# Controlling for Text in Causal Inference

emaadmanzoor.com

### Emaad Manzoor

### **ICWSM-22**





### 2016 - 20212021 - 2023

### emaadmanzoor.com/causal-inference-and-nlp-icwsm2022/

### Emaad Manzoor

2023 \_\_\_\_\_

# Empirically Establishing Causality

with

Unstructured Text as 





# Learning Outcomes

- In this tutorial, you will learn how to:
- 1. Formalize casual questions as causal estimands
- 2. Understand the process of causal identification
- 3. Estimate causal effects with text as a control

### What is Causal Inference?

### Using statistics & data to quantify the strength and existence of causal relationships

Causal Inference in Statistics: An Overview. Pearl (2009).

# Interdisciplinary Challenge



Vote t Cal

Click Ad





# Mask

Pay

### The Process of Causal Inference



### Causal Research Question

### Causal Estimand



### Causal Estimation Identification & Inference

### Causal Research Questions

- There are innumerable possible "effects of X on Y" questions in the world, only some are meaningful
- Meaningful causal research questions are typically motivated by social or economic theory — empirical studies are "tests of theory" (similar to physics)
- Many such tests → chances of false discovery: causal research questions need to be motivated



### Formal definition of "target of causal inference"

# Treatment a **Possible Actions** Eg. Vaccine or Placebo

### Causal Estimand





# Example Estimand: ITE

 $Y^{a=1}$ 

### Outcome had individual been vaccinated

Outcome had individual been given placebo

Ya=0

Individual Treatment Effect (ITE)

$$Y^{a=1} - Y^{a=0}$$

	$Y^{a=0}$	$Y^{a=}$
Rheia	0	1
Kronos	1	0
Demeter	0	0
Hades	0	0
Hestia	0	0
Poseidon	1	0
Hera	0	0
Zeus	0	1
Artemis	1	1
Apollo	1	0
Leto	0	1
Ares	1	1
Athena	1	1
Hephaestus	0	1
Aphrodite	0	1
Cyclope	0	1
Persephone	1	1
Hermes	1	0
Hebe	1	0
Dionysus	1	0

1

### ITEs Cannot Be "Identified"

Cannot be "expressed as a function of observable data" (without making strong assumptions)

 $Y^{a=1} = Y \text{ if vaccinated } \longrightarrow$  $Y^{a=0} = Y \text{ if placebo} \longrightarrow$ 

 $Y^{a=1}$  and  $Y^{a=0}$  not observable simultaneously

 $d \longrightarrow Y^{a=0} = ?$  $\longrightarrow Y^{a=1} = ?$ 

### **Example Estimand: ATE** N individuals i = 1, ..., N



# Example Estimand: ATE $Y_i^{A_i=1}$ $Y_i^{A_i=0}$ $\overline{Y_i^{\text{Rheia}}}$ Outcome had iOutcome had iOutcome had i

### been vaccinated

### Average Treatment Effect (ATE)

Most common causal estimand

 $Y_i^{a=0}$ 0 Poseidon been given placebo Hera Zeus Artemis Apollo Leto Ares  $\mathbb{E}[Y_{i}^{A_{i}=1}] - \mathbb{E}[Y_{i}^{A_{i}=0}]$ Athena Hephaestus Aphrodite Cyclope Persephone Hermes

Hebe

Dionysus



0

0

0

### Causal Identification

- Expressing causal estimands in terms of observable (not counterfactual) quantities
- Requires making identification assumptions - a good "identification strategy" minimizes the assumptions required
- When is ATE =  $\mathbb{E}[Y_{i}^{A_{i}=1}] \mathbb{E}[Y_{i}^{A_{i}=0}] =$  $E[Y_i | A_i = 1] - E[Y_i | A_i = 0]?$

	$A_{i}$	$Y_i$
Rheia	0	0
Kronos	0	1
Demeter	0	0
Hades	0	0
Hestia	1	0
Poseidon	1	0
Hera	1	0
Zeus	1	1
Artemis	0	1
Apollo	0	1
Leto	0	0
Ares	1	1
Athena	1	1
Hephaestus	1	1
Aphrodite	1	1
Cyclope	1	1
Persephone	1	1
Hermes	1	0
Hebe	1	0
Dionysus	1	0

