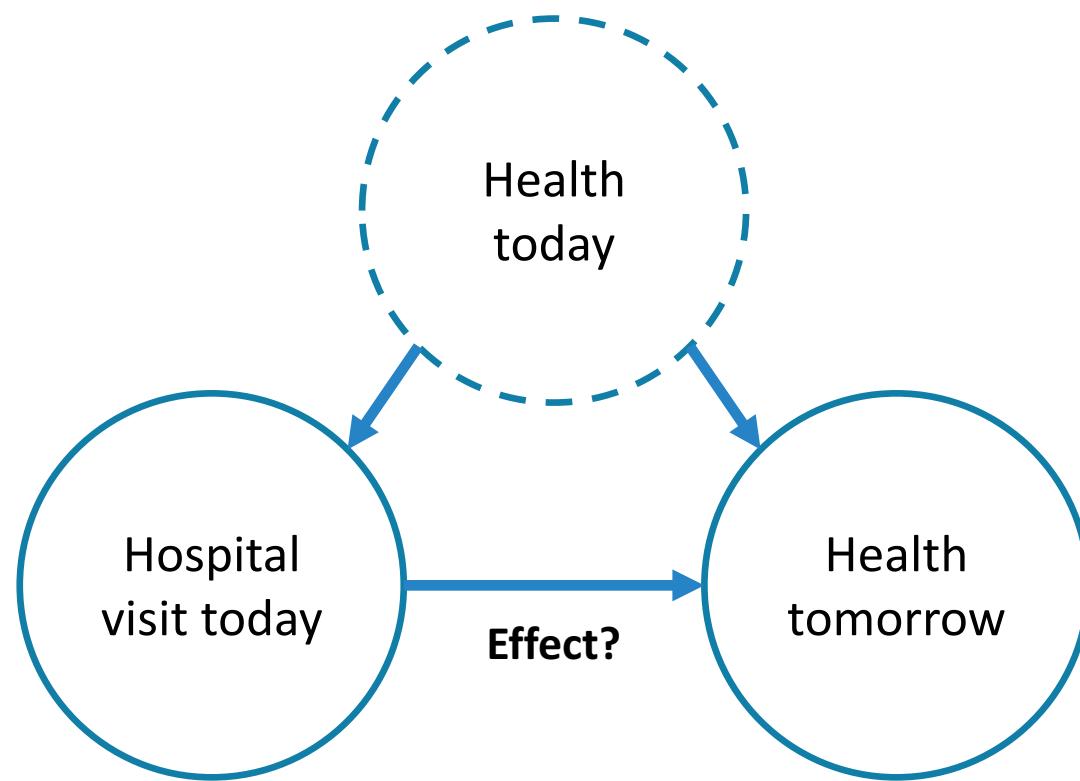


Arrow means “X causes Y”

# Confounders

---

The effect and cause might be *confounded* by a common cause, and be *changing together* as a result

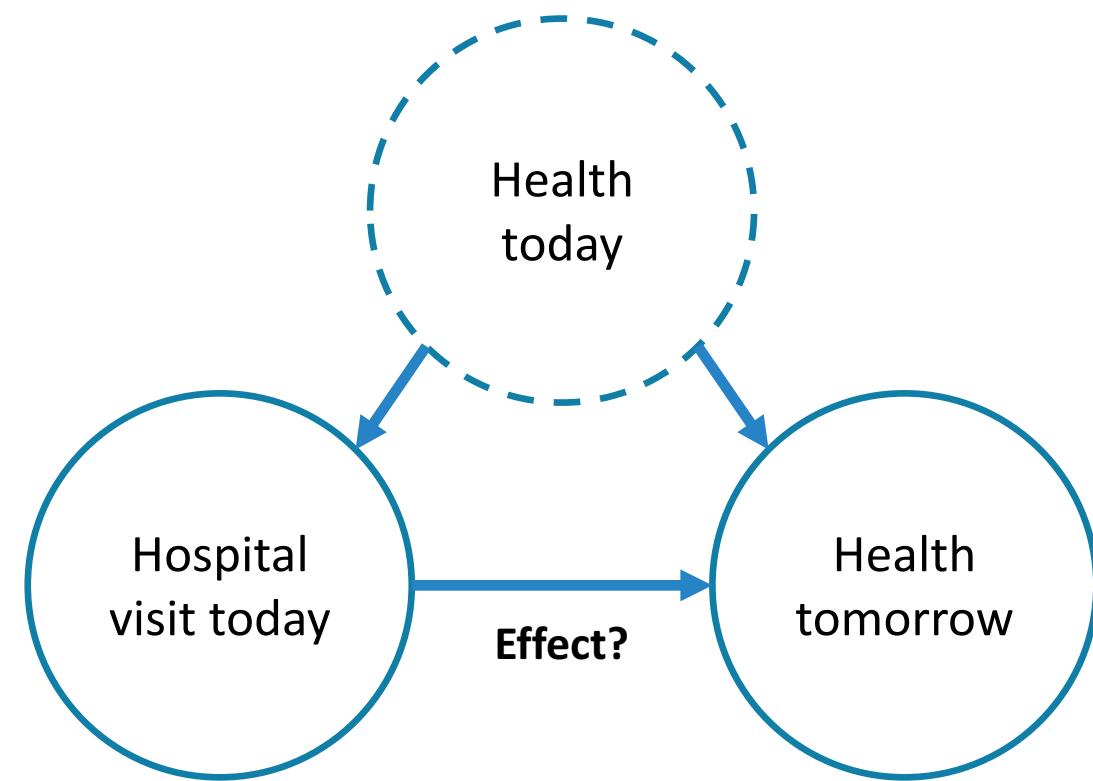


Dashed circle means “unobserved”

# Confounders

---

If we *only get to observe them changing together*, we can't estimate the effect of hospitalization changing alone



# Observational estimates

---

Let's say all sick people in our dataset went to the hospital today, and healthy people stayed home

The observed difference in health tomorrow is:

$$\Delta_{\text{obs}} = (\text{Sick and went to hospital}) - (\text{Healthy and stayed home})$$

# Observational estimates

---

Let's say all sick people in our dataset went to the hospital today, and healthy people stayed home

The observed difference in health tomorrow is:

$$\Delta_{\text{obs}} = [(\text{Sick and went to hospital}) - (\text{Sick if stayed home})] + \\ [(\text{Sick if stayed home}) - (\text{Healthy and stayed home})]$$

# Selection bias

---

Let's say all sick people in our dataset went to the hospital today, and healthy people stayed home

The observed difference in health tomorrow is:

$$\Delta_{\text{obs}} = \underbrace{[(\text{Sick and went to hospital}) - (\text{Sick if stayed home})]}_{\text{Causal effect}} + \underbrace{[(\text{Sick if stayed home}) - (\text{Healthy and stayed home})]}_{\text{Selection bias}}$$

(Baseline difference between those who opted in to the treatment and those who didn't)

# Basic identity of causal inference<sup>1</sup>

---

Let's say all sick people in our dataset went to the hospital today, and healthy people stayed home

The observed difference in health tomorrow is:

$$\text{Observed difference} = \text{Causal effect} - \text{Selection bias}$$

Selection bias is likely negative here, making the observed difference an underestimate of the causal effect

<sup>1</sup>Varian (2016)

# A fundamental problem across science and industry

- Knowledge (Science)
  - What is the effect of doing X?
  - Social sciences
  - Biology and medicine
- Solving problems (Industry)
  - When should I do X?
  - Predictive models in practice
  - Policies for better outcomes



*There is a gender gap in earnings for the alumni at every top university, although the size of the difference varies greatly. Credit Matt Rourke/Associated Press*

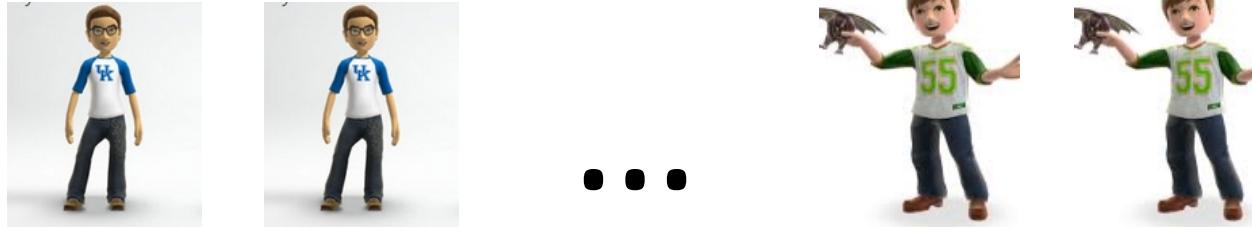
Gaps in Earnings Stand Out in Release of College Data

SEPT. 13, 2015  
[nytimes.com](http://nytimes.com) | Sept. 13, 2015



# How do predictive systems work?

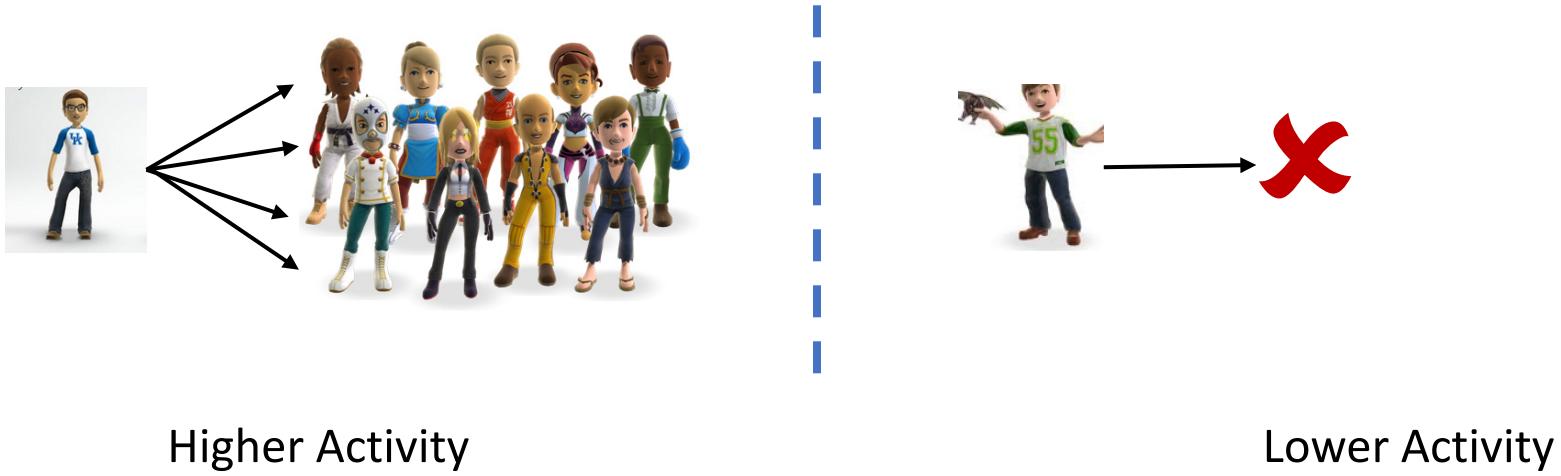
**Aim:** Predict future activity for a user.



We see data about their user profile and past activity.

E.g., for any user, we might see their age, gender, past activity and their social network.

# From data to prediction



Use these correlations to make a predictive model.

Future Activity ->

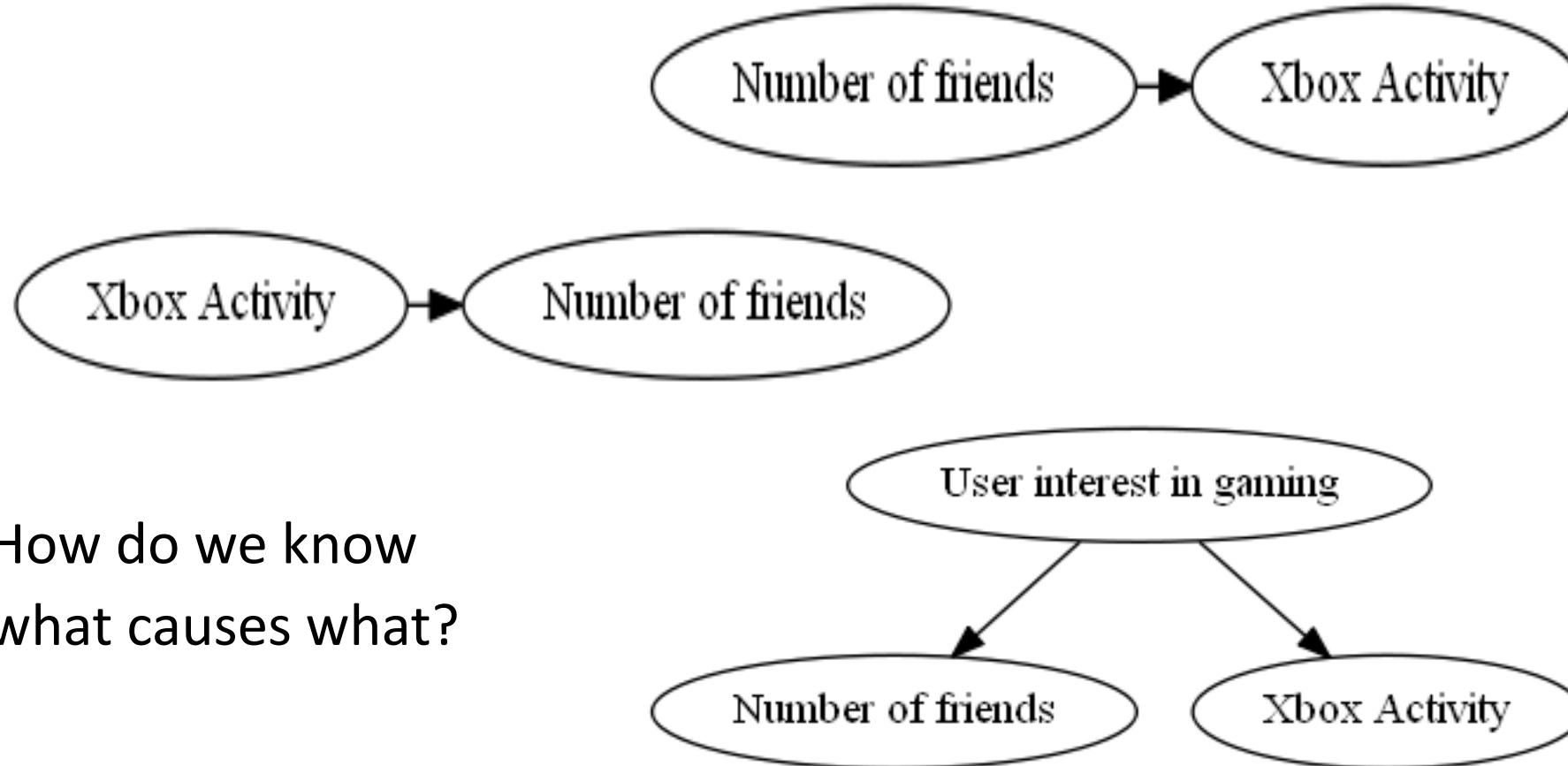
$$f(\text{number of friends, logins in past month})$$

# From data to “actionable insights”

Number of friends can predict activity with high accuracy.

