

# Causality

B9145: Reliable Statistical Learning

# “What if” questions

You’ll often want to know the causal effect of some action.

You’ll want to answer “**what if**” questions:

- Will patients get better **if** they take this drug?
- Will fewer people smoke **if** we add a cigarette tax?
- Will more people buy this product **if** they see this new advertisement?
- Will I have a better outcome **if** I went to surgeon A vs surgeon B?

# Answering “what if” questions

Will patients get better if they take this drug?

- Let's say all sick patients take the drug, and many get better.  
Is this enough information?
- Let's say that we also know that if the sick patients were to not take the drug, they would not get better.
- In real life, you can't have patients take the drug and also not take the drug.  
These are **counterfactuals**

# Causal inference

Causal inference is a multidisciplinary field built across economics, epidemiology, and statistics.

Focus is on questions about **counterfactuals**

# Drug example

- There's a disease. Some people with the disease get better on their own.
- You develop a drug. You recruit people to try out the drug.
  - Anyone who wants the drug gets the drug.
- You find that a larger fraction of the recruited people, who take the drug, get better.
- However, years later, you find out the drug does not work. What could have gone wrong?

# Drug example

Perhaps the disease is debilitating, so that only people with milder symptoms were able to try out the drug, and people with milder symptoms are more likely to get better.

Or perhaps people who are more affluent are more likely to get better on their own, and also more likely to try the drug.

## **Problem: we did not observe counterfactuals**

- For each person, what would have happened if they took the drug vs didn't?

# Secret to life

**The New York Times**

## ***Another Benefit to Going to Museums? You May Live Longer***

Researchers in Britain found that people who go to museums, the theater and the opera were less likely to die in the study period than those who didn't.

**WIRED**

Technology

Science

Culture

Video

Reviews

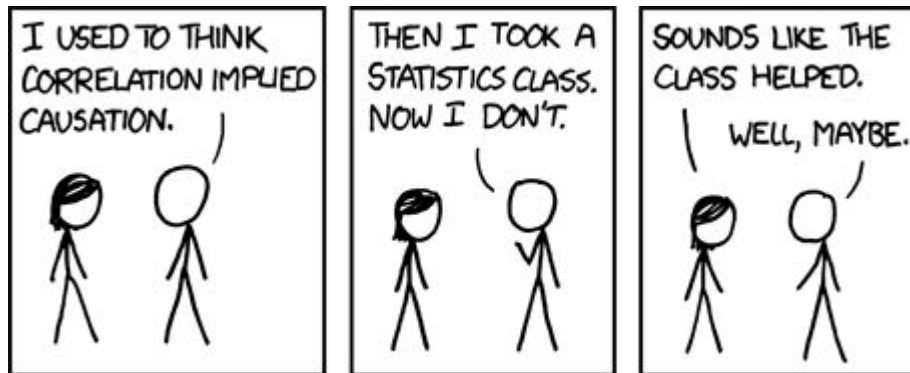
Magazine

## **Liking curly fries on Facebook reveals your high IQ**

By **PHILIPPA WARR**

12 Mar 2013





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



# Prediction and causation

Lots of ML is about prediction.

- How and when can we use prediction to estimate causal effects?
- What structure of the data do we need?

Causal inference is a multidisciplinary field built across economics, epidemiology, and statistics.

# Binary actions

- Today we will focus on the setting with two actions
  - One action represents treatment (1), the other is control (0)
- This is still foundational
  - Key difficulties still persist here despite the simplicity
  - Core technical insights will translate to more general settings
- In complex problems, this is often the standard
  - Control is status quo, treatment is the new program
  - Throughout economics, medicine, and tech, it requires a tremendous amount of domain knowledge and effort to come up with an alternative to the current system

# Potential outcomes

Also known as **bandit feedback**

Framework for explicitly modeling counterfactuals

- A: binary treatment assignment (1: treated, 0: control)
- $Y(1)$  and  $Y(0)$  are potential outcomes under treatment and control, respectively.
  - Assume we observe  $Y=Y(A)$ .
- X is observed covariates

**First goal:** Estimate average treatment effect

$$\tau := \mathbb{E}[Y(1) - Y(0)]$$

**Problem:** We only observe  $Y=Y(A)$

It's a missing data problem!