### Assumption 1: $Y_i^{A_i=a} = Y_i$ if $A_i = a$

	A	Y
Rheia	0	0
Kronos	0	1
Demeter	0	0
Hades	0	0
Hestia	1	0
Poseidon	1	0
Hera	1	0
Zeus	1	1
Artemis	0	1
Apollo	0	1
Leto	0	0
Ares	1	1
Athena	1	1
Hephaestus	1	1
Aphrodite	1	1
Cyclope	1	1
Persephone	1	1
Hermes	1	0
Hebe	1	0
Dionysus	1	0

### Assumption 1: $Y_i^{A_i=a} = Y_i$ if $A_i = a$

Violated when there is more than one "version" of the treatment (for example, if  $A_i = 1$  implies vaccination by Moderna or Pfizer)

Identifying the ATE

# $Y_i^{A_i=1} = M$ for Moderna $Y_i^{A_i=1} = P$ for Pfizer



### Assumption 1: $Y_i^{A_i=a} = Y_i$ if $A_i = a$

Violated when there is "interference" or "spillover" (for example, if vaccinating individual i makes individual *i* get the vaccine)

Identifying the ATE

 $Y_i^{A_i=0} = P \text{ when } A_j = 0$  $Y_i^{A_i=0} = Q \text{ when } A_j = 1$ 



### Assumption 1: $Y_i^{A_i=a} = Y_i$ if $A_i = a$

	A	Y	$Y^0$	
Rheia	0	0	0	
Kronos	0	1	1	
Demeter	0	0	0	
Hades	0	0	0	
Hestia	1	0	?	
Poseidon	1	0	?	
Hera	1	0	?	
Zeus	1	1	?	
Artemis	0	1	1	
Apollo	0	1	1	
Leto	0	0	0	
Ares	1	1	?	
Athena	1	1	?	
Hephaestus	1	1	?	
Aphrodite	1	1	?	
Cyclope	1	1	?	
Persephone	1	1	?	
Hermes	1	0	?	
Hebe	1	0	?	
Dionysus	1	0	?	

### Assumption 1: $Y_i^{A_i=a} = Y_i \text{ if } A_i = a$ Assumption 2: $Y_i^a \perp A_i$

	A	Y	$Y^0$	
Rheia	0	0	0	
Kronos	0	1	1	
Demeter	0	0	0	
Hades	0	0	0	
Hestia	1	0	?	
Poseidon	1	0	?	
Hera	1	0	?	
Zeus	1	1	?	
Artemis	0	1	1	
Apollo	0	1	1	
Leto	0	0	0	
Ares	1	1	?	
Athena	1	1	?	
Hephaestus	1	1	?	
Aphrodite	1	1	?	
Cyclope	1	1	?	
Persephone	1	1	?	
Hermes	1	0	?	
Hebe	1	0	?	
Dionysus	1	0	?	

### Assumption 1: $Y_i^{A_i=a} = Y_i$ if $A_i = a$ Assumption 2: $Y_i^a \perp A_i$

In general,  $Y_i^a \perp A_i$  is formally assessed under assumptions using causal directed acyclic graphs and do-calculus (Pearl, 2008)



Causal effect of A on Yis identified if pathways between U and A or Uand Y are blocked



### Assumption 1: $Y_i^{A_i=a} = Y_i$ if $A_i = a$

Assumption 2:  $Y_i^a \perp A_i$ 

### Identification Proof:

 $ATE = \mathbb{E}[Y_i^{A_i=1}] - \mathbb{E}[Y_i^{A_i=0}]$ 

 $= \mathbb{E}[Y_i | A_i = 1] - \mathbb{E}[Y_i | A_i = 0] \text{ (assumption 1)}$ 





# Randomized Experiments (RCTs)

### Randomly assign individuals to treatment actions

### (and there is no attrition / selection bias)

 $Y_i^a \perp A_i$  by design if treatment is randomized

Issues with randomized experiments: Ethicality, feasibility, cost, generalizability



# Causality with Observational Data

- Strategies to argue for  $Y_i^a \perp A_i$ :
- 1. Control for <u>observed</u> confounders
- 2. Block causal pathways between unobserved confounders and treatment/outcome
- 3. Find <u>natural or quasi-experiments</u> to reduce the assumptions required





Causal effect of A on Yis identified if pathways between U and A or Uand Y are blocked





Does visible status make you more persuasive in online conversations?