How do we increase activity of users?

# Different explanations are possible



How do we know  
what causes what?

**Decision:** To increase activity, would it make sense to launch a campaign to increase friends?

## Another example: Search Ads

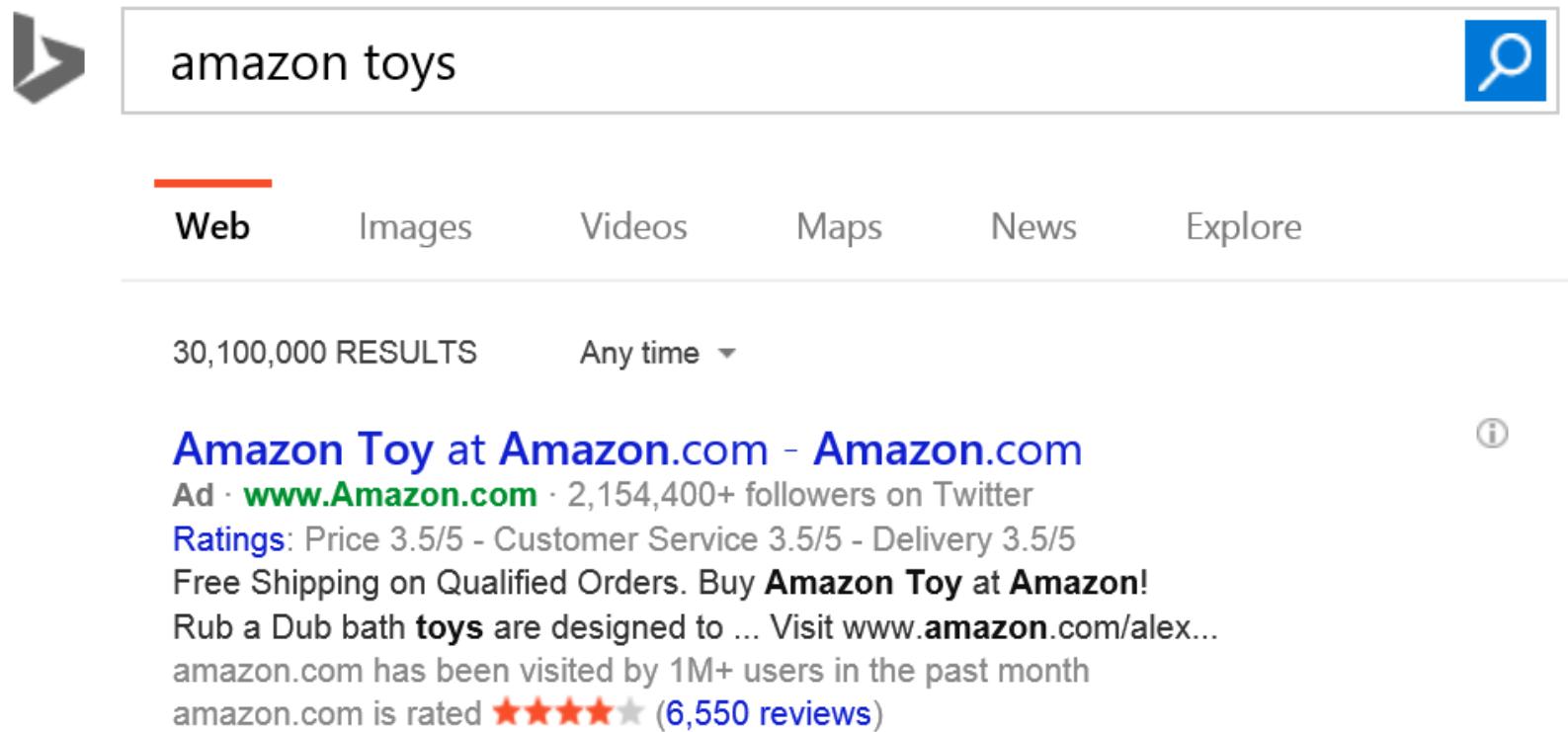


Search engines use ad targeting to show relevant ads.

Prediction model based on user's search query.

Search Ads have the highest click-through rate (CTR) in online ads.

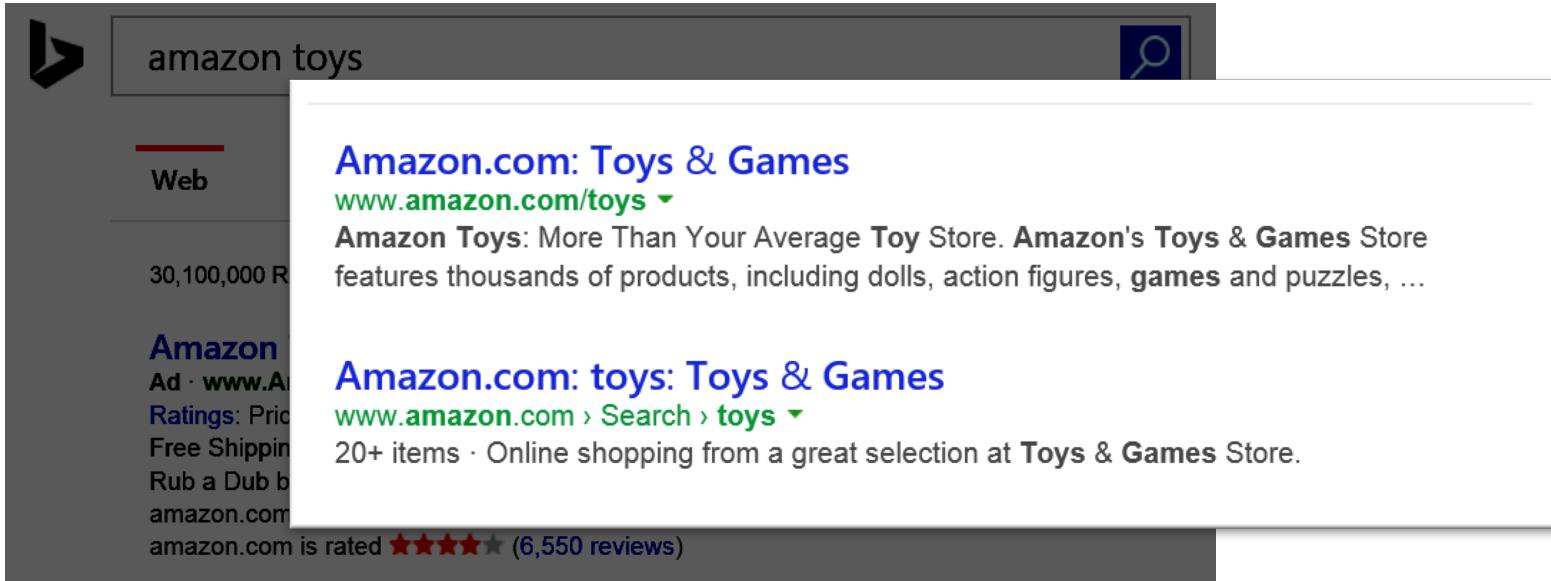
# Are search ads really that effective?



A screenshot of a Bing search results page. The search bar at the top contains the query "amazon toys". Below the search bar is a navigation menu with tabs: Web (which is underlined in red), Images, Videos, Maps, News, and Explore. Underneath the navigation menu, it says "30,100,000 RESULTS" and "Any time ▾". The first search result is a link to "Amazon Toy at Amazon.com - Amazon.com". To the right of the link is an information icon (an "i" inside a circle). Below the link, it says "Ad · www.Amazon.com · 2,154,400+ followers on Twitter". It also displays ratings: "Ratings: Price 3.5/5 - Customer Service 3.5/5 - Delivery 3.5/5", "Free Shipping on Qualified Orders. Buy Amazon Toy at Amazon!", and a snippet of text: "Rub a Dub bath toys are designed to ... Visit www.amazon.com/alex...". It also mentions "amazon.com has been visited by 1M+ users in the past month" and "amazon.com is rated ★★★★☆ (6,550 reviews)".

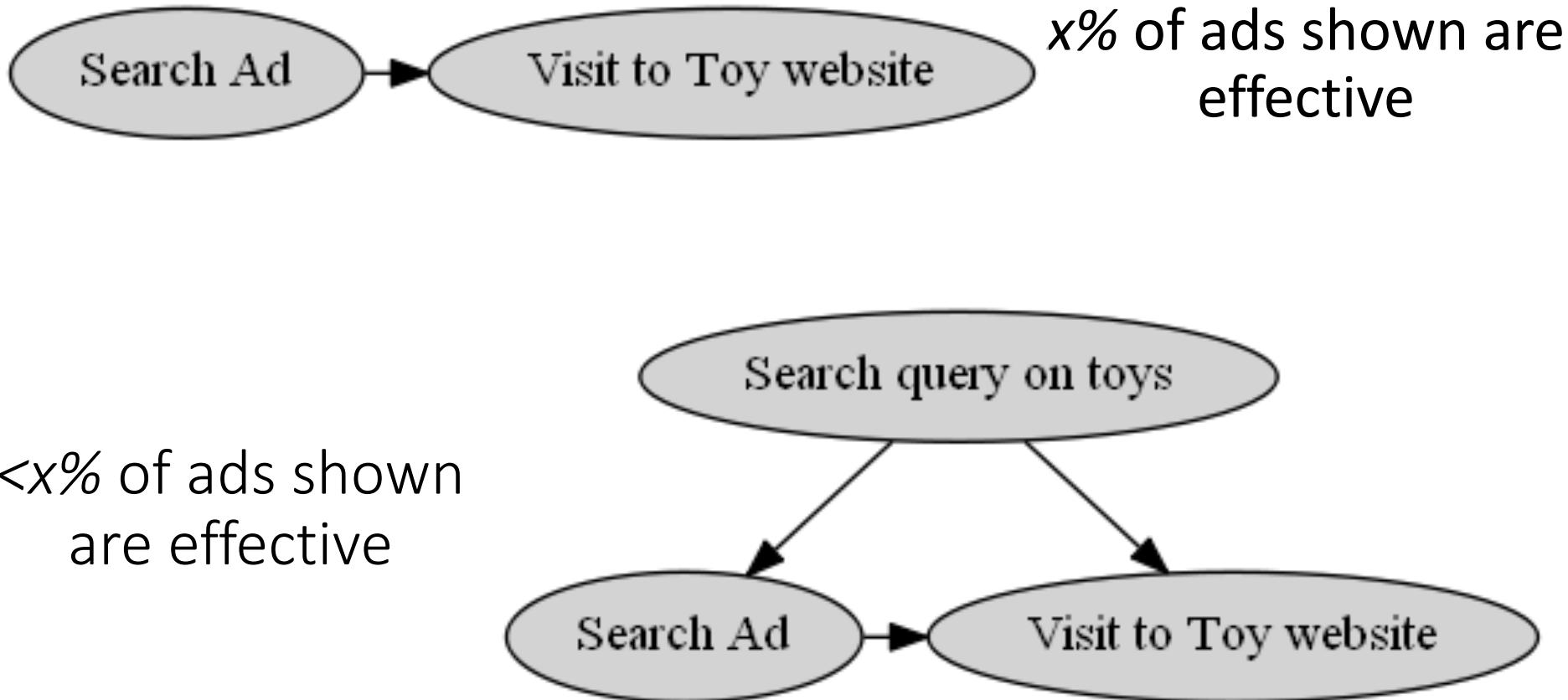
Ad targeting was highly accurate.

# But search results point to the same website



**Counterfactual question:** Would I have reached Amazon.com anyways, without the ad?

Without reasoning about causality, may overestimate effectiveness of ads



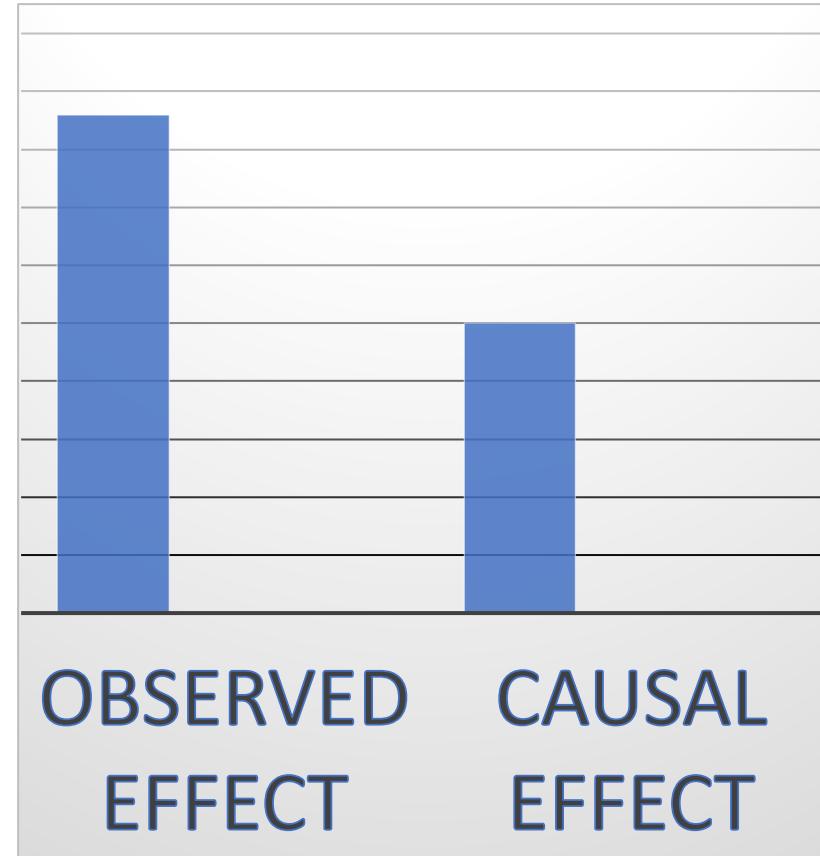
Okay, search ads have an explicit intent.  
Display ads should be fine?