# Potential outcomes

Also known as **bandit feedback**

Framework for explicitly modeling counterfactuals

- A: binary treatment assignment (1: treated, 0: control)
- $Y(1)$  and  $Y(0)$  are potential outcomes under treatment and control, respectively.
  - Assume we observe  $\mathbf{Y}=\mathbf{Y}(\mathbf{A})$ .
- X is observed covariates



**SUTVA assumption**

**First goal:** Estimate average treatment effect

$$\tau := \mathbb{E}[Y(1) - Y(0)]$$

**Problem:** We only observe  $\mathbf{Y}=\mathbf{Y}(\mathbf{A})$

# SUTVA assumption: $Y=Y(A)$

**SUTVA** = Stable Unit Treatment Value Assumption

## 1. No interference between units

The potential outcomes for any unit do not vary with the treatments assigned to other units.

## 2. No hidden variation of treatment

For each unit, there are no different forms or versions of each treatment level, which lead to different potential outcomes.

## When might these assumptions not hold?

- Network effects: effect of vaccines on disease prevalence
- Two-sided platforms: ridesharing, ad auctions
- Equilibria: if everyone gets a job training, it won't increase everyone's income
- Different ways of administering a drug, expired vs un-expired medication

# Average treatment effect (ATE)

**First goal:** Estimate average treatment effect

$$\tau := \mathbb{E}[Y(1) - Y(0)]$$

- We only observe  $Y := Y(A)$
- What could go wrong?

Person	A	Y(0)	Y(1)	Y(1) - Y(0)
1	1	0	0	0
2	1	0	0	0
3	1	0	0	0
4	1	0	0	0
5	0	1	1	0
6	0	1	1	0
7	0	1	1	0
8	0	1	1	0

# Randomized controlled trials (RCT)

also called A/B testing, (randomized) experiments

- First try: let's randomize treatment assignments

$$Y(1), Y(0) \perp A$$

- By randomized assignments and then SUTVA, we have

$$\begin{aligned}\tau &= \mathbb{E}[Y(1) - Y(0)] = \mathbb{E}[Y(1) \mid A = 1] - \mathbb{E}[Y(0) \mid A = 0] \\ &= \mathbb{E}[Y \mid A = 1] - \mathbb{E}[Y \mid A = 0] \end{aligned}$$

← observable

- We can estimate final line from i.i.d. data  $(Y_i, A_i)$   
→ **difference in means estimator**

# Difference in means estimator



# Linear regression adjustments

Randomization might not be perfect in small samples. Can we do better?



# Linear regression adjustments

# Beyond RCTs: observational studies

- Randomization is sometimes infeasible or prohibitively expensive
  - post-market drug surveillance
  - effect of air pollution on long-term health outcomes
  - effect of a government policy on some economic outcome
- May want to use existing data collected under a previous data generating policy

# Simpson's paradox

# successful procedures / # total procedures

	actual surgery	band-aid	total
	0 / 1	99 / 100	99 / 101
	40 / 50	1 / 1	41 / 51

**Who is the better doctor?**

or, what is the effect of doctor choice on procedure success?

Y: procedure success, A: doctor, X: type of procedure

# Berkeley admissions

- Berkeley was sued for gender bias in admissions based on 1973 numbers: 44% of men were admitted but only 35% of women
- But individual department's admissions record showed no evidence of such gender-based discrimination
- Turns out women systematically applied to more competitive majors

**Y, A, X?**

# No unobserved confounding

- Previous regression-based direct method still works if there are no unobserved confounders (also called **ignorability**)

**Assumption:**  $Y(1), Y(0) \perp A \mid X$

- Observed treatment assignments are based on covariate information alone (+ random noise)
- Treatment assignment does not use information about counterfactuals
- Strong assumption, often violated in practice.
  - e.g. doctors often use unrecorded info to prescribe treatments

# No unobserved confounding: choice of $X$

How to choose  $X$  so that  $Y(1), Y(0) \perp A \mid X$  holds?