Based on "Influence via Ethos: On the Persuasive Power of Reputation in Deliberation Online" (*arxiv.org/abs/2006.00707*)

### Challenger

Reputation

miguelguajiro 110 $\Delta$  Score hidden  $\cdot$  12 hours ago

By responsible, do you mean sustainable? And how do you conclude that most people believe their lives on the whole are environmentally sustainable? Could it be that people make the easy responsible choices while also aware that their lives as a whole aren't sustainable?

Reply Give Award Share Report Save

1

togtogtog 4 🖉 🔊 Score hidden · 11 hours ago

Now that is a good point. Maybe people simply don't think they are living sustainable lives and also, many people simply don't think about it one way or the other.

I guess I meant that those of us who *do* think we are living in an environmentally friendly way simply are NOT living sustainably by any means. But I wasn't very clear in how I expressed this.

Δ

### Indicator of successful persuasion



Does visible status make you more persuasive in online conversations?

Based on "Influence via Ethos: On the Persuasive Power of Reputation in Deliberation Online" (*arxiv.org/abs/2006.00707*)



Would having a theorem improve a paper's rating?

Setting: Recommender system provides a small paper list to each reviewer based on reviewer preferences and the paper text









Would having a theorem improve a paper's rating?

Treatment ( $A_j = 0,1$ ): Paper j has theorem ( $A_j = 1$ ) or not

Outcome  $(r_{ij} = 1, ..., 5)$ : Reviewer *i*'s rating for paper *j* 



Would having a theorem improve a paper's rating?

Target Estimand: ATE for each reviewer *i*, over all papers *j* 

**ATE<sub>i</sub> =**  $E[r_{ij}^{A_j=1}] - E[r_{ij}^{A_j=0}]$ 



Would having a theorem improve a paper's rating?

Is treatment assigned randomly? No. For each reviewer, some papers more likely to be recommended than others

 $E[r_{ii}^{A_j=a}] \neq$  $E[r_{ij} \mid A_j = a]$ 

Would having a theorem improve a paper's rating?

For a given reviewer, if I <u>fix</u> <u>the research topic</u>, any paper is equally likely to be recommended (<u>random</u>)

$$E[r_{ij}^{A_j=a} | \text{Topic}_j] =$$
$$E[r_{ij} | A_j = a, \text{Topic}_j]$$

### Conditional Randomization

Would having a theorem improve a paper's rating?

Since <u>each paper's topic can</u> <u>be fully inferred from its</u> <u>text</u>, I can simply control for each paper's text

$$E[r_{ij}^{A_j=a} | \text{Text}_j] =$$
$$E[r_{ij} | A_j = a, \text{Text}_j]$$

Conditional Randomization

# The Estimation Challenge

- Text is inherently unstructured, high dimensional
- Several ad-hoc ways to structure text and reduce its dimensionality: Topic modeling (LDA, NMF), document embeddings, hand-coding features
- Key issue 1: No guarantee confounders are retained
- Key issue 2: Brittle (which representation is the best?)
- Key issue 3: Inference is generally invalid



# The Estimation Challenge Common approach: Fixed g(.) (eg. topics)

Mechanics:

1. Apply g(.) to text to obtain text covariates

2. Regress Y on T and text covariates

### $Y = \theta T + g(text) + \epsilon$ Treatment Effect Fixed function

oftext





# Demo: Controlling for Words

File	Edit	View	Insert	Cell	Kernel	Widgets	Help		Trusted	Python 3 O
	•	2	↑ ↓	▶ Run	C	Markdov	vn 🗸 🖾	<u></u>		
		Trea	Itment	t Effec	t Esti	mation	3: Regre	ss $Y_i$ on	$Z_i$ and	<b>X</b> <sub>i</sub>
		Do not	run this, i	t takes to	o long, a	nd may not e	ven converge!			
I	n [43]	: # %%t # y = # X = # mod # res # res	ime simula sm.add del = sm = mode s.summar	ted_data constan .OLS(end l.fit(mo y(yname=	a[:, 0] nt(np.h dog=y, ethod=' ="Y", X	stack((sin exog=X) pinv", max name=["con	nulated_dat kiter=200) hst", "Z"])	ta[:, 1:2],	simulate	d_data[:,