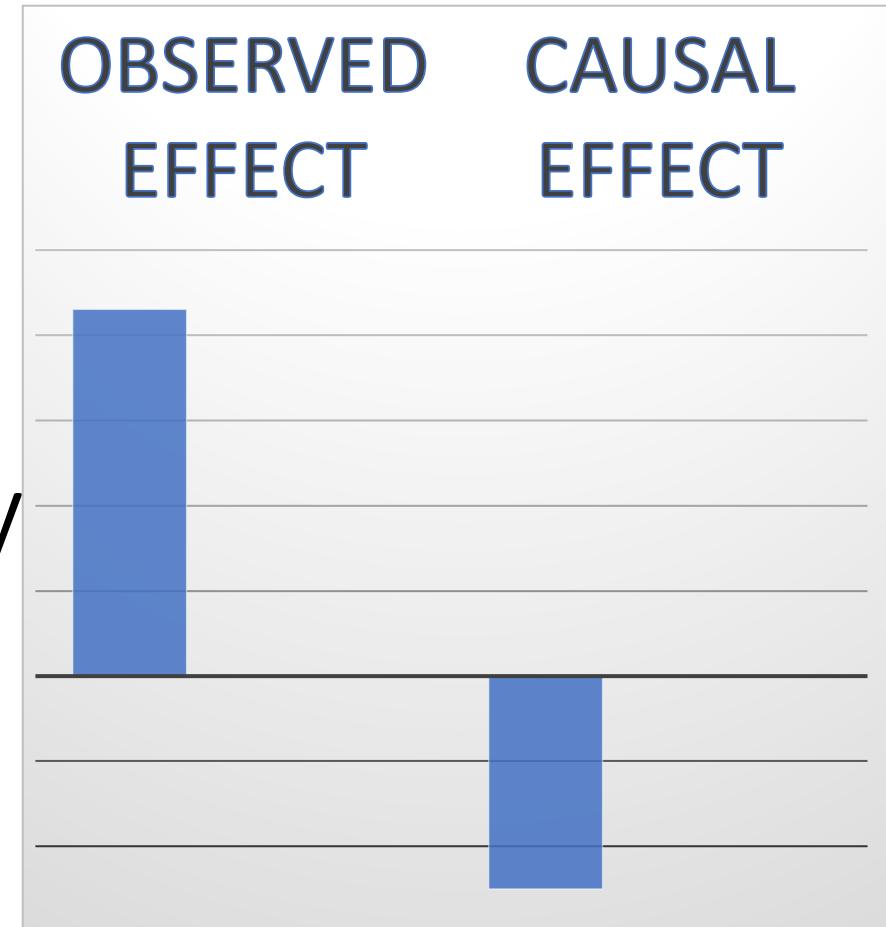
**Probably not.**

There can be many hidden causes for an action, some of which may be hard to quantify.

So far, so good. Be  
mindful of hidden  
causes, or else we  
might  
overestimate  
causal effects.



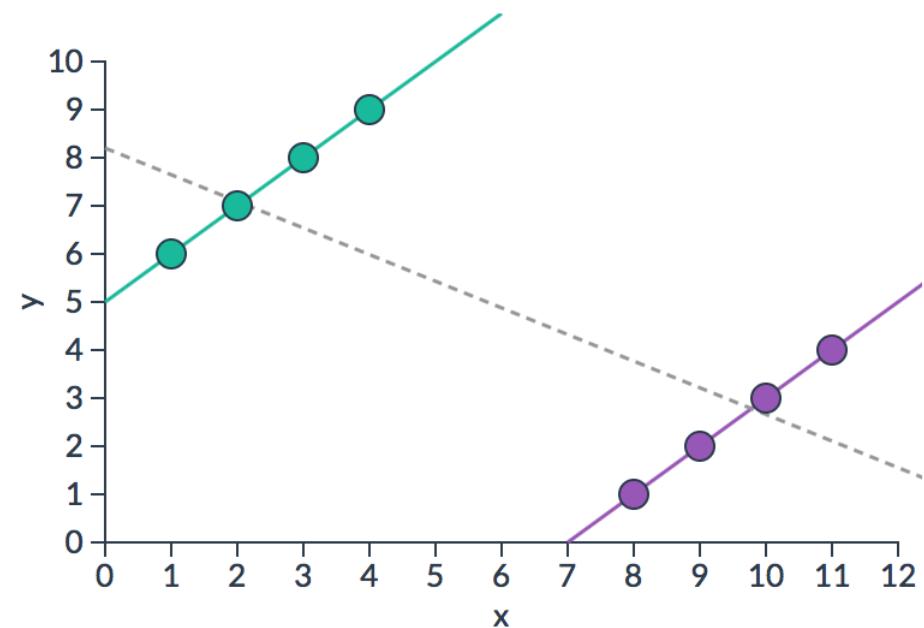
(But)  
Ignoring hidden  
causes can also  
lead to completely  
wrong  
conclusions.



# Simpson's paradox

---

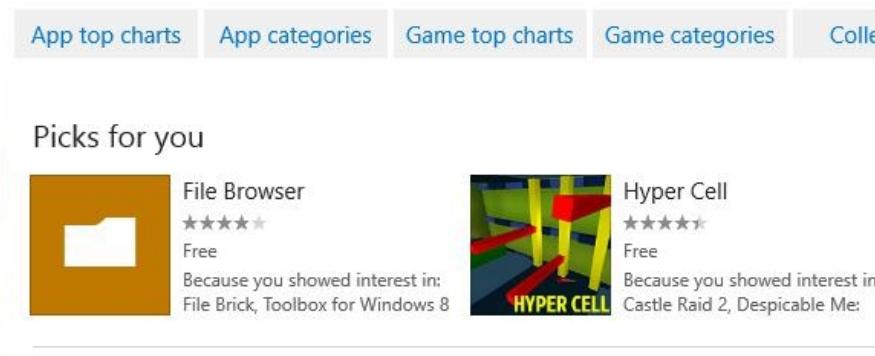
Selection bias can be so large  
that *observational and causal*  
*estimates give opposite effects*  
(e.g., going to hospitals makes  
you less healthy)



# Example: Which algorithm is better?

Have a current production algorithm. Want to test if a new algorithm is better.

Say recommendations on app store.



Algorithm A

?

Algorithm B

# Comparing old versus new algorithm

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

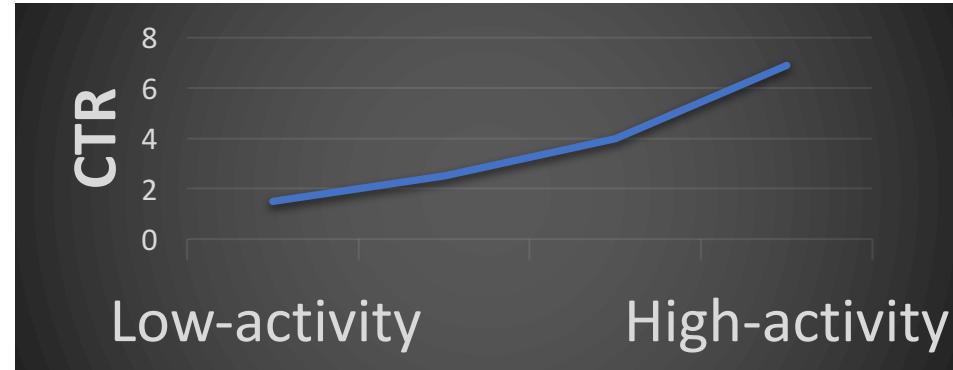
From system logs, collect data for 1000 sessions for each. Measure CTR.

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

New algorithm is better?

# Frequent users of the Store tend to be different from new users

So let us look at  
CTR separately.



Old Algorithm (A)	New Algorithm (B)	Low-activity Users
10/400 ( <b>2.5%</b> )	4/200 ( <b>2%</b> )	

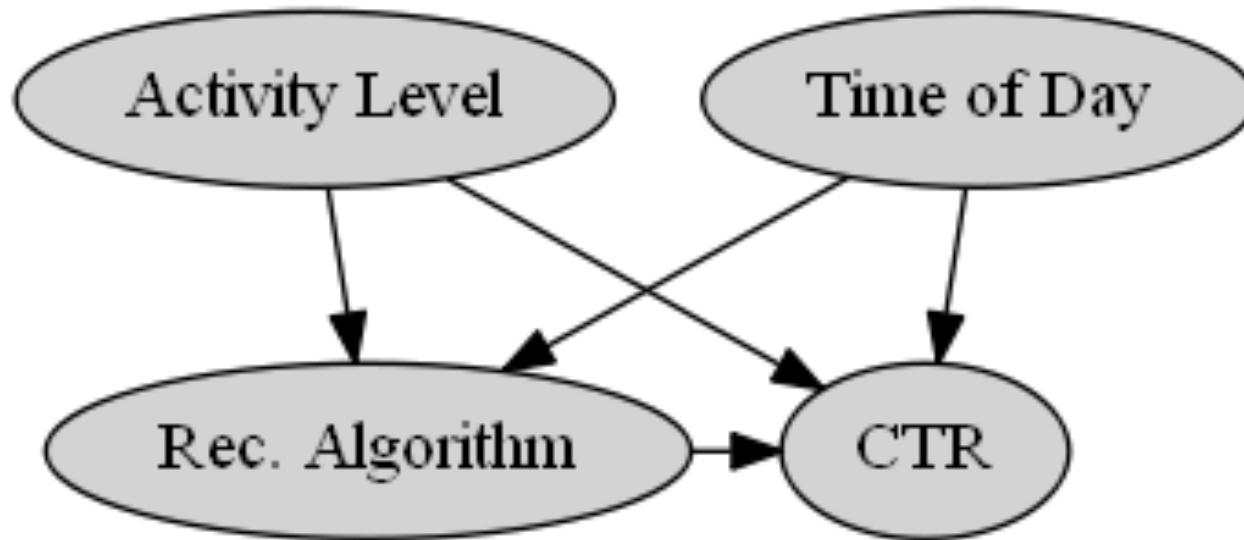
Old Algorithm (A)	New Algorithm (B)	High-activity Users
40/600 ( <b>6.6%</b> )	50/800 ( <b>6.2%</b> )	

# The Simpson's paradox

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

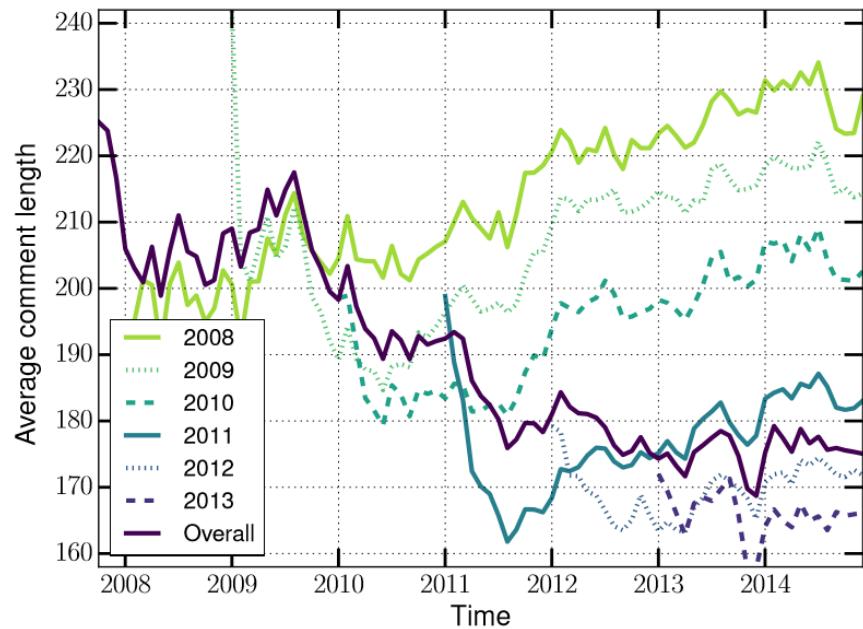
- Is Algorithm A better?

Answer (as usual): May be, may be not.



- E.g., Algorithm A could have been shown at different times than B.
- There could be other hidden causal variations.

# 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

# Intro to Experiments

Andrew #1

“To find out what happens when you change something, it is necessary to change it.”

---

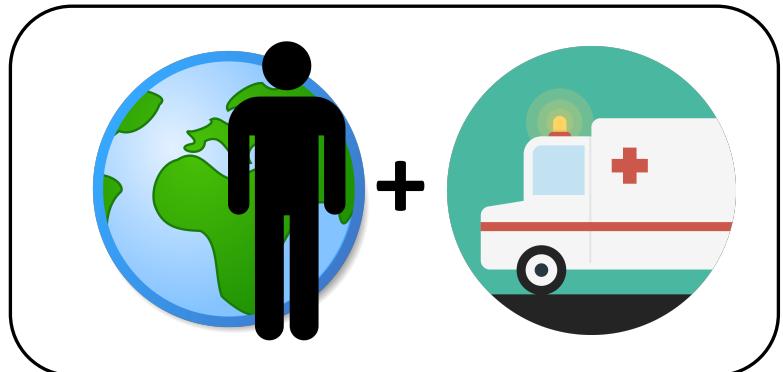
-GEORGE BOX

# Counterfactuals

---

To isolate the causal effect, we have to *change one and only one thing* (hospital visits), and compare outcomes

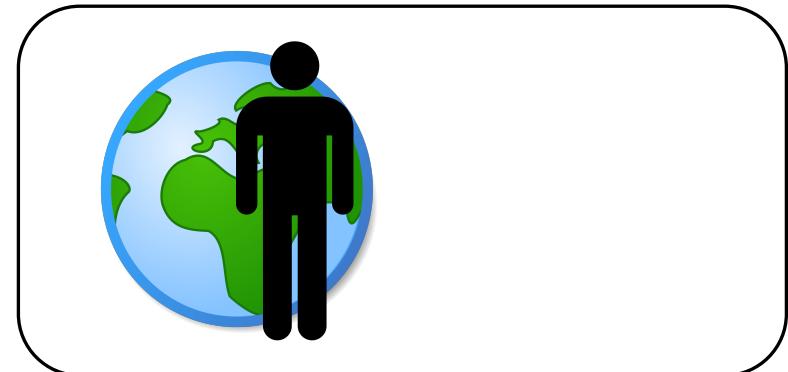
**Reality**



(what happened)

**vs**

**Counterfactual**



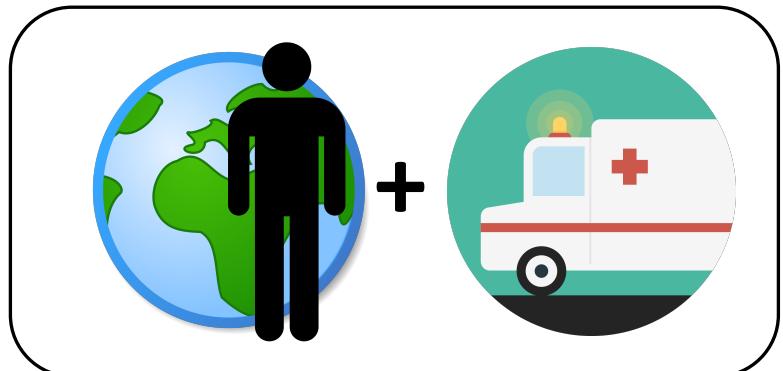
(what would have happened)

# Counterfactuals

---

We never get to observe *what would have happened if we did something else*, so we have to estimate it

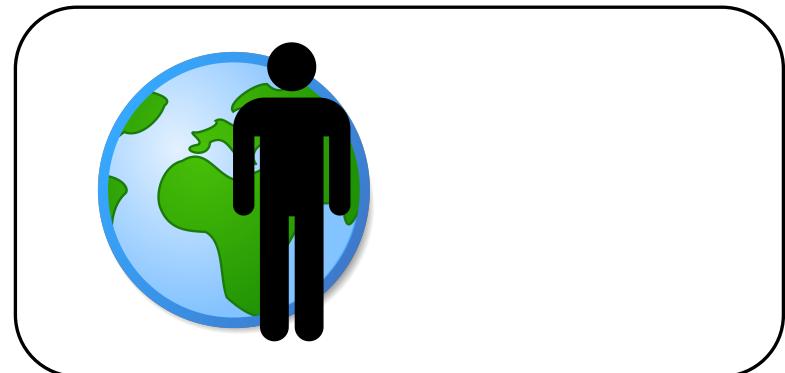
**Reality**



(what happened)