**Not this:**

- A drug

# Direct method

- By no unobserved confounding (and then SUTVA),

$$\begin{aligned}\mu_a^\star(X) &:= \mathbb{E}[Y(a) \mid X] \\ &= \mathbb{E}[Y(a) \mid X, A = a] \\ &= \mathbb{E}[Y \mid X, A = a] \leftarrow \text{observable}\end{aligned}$$

- Fit **outcome models** by loss minimization

$$\text{minimize}_{\mu_a \in \mathfrak{M}_a} \mathbb{E}[(Y(a) - \hat{\mu}_a(X))^2 \mid A = a]$$

- ATE estimator

$$\hat{\tau}_{\text{DM}} := \frac{1}{n} \sum_{i=1}^n \hat{\mu}_1(X_i) - \hat{\mu}_0(X_i)$$

- Good if outcome models are easy to learn
- Similar to RCT adjustment



# Overlap

- We need enough samples for both control and treatment throughout the covariate space (i.e. for each  $X$ )
  - You want to compare treatment and control over the same  $X$ 's
  - If treatment contains  $X$ 's that are rare or nonexistent in the control, it's hard to compare
  - Overlap governs effective sample size
- **Propensity score**  $e^*(X) := \mathbb{P}(A = 1 \mid X)$
- Assume that there exists  $\epsilon > 0$  such that  $\epsilon \leq e^*(X) \leq 1 - \epsilon$  almost surely
  - This means I have at least  $\epsilon n$  number of samples for fitting the two outcome models

# Overlap

- Overlap breaks if data is generated by a deterministic policy
  - e.g. always assign the drug (treatment) when age  $> 50$
- We need sufficient amount of randomness in treatment assignment in all covariate regions
- Often violated in practice

# Inverse propensity weighting

- What if the outcome models are very complex and difficult to estimate?
  - Direct method less good
- A natural approach is to reweight samples to correct for confounding bias
  - Essentially importance sampling
- First, estimate the propensity score
  - e.g. run logistic regression to predict A given X  $e^{\star}(X) := \mathbb{P}(A = 1 \mid X)$

$$e^\star(X) := \mathbb{P}(A = 1 \mid X)$$

# Inverse propensity weighting

- Estimator

$$\hat{\tau}_{\text{IPW}} := \frac{1}{n} \sum_{i=1}^n \left( \frac{A_i}{\hat{e}(X_i)} Y_i - \frac{1 - A_i}{1 - \hat{e}(X_i)} Y_i \right)$$

- Can work well if propensity score is simple to estimate
- But estimating this well over the entire covariate space can be difficult
  - Calibration is hard, especially in high-dimensions
- When overlap doesn't hold, importance weights blow up

# Inverse propensity weighting

$$e^{\star}(X) := \mathbb{P}(A = 1 \mid X)$$

$$\hat{\tau}_{\text{IPW}} := \frac{1}{n} \sum_{i=1}^n \left( \frac{A_i}{\hat{e}(X_i)} Y_i - \frac{1 - A_i}{1 - \hat{e}(X_i)} Y_i \right)$$

# Recap of assumptions

- **SUTVA:**  $Y=Y(A)$
- **Ignorability** / no unobserved confounding:  $Y(1), Y(0) \perp A \mid X$
- **Overlap:**  $\epsilon > 0$  such that  $\epsilon \leq e^*(X) \leq 1 - \epsilon$

# Conditional Average Treatment Effect (CATE)

**Second goal:** Estimate **conditional average treatment effect**

$$\tau(X) := \mathbb{E}[Y(1) - Y(0) \mid X]$$

- Estimate **personalized** treatment effects
  - E.g. a drug is more effective in some age groups than in others
- As before, missing data: we may not observe both  $Y(0)$  and  $Y(1)$

# Estimating heterogeneous treatment effect

There are a few ways to do it, including

- T-Learner
- S-Learner
- R-Learner



# T-learner

Separate models for treatment and control

# S-Learner

Shared feature representation, similar model class for both treatment and control

R-learner