# Demo: Control for Topics

File Edit	View Insert Cell Kernel Widgets Help Trusted Python 3 O
₽ + ≈ 2	Image: Second secon
	Treatment Effect Estimation 4: Regress $Y_i$ on $Z_i$ and
	Document-Topic Weights
	Play around with the number of topics.
In [77]:	<pre>%time</pre>
	<pre>ef print_top_words(model, feature_names, n_top_words): fem topig idu _topig in enumerate(model_gemmenents_);</pre>
	<pre>message = "Topic #%d: " % topic idx</pre>
	<pre>message += " ".join([feature_names[i]</pre>
	<pre>for i in topic.argsort()[:-n_top_words - 1:-1]]</pre>
	<pre>print(message)</pre>
	princ()
	umtopics = 50
	<pre>mf = NMF(n_components=numtopics).fit(tfidf_vectors)</pre>
	fidf foature names - westeriger set foature names()



Directed Acyclic Graph: Arrows represent possible causality, no arrow represents no causality

Recall: Confounder is common cause of treatment and outcome



# Can view text as 4 logical components



### Only need to somehow find and control for component *a* Needle in a haystack



### Alternative to finding this needle without using dimensionality reduction

Measure and combine correlation between text, treatment, and outcome



Measuring correlations

How well can I predict the treatment status / outcome value from the text?

### Double ML Mechanics

### Unknown g(.) jointly estimated with $\hat{\theta}$

Mechanics (Robinson, 1988):

- 1. Measure prediction errors
  - $\tilde{Y} = Y m_1(\text{text}),$
  - $\tilde{T} = T m_2(\text{text})$
- 2. Regress  $\tilde{Y}$  on  $\tilde{T}$

### $Y = \theta T + g(text) + \epsilon$ Treatment Effect function of text





# Double ML Theory

- Issue with (Robinson, 1988):  $\sqrt{n}$  consistency of  $\hat{\theta}$  requires  $m_1, m_2$  to be <u>sieve estimators</u> — poor models of text
- Double machine learning (DML):  $\sqrt{n}$  consistency possible if  $m_1, m_2$  are regularized <u>ML models</u> trained on held-out data — great for text! (must converge at  $n^{-1/4}$  or better)
- Valid asymptotic <u>confidence intervals</u>
- <u>General recipe (extensible beyond partially-linear models)</u>





### Double ML + Text Demo



### Biases Eliminated by Double ML Regularization bias and overfitting bias $Y = D\hat{\theta}_0$ $\sqrt{n}(\hat{\theta}_0 - \theta_0) = \left(\frac{1}{n}\sum D_i^2\right)^{-1}$ i∈I ı∈I :=a $+ \left(\frac{1}{n}\sum_{i\in I}D_i^2\right)^{-1}\frac{1}{\sqrt{n}}\sum_{i\in I}D_i(g_0(X_i) - \hat{g}_0(X_i))$ :=b

 $g_0$  is an ML method – b goes to 0 too slowly!

$$+ \hat{g}_0(X) + \hat{U}_i$$
$$\frac{1}{\sqrt{n}} \sum_{i \in I} D_i U_i$$

# Biases Eliminated by Double ML

Eliminating Regularization Bias: Partialling-out procedure as we saw earlier, also called orthogonalization

$$b^* = (E[V^2])^{-1} rac{1}{\sqrt{n}} \sum_{i \in I} \underbrace{(\widehat{m}_0)}_{\widehat{m}_0}$$

Bias after orthogonalization is the product of 2 errors – goes to zero more quickly (also source of  $n^{-1/4}$  convergence rate requirement)

 $(X_i) - m_0(X_i))(\hat{g}_0(X_i) - g_0(X_i))$  $\hat{g}_0$  estimation error estimation error



# Biases Eliminated by Double ML

- Eliminating Overfitting Bias: Cross-fitting
- Split sample into train and estimation subsets
- Train ML models on the train subset
- Impute prediction errors on the estimation subset
- Estimate causal effect on the estimation subset
- Repeat after flipping subsets, average estimates

# Alternative Approaches

- 1. Causal Forests: Restricted to tree-structured models, double ML permits using <u>neural networks</u>
- 2. Causal BERT, DragonNet, etc.: Do not have consistency guarantees, ways to do inference
- 3. Targeted learning / TMLE (van der Laan et al): Better finite sample properties



### Next Steps

- Preprint using double ML to control for text: On the

• Survey (preprint): Causal Inference in Natural Language Processing: Estimation, Prediction, Interpretation and Beyond

• Survey (ACL '20): Text and Causal Inference: A Review of Using Text to Remove Confounding from Causal Estimates

Persuasive Power of Reputation in Deliberation Online