**vs**

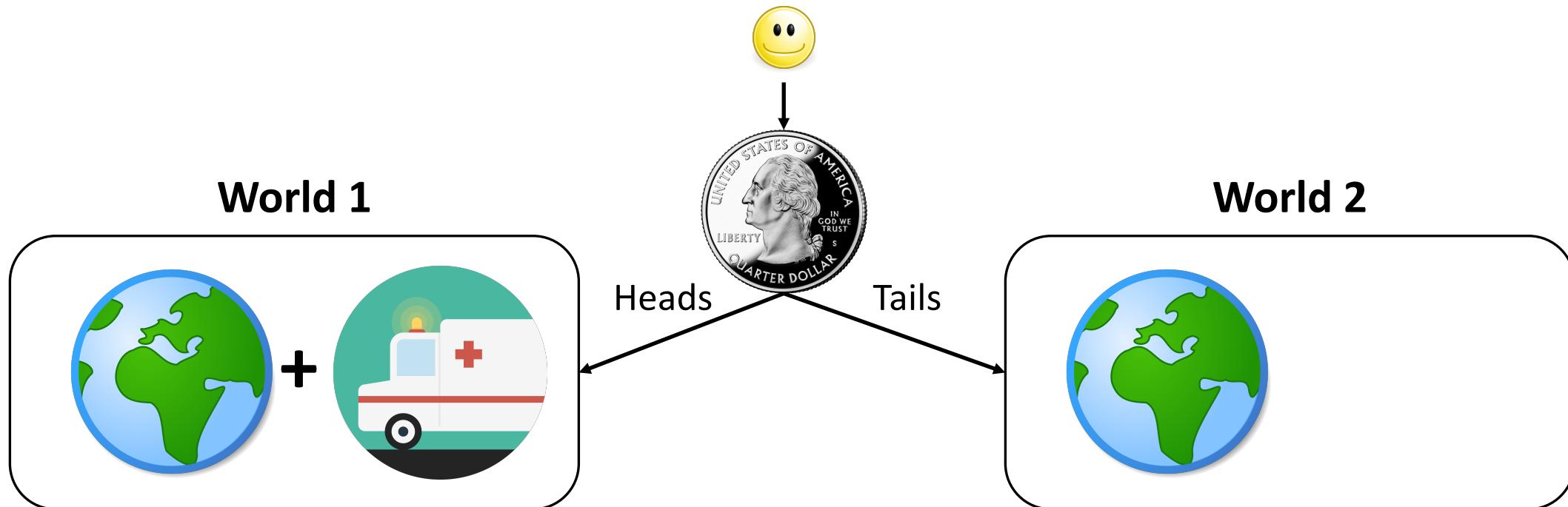
**Counterfactual**



(what would have happened)

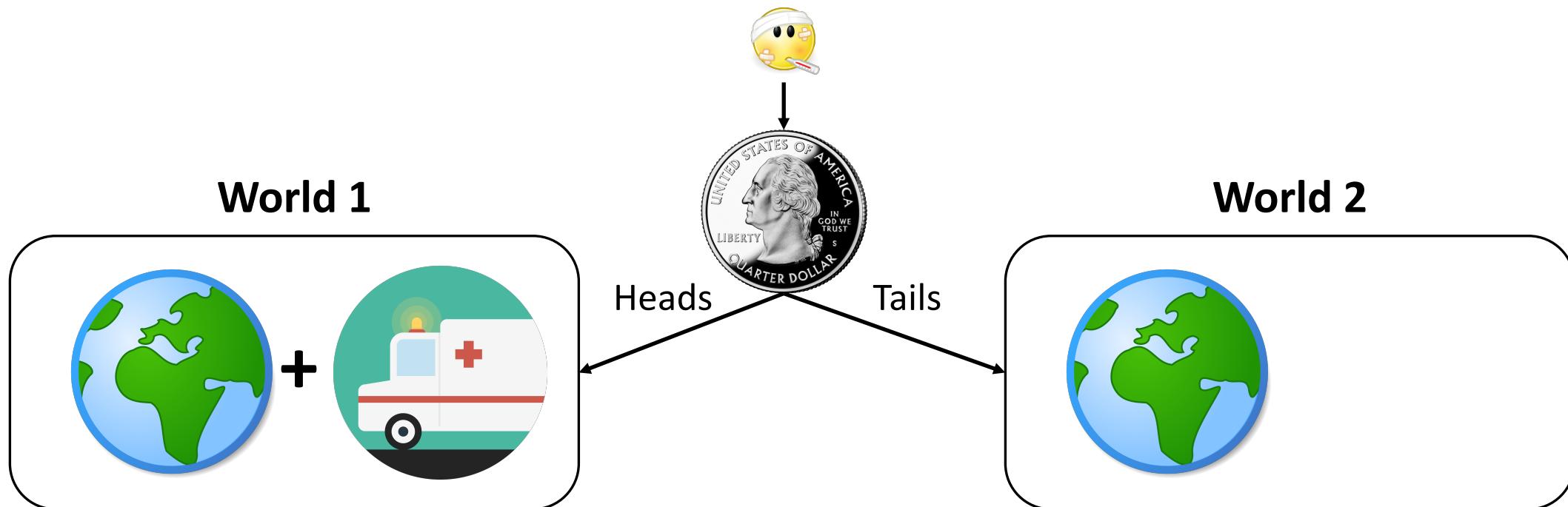
# Random assignment

We can use randomization to create two groups that differ only in which treatment they receive, restoring symmetry



# Random assignment

We can use randomization to create two groups that differ only in which treatment they receive, restoring symmetry

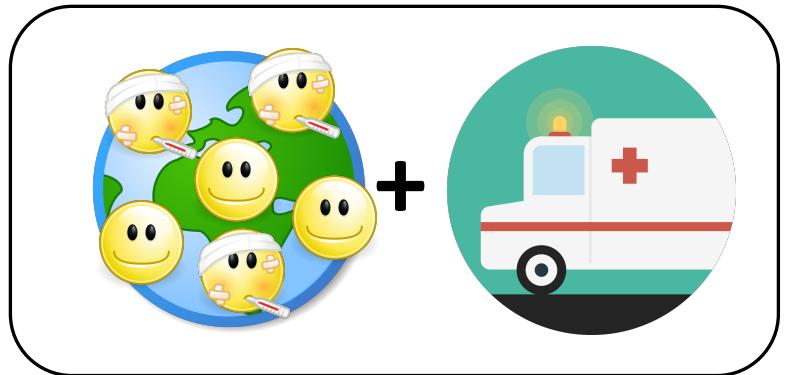


# Random assignment

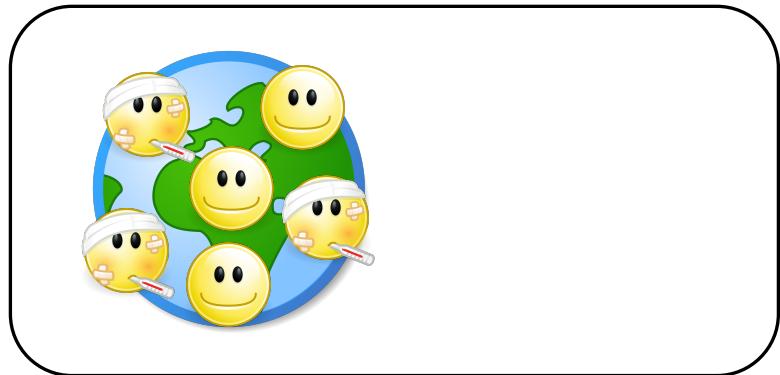
---

We can use randomization to create two groups that differ only in which treatment they receive, restoring symmetry

**World 1**



**World 2**



# Basic identity of causal inference

---

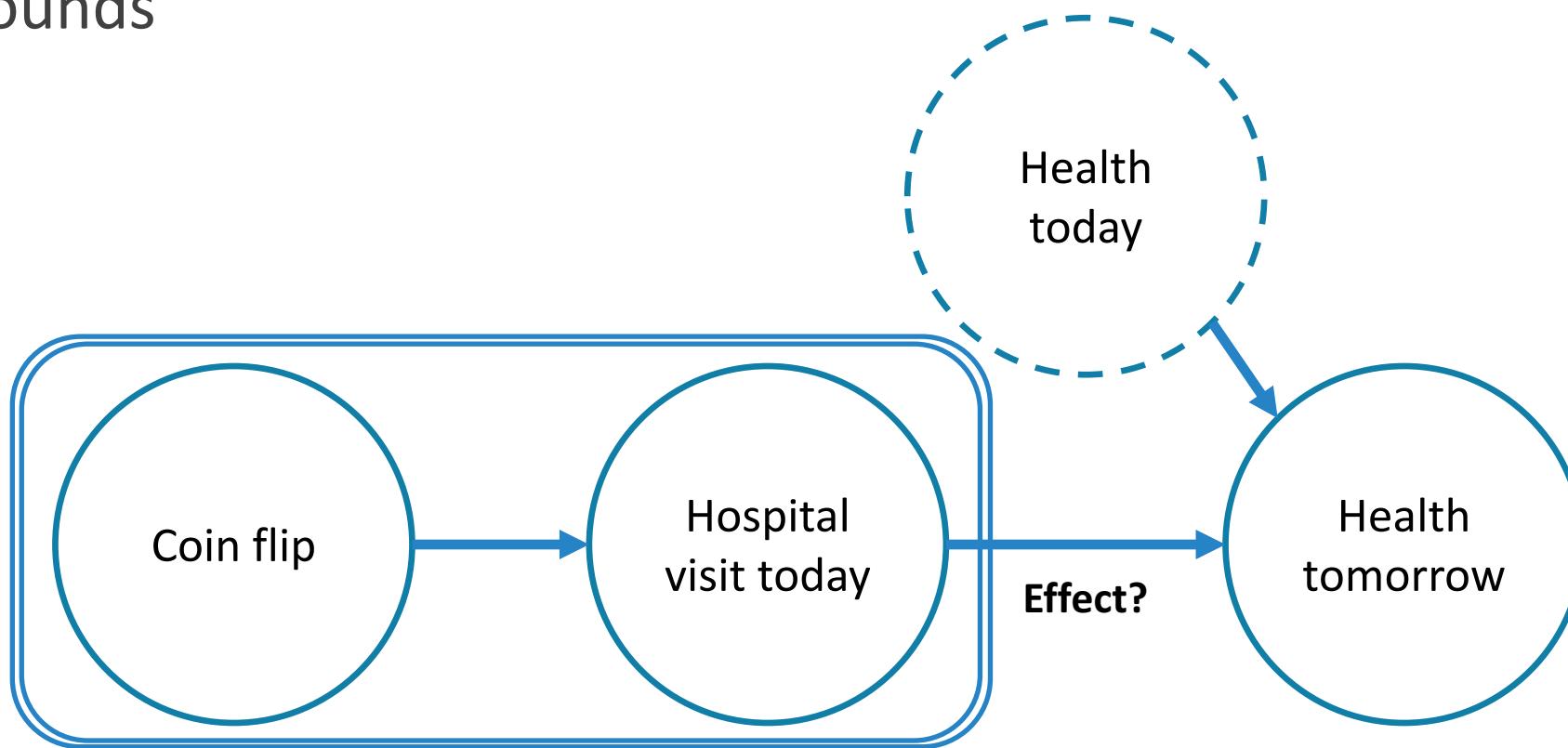
The observed difference is now the causal effect:

$$\begin{aligned}\text{Observed difference} &= \text{Causal effect} - \text{Selection bias} \\ &= \text{Causal effect}\end{aligned}$$

Selection bias is zero, since there's no difference, on average, between those who were hospitalized and those who weren't

# Random assignment

Random assignment determines the treatment independent of any confounds



Double lines mean  
“intervention”

# Experiments: Caveats / limitations

---

Random assignment is the “gold standard” for causal inference, but it has some limitations:

- Randomization often isn’t feasible and/or ethical
- Experiments are costly in terms of time and money
- It’s difficult to create convincing parallel worlds
- Inevitably people deviate from their random assignments

**Anyone can flip a coin, but it’s difficult to create convincing parallel worlds**

# Two goals for experiments

## **Internal validity:**

Could anything other than the treatment (i.e. a confound) have produced this outcome?

Did doctors give the experimental drug to some especially sick patients (breaking randomization) hoping that it would save them?

## **External validity (Generalization)**

Do the results of the experiment hold in settings we care about?

Would this medication be just as effective outside of a clinical trial, when usage is less rigorously monitored?

# How we conduct behavioral experiments



Lab Experiment

Field Experiment



**Better internal validity** (correctness):

- Greatest procedural control
- Can carefully curate situations

**But** less external validity (generalization):

- Artificial context, simple tasks
- Demand effects
- Homogeneous (WEIRD) subject pools
- Time/scale limitations

**Better generalization**

- Experiment findings apply to at least one real-world setting

**But:**

- Less control, more potential confounds
- Demand of experiment conflict with goals of real organizations
- More effort to conduct and manage
- More room for error

# Examples of Causal Inference

Amit #2

# Natural experiments

---

Sometimes we get lucky and nature effectively runs experiments for us, e.g.:

- As-if random: People are randomly exposed to water sources
- Instrumental variables: A lottery influences military service
- Discontinuities: Star ratings get arbitrarily rounded
- Difference in differences: Minimum wage changes in just one state

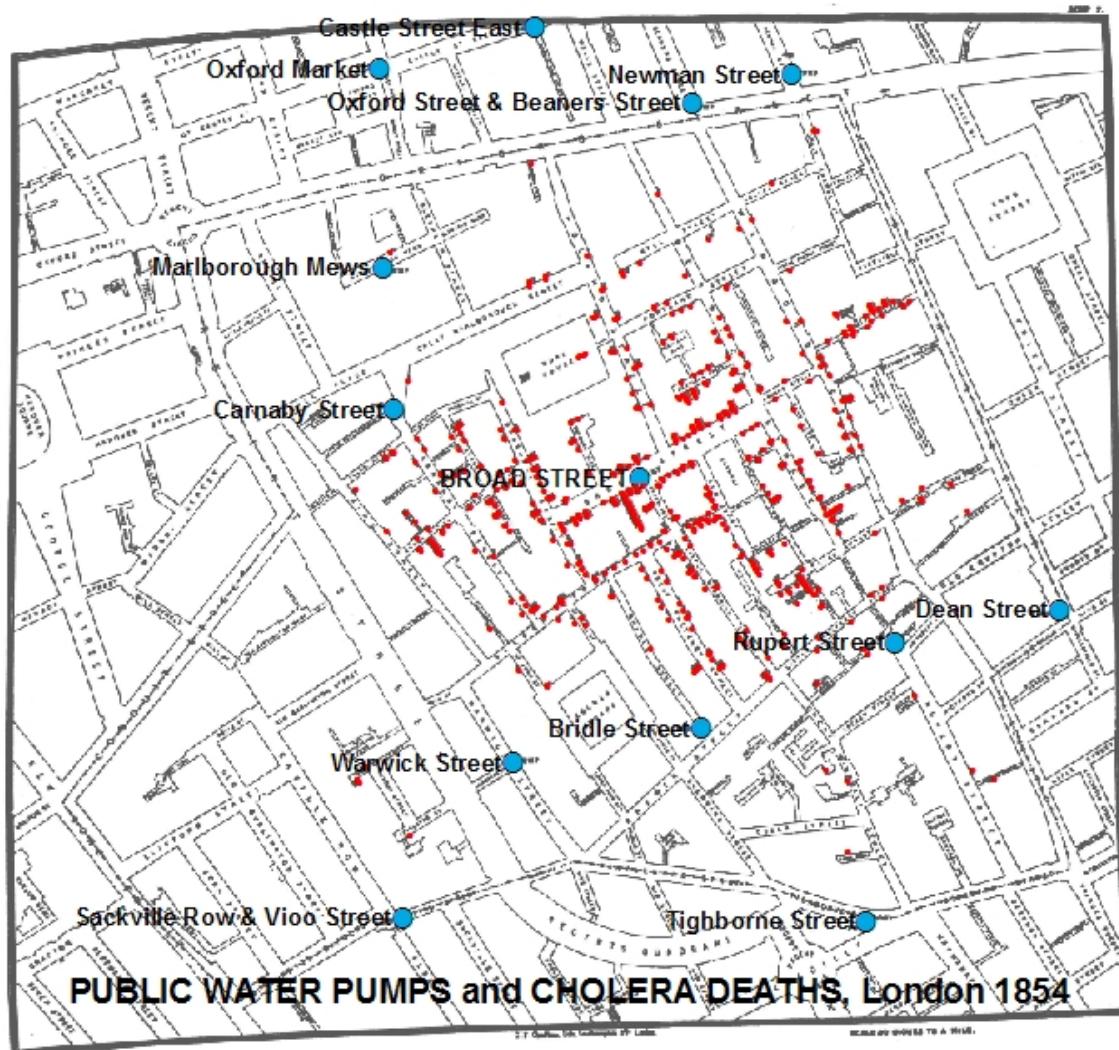
# Natural experiments

---

Sometimes we get lucky and nature effectively runs experiments for us, e.g.:

- As-if random: People are randomly exposed to water sources
- Instrumental variables: A lottery influences military service
- Discontinuities: Star ratings get arbitrarily rounded
- Difference in differences: Minimum wage changes in just one state

Experiments happen all the time, we just have to notice them



## As-if random

---

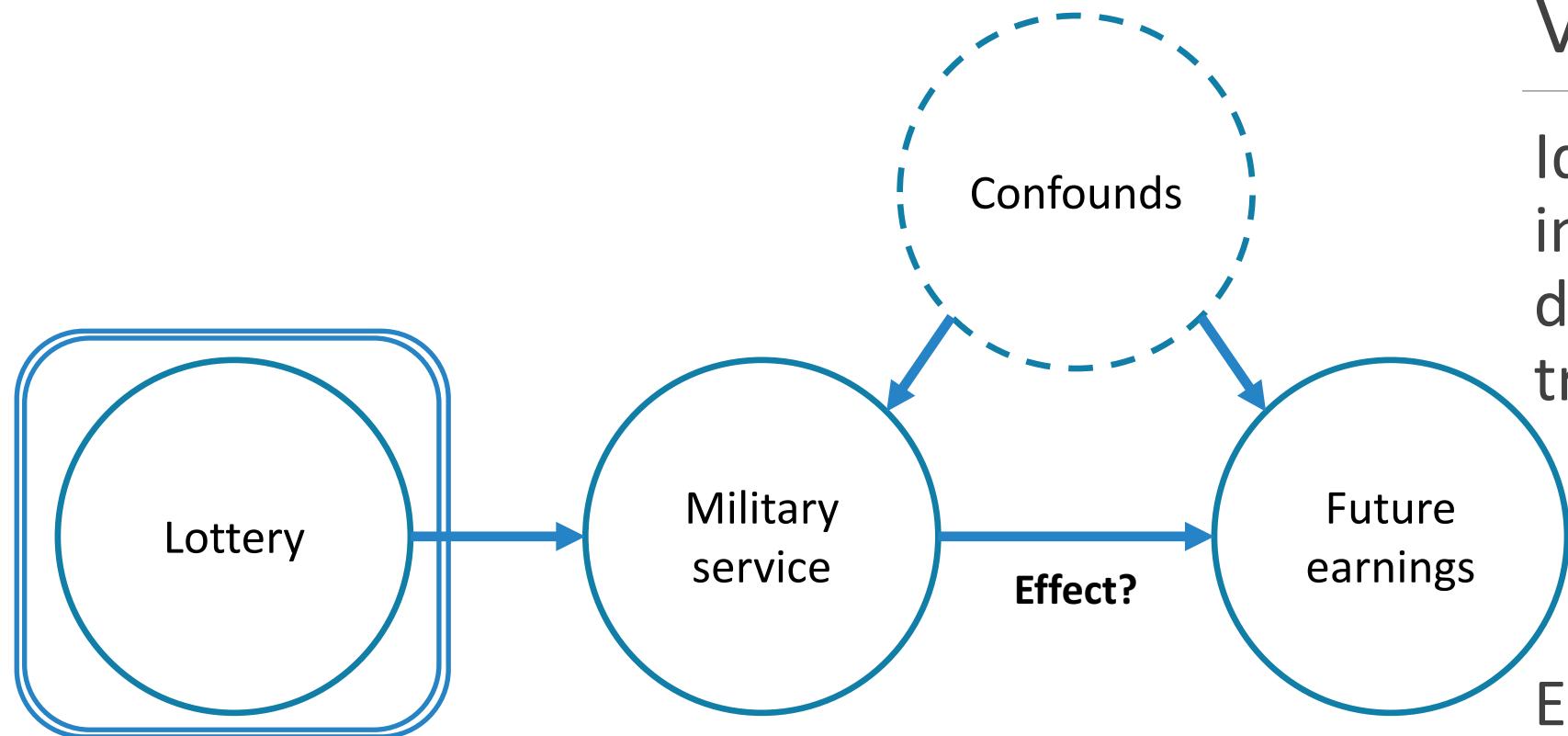
Idea: Nature randomly assigns conditions

Example: People are randomly exposed to water sources (Snow, 1854)

# Instrumental variables

---

Idea: An instrument independently shifts the distribution of a treatment

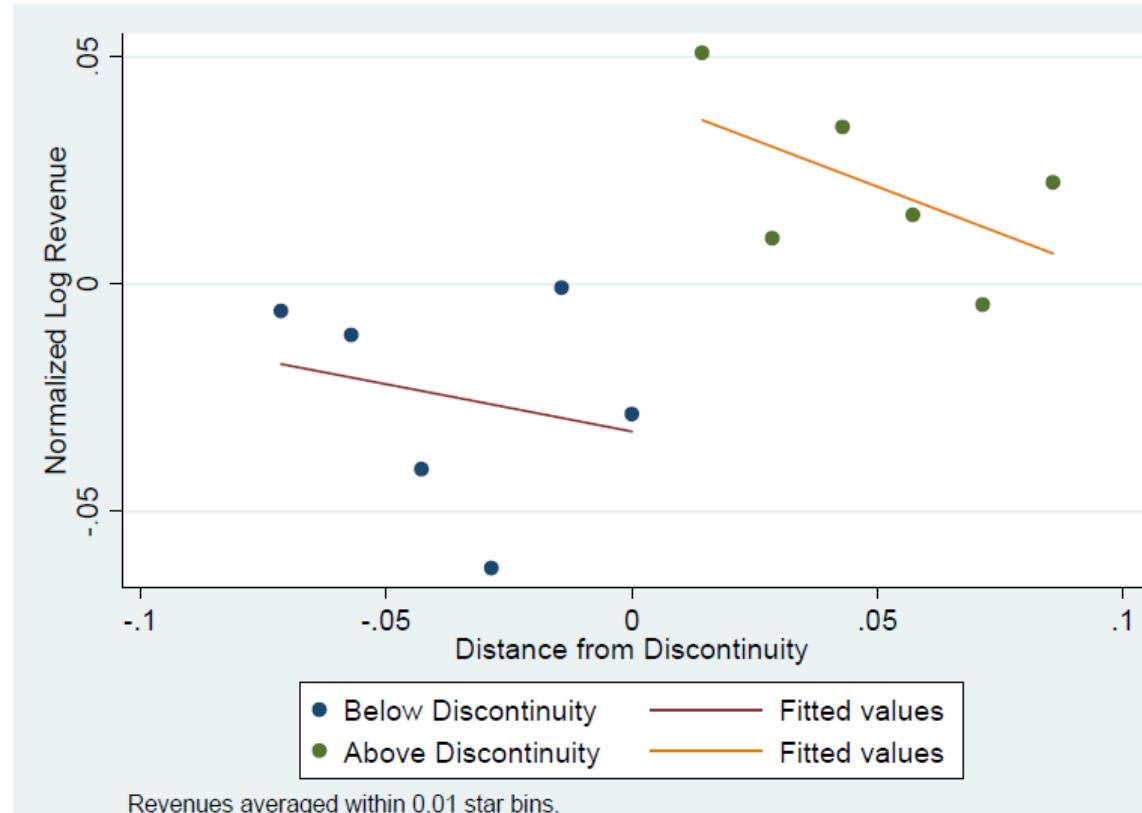


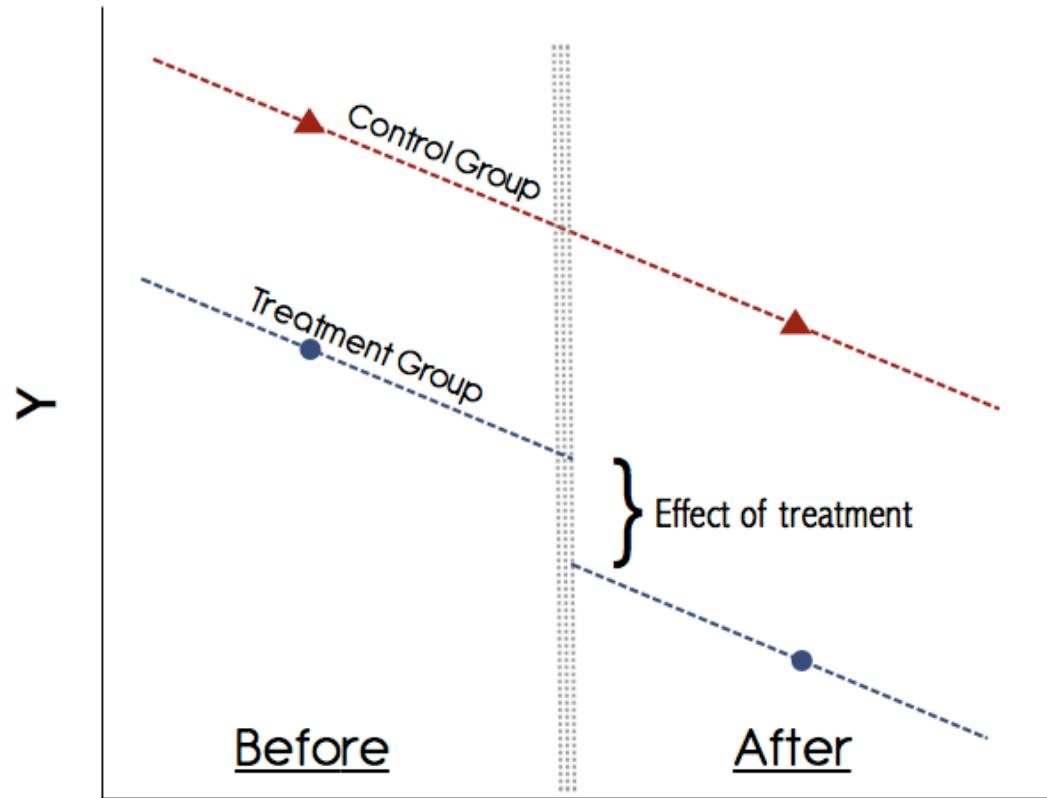
Example: A lottery influences military service (Angrist, 1990)

# Regression discontinuities

Idea: Things change around an arbitrarily chosen threshold

Example: Star ratings get arbitrarily rounded (Luca, 2011)



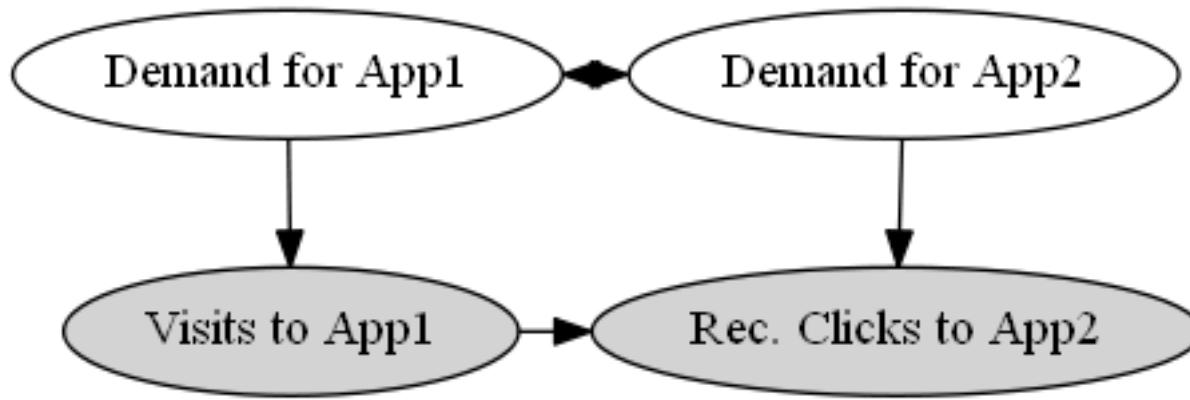


# Difference in differences

Idea: Compare differences after a sudden change with trends in a control group

Example: Minimum wage changes in just one state (Card & Krueger, 1994)

# Cont. example: Effect of store recommendations

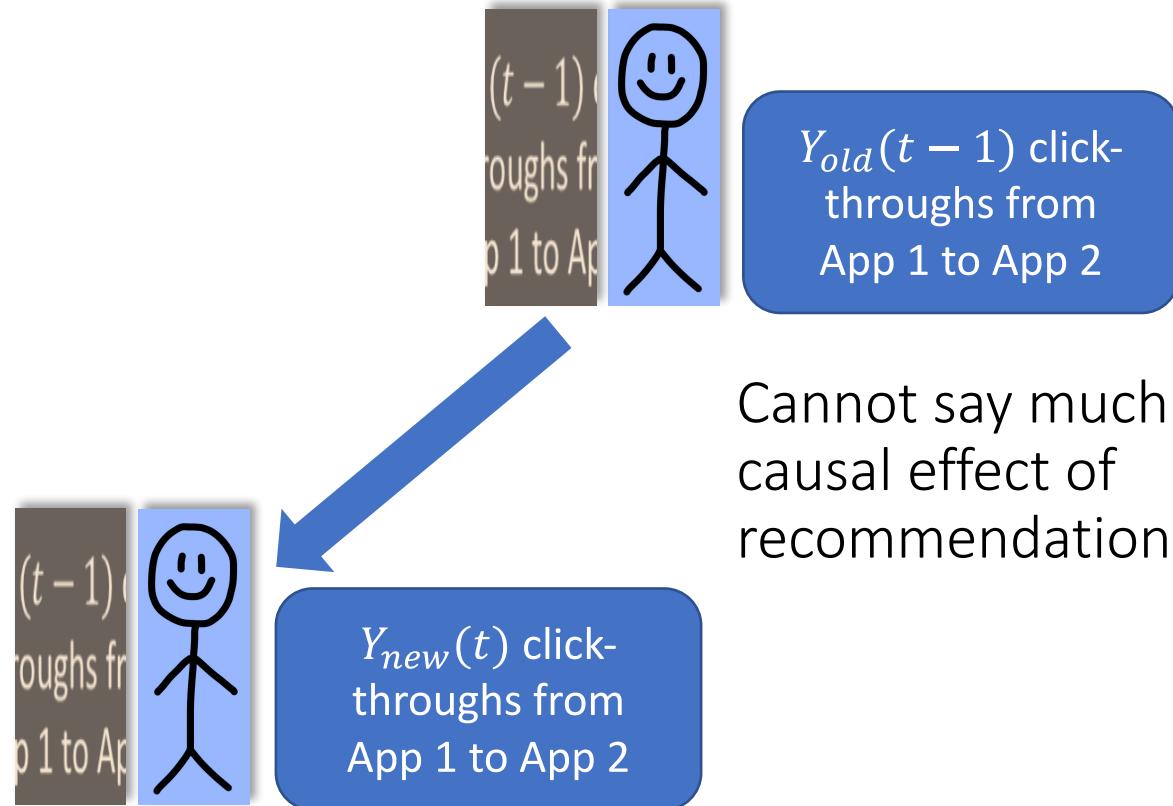


How many new visits are *caused by* the recommender system?

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

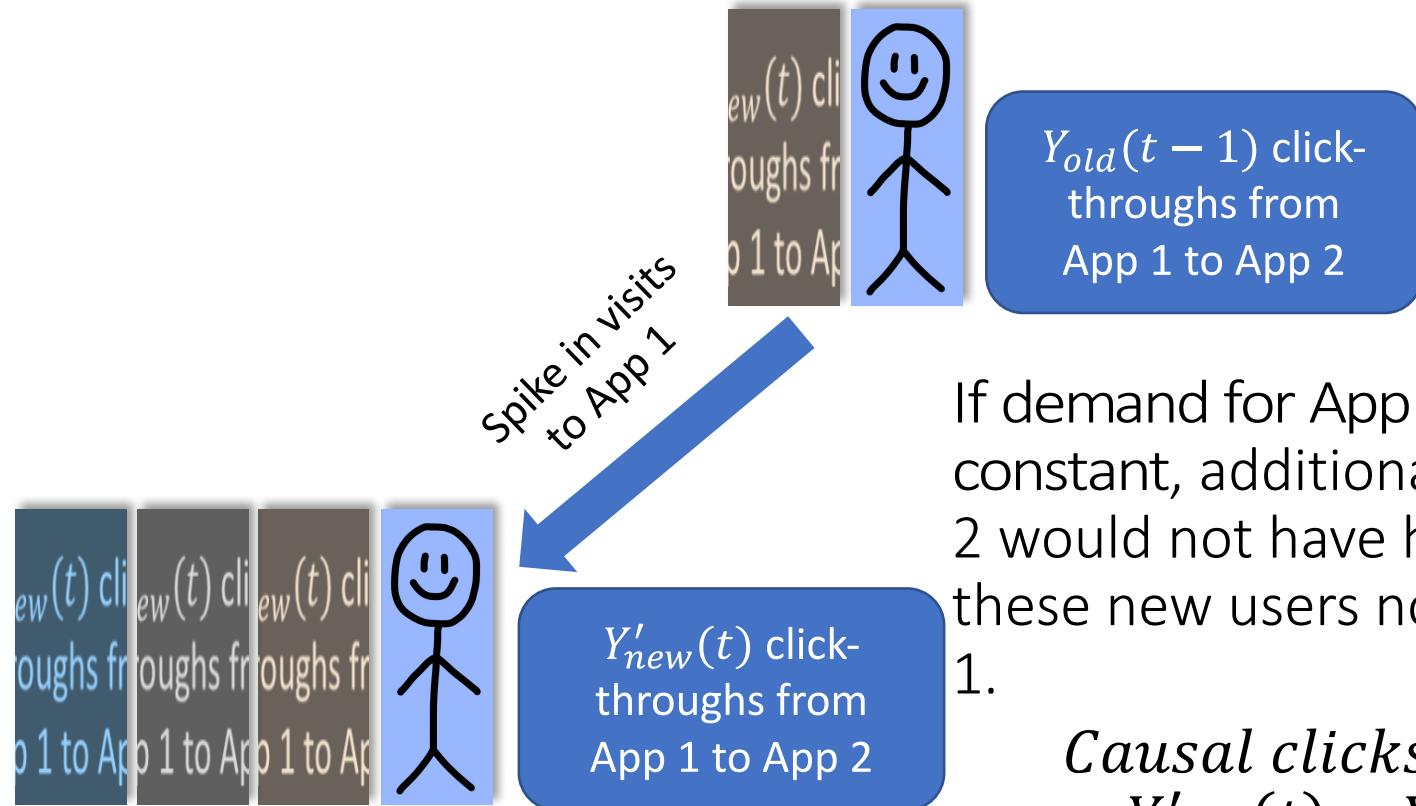
⇒ Users would most likely have visited App 2 even without recommendations.

# Traffic on normal days to App 1



Cannot say much about the causal effect of recommendations from App 1.

# External shock brings as-if random users to App1



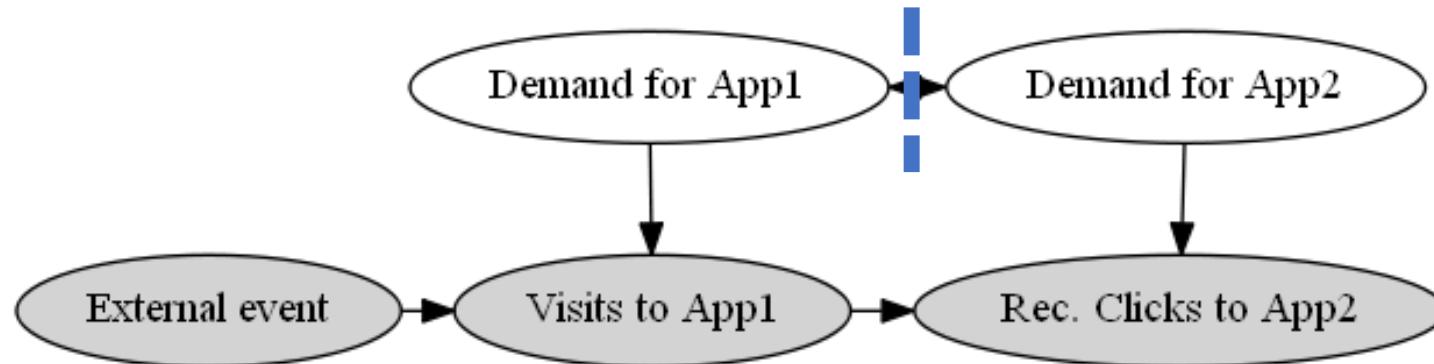
If demand for App 2 remains constant, additional views to App 2 would not have happened had these new users not visited App 1.

$$\begin{aligned} \text{Causal clicks} \\ = Y'_{\text{new}}(t) - Y_{\text{old}}(t - 1) \end{aligned}$$

# 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 =  $\frac{\Delta(\text{Rec. Click-throughs from App1 to App2})}{\Delta(\text{Visits to App1})}$



# Natural experiments: Caveats

---

Natural experiments are great, but:

- Good natural experiments are hard to find
- They rely on many (untestable) assumptions
- The treated population may not be the one of interest

# Natural experiments: Caveats

---

Natural experiments are great, but:

- Good natural experiments are hard to find
- They rely on many (untestable) assumptions
- The treated population may not be the one of interest

Sometimes we can use *additional data + algorithms* to  
*automatically find* natural experiments

# What can we do without an experiment or natural experiment?

- Condition on the “right” variables
- How to find the right variables?
  - No good answer
  - Depends on domain knowledge
  - Many different frameworks (e.g. Pearlian graphical models and do-calculus)

# General method: Conditioning on variables

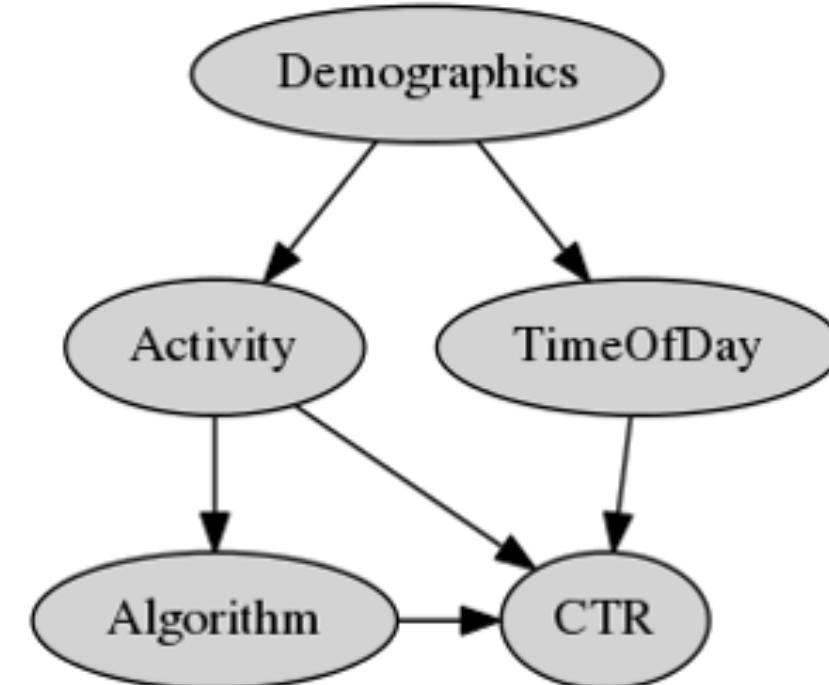
**Intuition:** Compare effect of algorithm on similar users.

⇒ Compare users with the same activity level.

Steps:

1. Stratify log data based on activity levels.
2. Compare CTR of different algorithms within these strata.

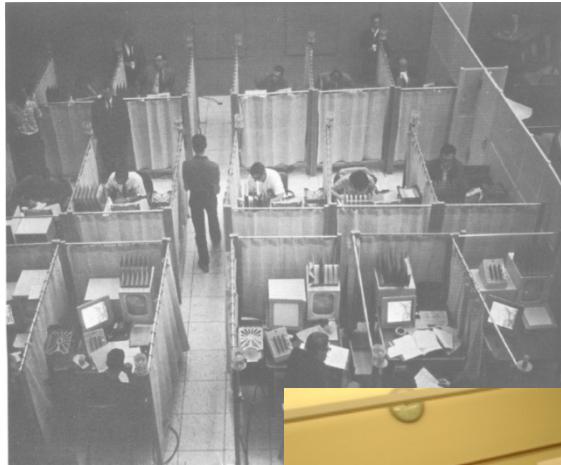
High-dimension? Can derive many complicated methods from this intuition.



# Examples of Experiments and Demo

Andrew #2

# Behavioral science labs are very limiting



ca. 1960s



ca. 2000s

- High degree of procedural control
- **Optimized for causal inference**

But, **many limitations:**

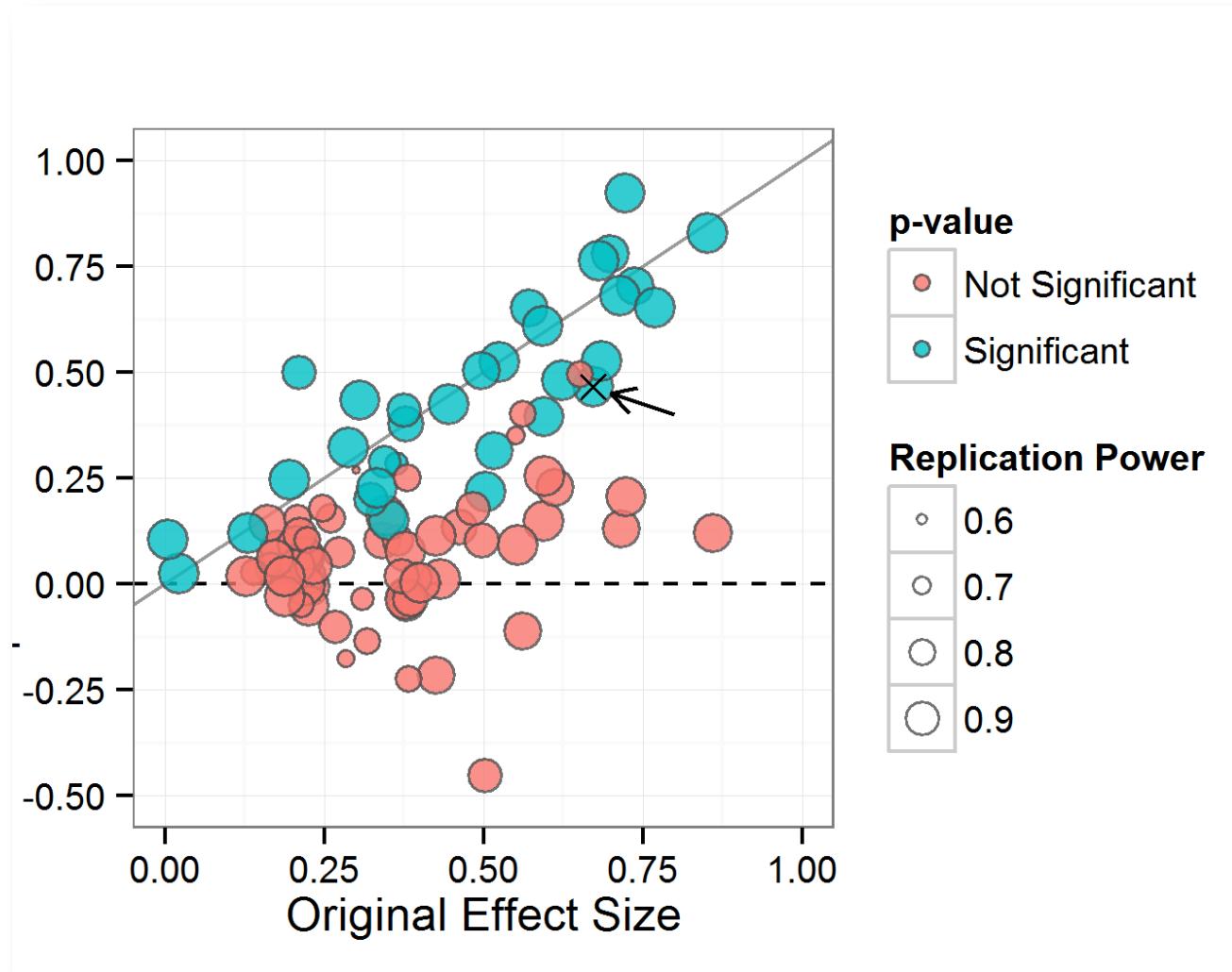
- Artificial environment
- Simple tasks, demand effects
- Homogeneous (WEIRD)\* subject pools
- Time/scale limitations
- Expensive, difficult to set up

**Poor generalization, expensive, slow**

\* Western, Educated, Industrialized, Rich, Democratic [Henrich et al. 2010]

Experiments are underpowered

Two-thirds of psychology studies don't replicate!



Estimating the Reproducibility  
of Psychological Science (2015)

# Most social science experiments aren't social

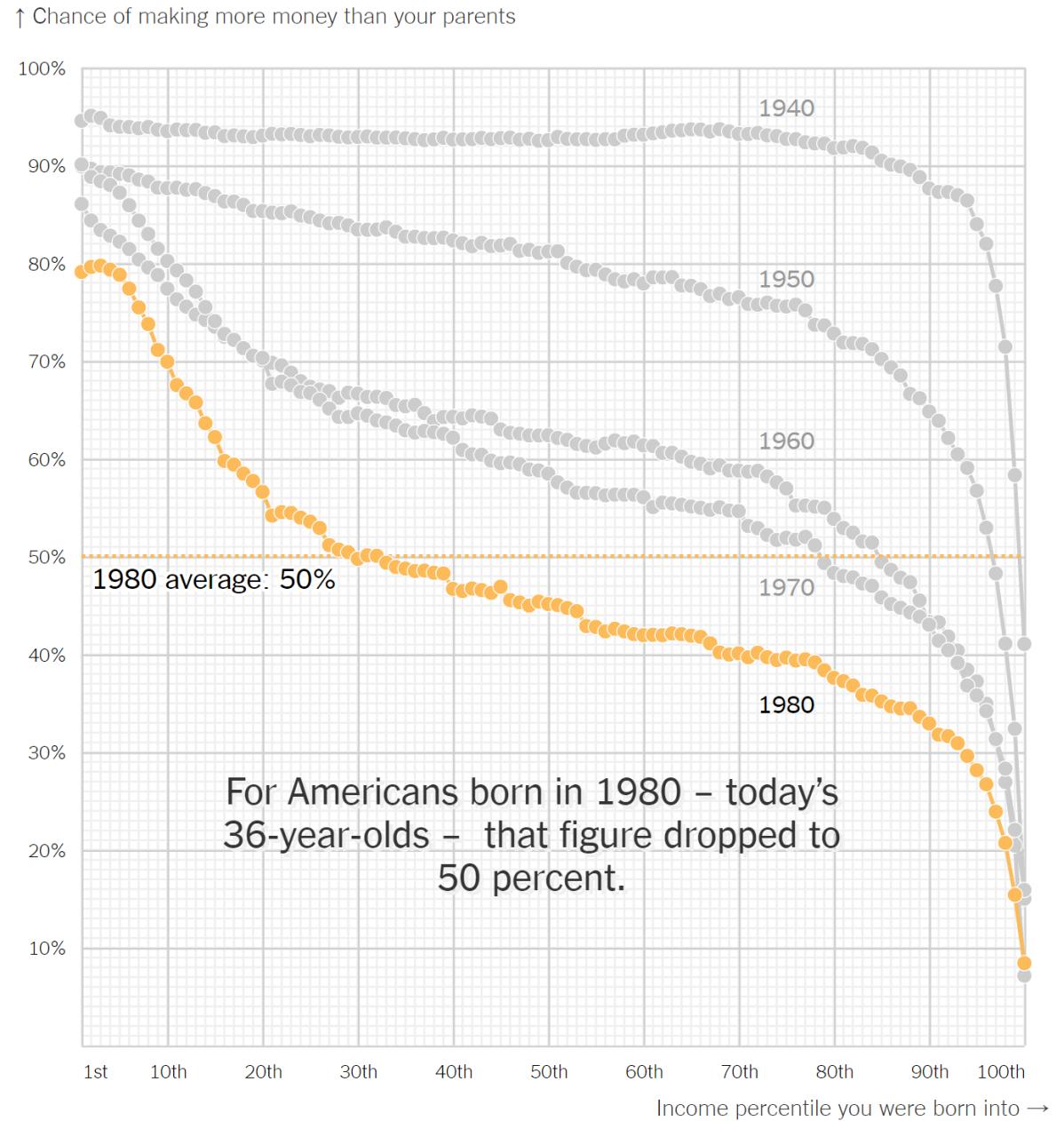
- Vast majority of experiments done with individuals, focused on individual behavior: logically, it's just easier
- **We actually don't know the causal effects of policies in many large-scale, collective behavior settings!**



# Economic Inequality

Social mobility: if you work hard, can you be more successful than your parents?

- How do social safety nets, economic redistribution, etc. affect social mobility and resulting inequality?
- What **policies can we enact** that will improve social mobility or reduce inequality?



*The American Dream, Quantified at Last.* NYTimes, Dec. 8 2016.

# Systems of Governance



*I think we should look to countries like Denmark, like Sweden and Norway, and learn from what they have accomplished for their working people*

What would happen if we suddenly replaced our current political system with a Scandinavian “socialist state”?

Would everyone be happier?

Or would it be culturally untenable?

# The black box of macroeconomic policy

There are no clear theories behind whether **lowering** interest rates, **tightening** interest rates, or **quantitative easing** work as intended.

Given a sufficiently large simulated economy, can we study these experimentally?

ECONOMY

## *Fed Raises Interest Rates for Third Time Since Financial Crisis*

By BINYAMIN APPELBAUM MARCH 15, 2017



Janet Yellen, the Federal Reserve chairwoman, announced the board's decision on interest rates. March 15, 2017. Photo by Al Drago/The New York Times. [Watch in Times Video »](#)



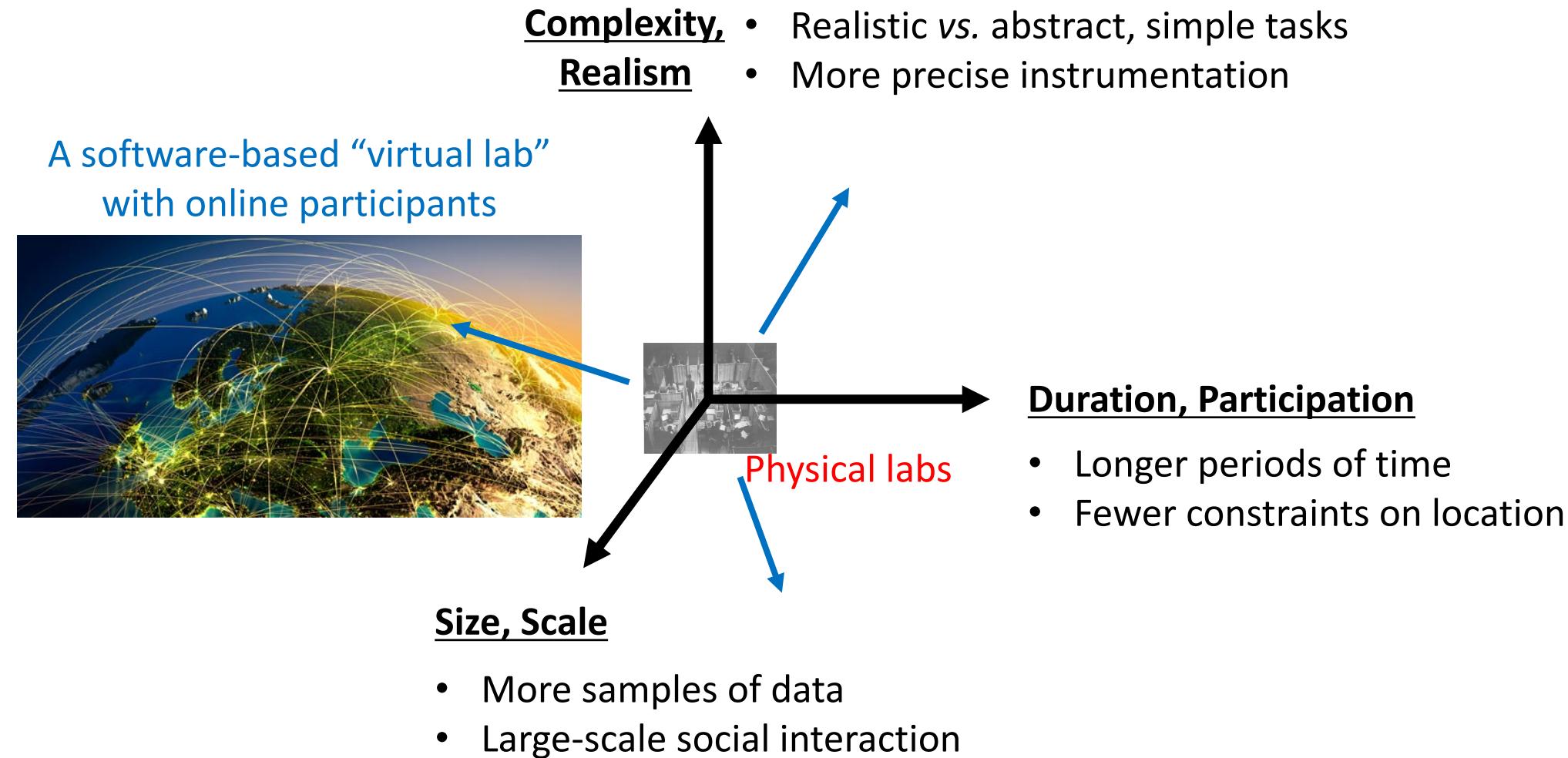
The [Federal Reserve](#), which raised its benchmark rate on Wednesday for the

# So, about them experiments...

- They're costly to run
- Yet limited in the types of questions they can answer
- Published experimental research is probably wrong
- They're far from answering many big important questions

How do we fix this?

# Expanding the experiment design space



# Behavioral experiments + the Internet



Control

The Virtual Lab

TEST MY BRAIN

amazon mechanical turk™  
Artificial Artificial Intelligence

Volunteer web labs

Inside the Largest Virtual Psychology Lab in the World

Riot Games wants you to behave yourself when you play League of Legends, so it's turned the game into a virtual lab

A central collage of images. On the left, there's a stylized brain icon with the text "TEST MY BRAIN" and "Volunteer web labs". Next to it is the "amazon mechanical turk™" logo with the tagline "Artificial Artificial Intelligence". Below these are two screenshots from the video game League of Legends, showing a match in progress. Overlaid on the bottom right of the collage is the text "Inside the Largest Virtual Psychology Lab in the World" and "Riot Games wants you to behave yourself when you play League of Legends, so it's turned the game into a virtual lab".

f facebook.

Google

Unwitting  
labs in the  
field

bing

Generalization

# Prisoner's Dilemma



Cooperate

Defect

Split – Split	50% - 50%
Steal – Split	100% - 0%
Steal – Steal	0% - 0%

Q: What happens to people's behavior when they are in a social dilemma for a long time?

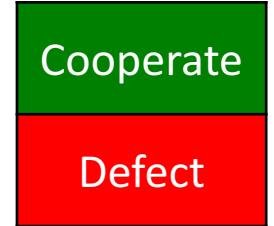
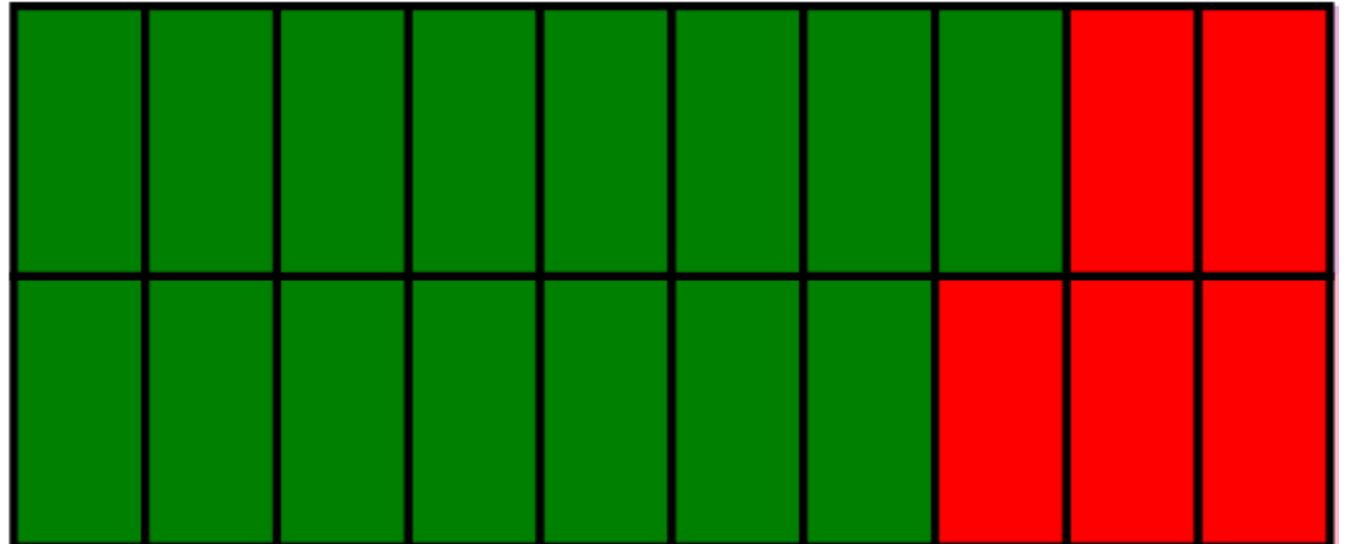
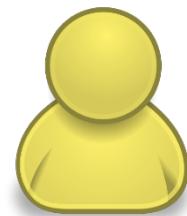
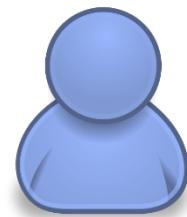
Also known as: real life.

A: According to lab experiments, we don't know. It seems like people start stabbing each other in the back when they play this game

...but the world seems to (mostly) be doing okay.

# Longitudinal behavior in a social dilemma

anonymous  
partners

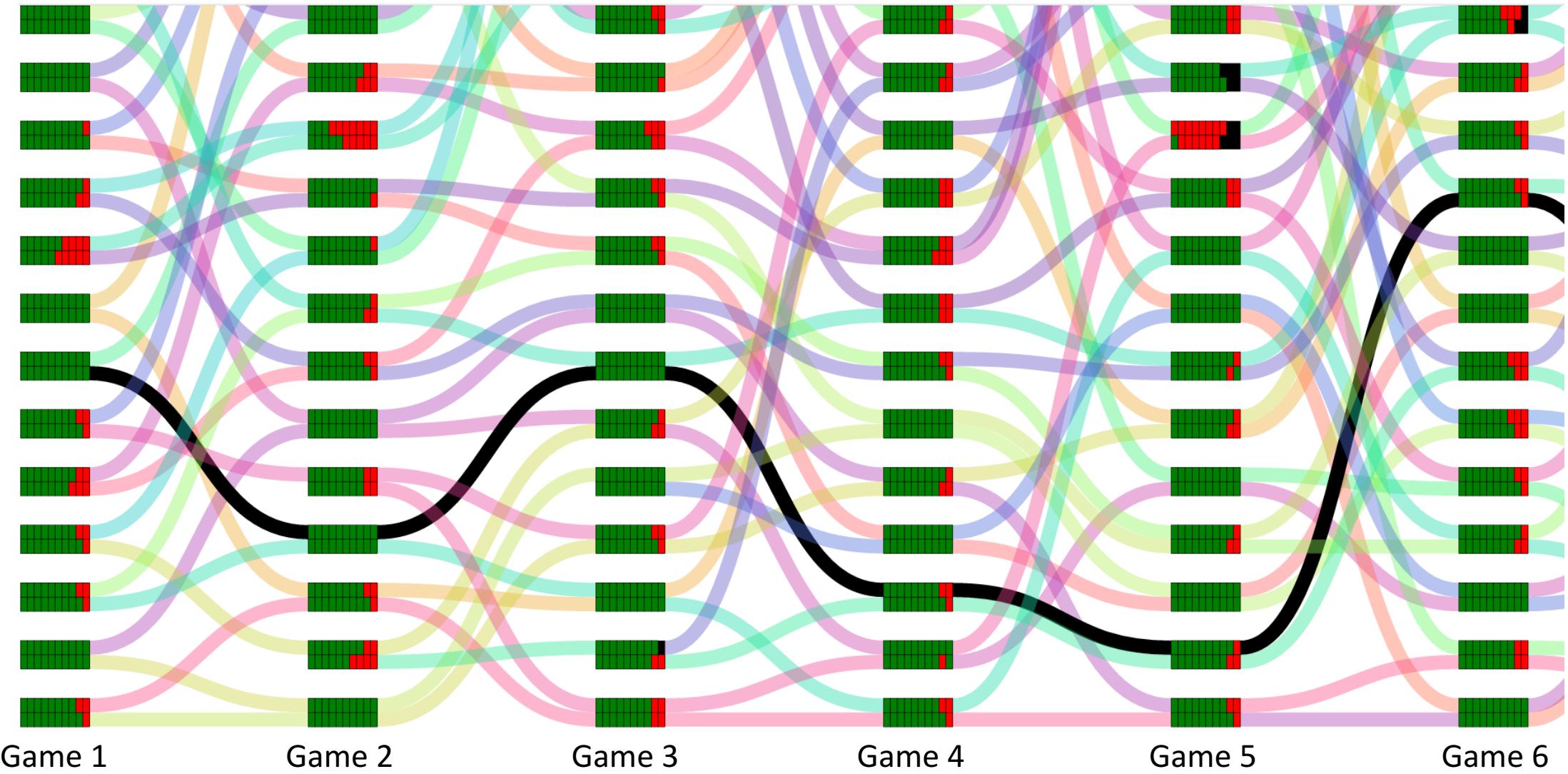


# Random rematching across games

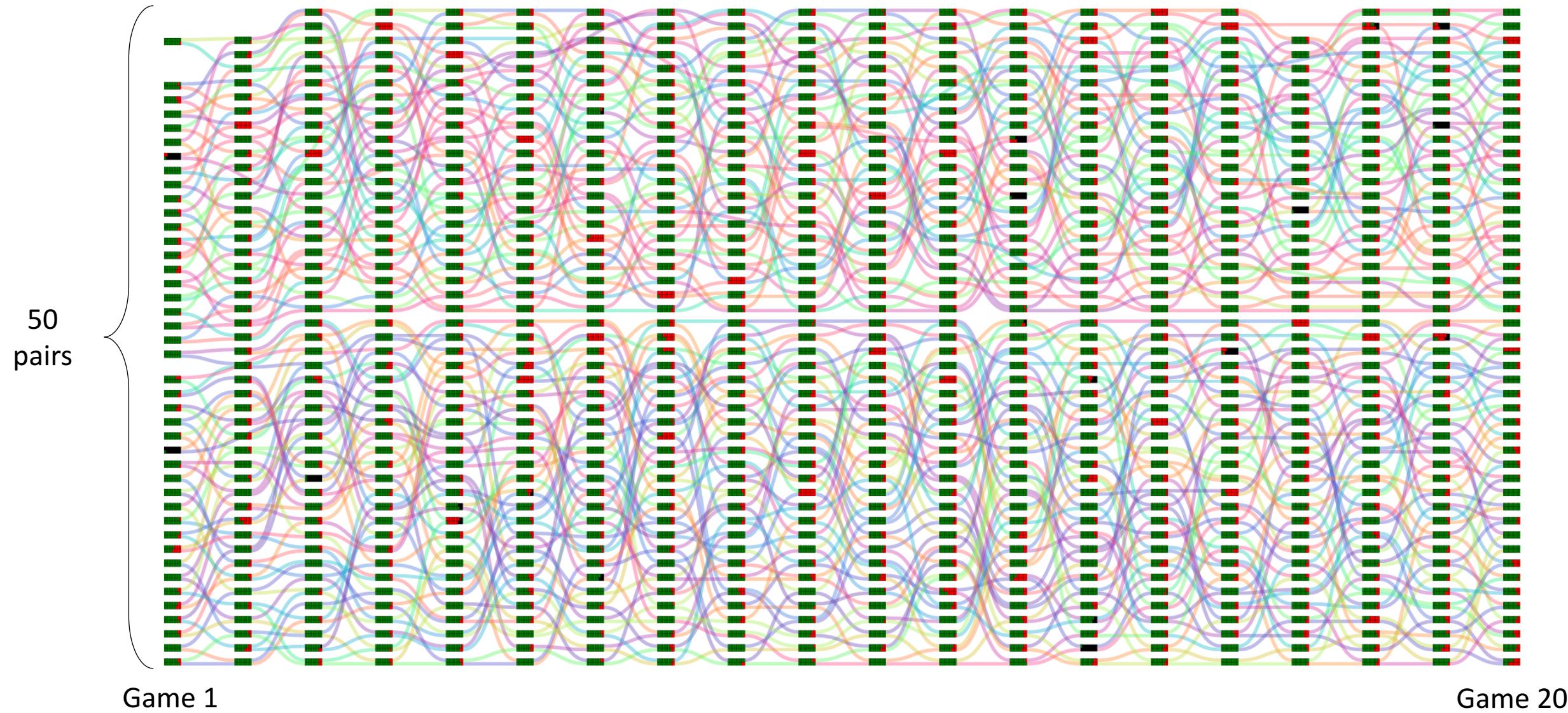


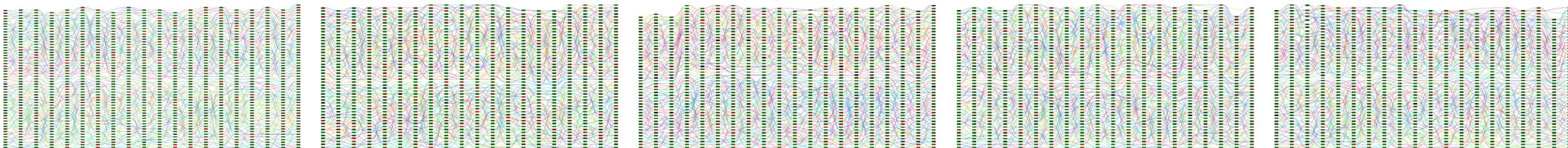
Game 1

# Random rematching across games

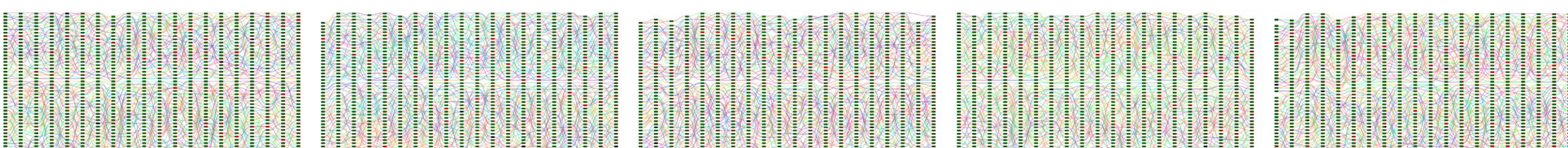
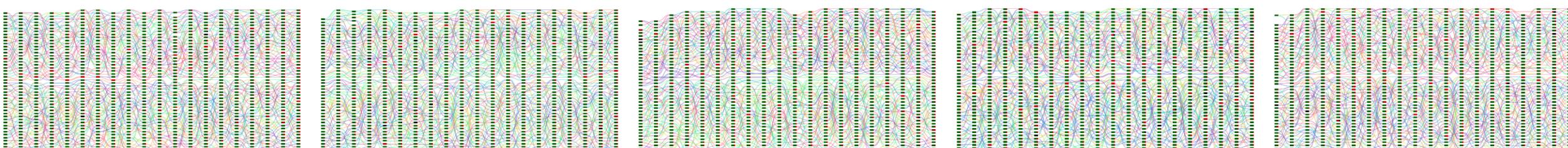
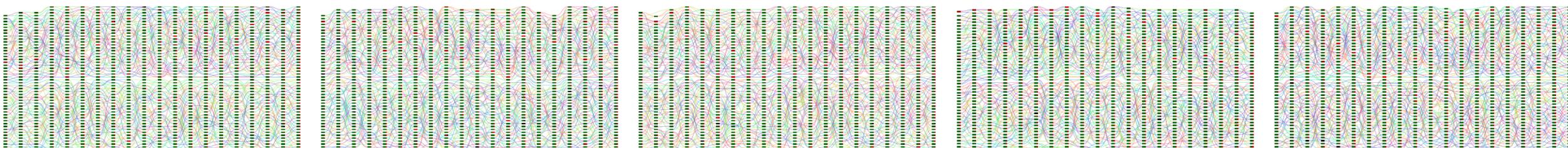


# One experiment session – 20 games





Aug 4, 2015 – Day 1



Aug 31, 2015 – Day 20

# Demo time!

You get to play prisoner's dilemma with each other!

The person with the most money at the end wins.

Navigate your laptop or mobile browser to:

<http://skygeirr.andrewmao.net:3000>

# Demo

- You will be the participants in this prisoner's dilemma experiment!
- Log in using the menu, choose “main”
- Instructions: you will be playing IPD with repeated partners
  - Loses some context, so I will explain. (wait a bit)
  - Press begin when you are ready
- Admin console; shows all user connections
- Lobby, starting a game
  - Play some rounds with a partner; timeouts and stuff
- Live experiments timer and graphs
- Disconnections, slow players
- One-way mirror of slow games – very useful to understand behavior
- Visualization of data at end

# A “superollider” for social science



**Realistic settings for collective behavior...**

**with precise instrumentation and measurement...**

**studied longitudinally over periods of time...**

**with a high degree of experimental control.**

# Closing thoughts

---

Large-scale *observational data* is useful for building *predictive models* of a *static world*

# Closing thoughts

---

But without appropriate *random variation*, it's hard to  
*predict what happens when you change something* in the  
world

# Closing thoughts

---

*Randomized experiments* are like *custom-made datasets* to answer a specific question

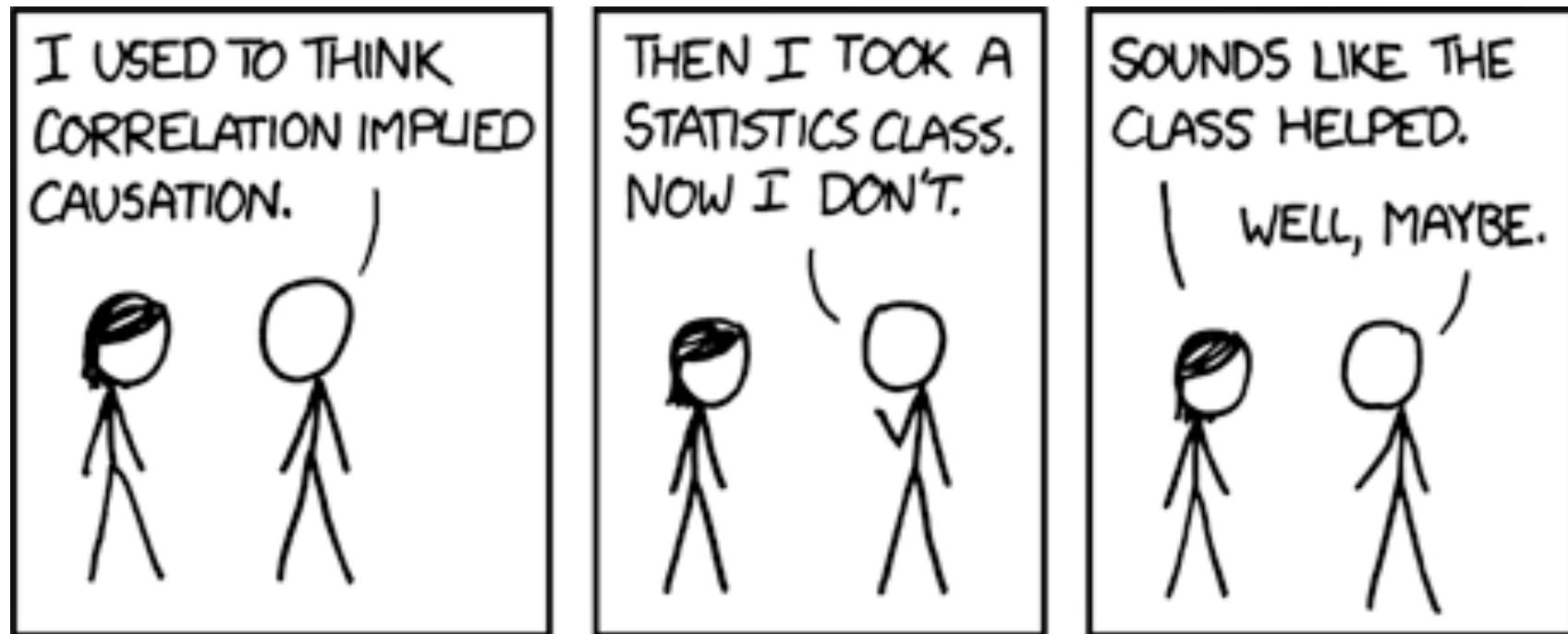
# Closing thoughts

---

*Additional data + algorithms can help us discover and analyze these examples in the wild*

# Causality is tricky!

---



“Correlation doesn’t imply causation, but it does waggle its eyebrows suggestively and gesture furtively while mouthing ‘look over there’”

<https://www.xkcd.com/552/>