Controlling for Text in Causal Inference

NLP+CSS 201, Fall 2021

emaadmanzoor.com

Emaad Manzoor

Outline

- 1. Causal inference primer
- 2. Causality from non-experimental data
- 3. Text as a control
- 4. Double machine learning
- 5. Next steps

What is Causal Inference?

Process of (i) <u>establishing</u>, and (ii) <u>quantifying</u> causal relationships empirically (using <u>statistics</u> + <u>data</u>)

Interdisciplinary Challenge



Vote t Cal

Click Ad





Mask

Pay

What is a Causal Effect? Typical Setup

Treatment a

Possible Actions Eg. Vaccine or Placebo





Causal Estimand: ITE

 $Y^{a=1}$

Outcome had individual been vaccinated

Out indiv give

Individual Treatment Effect (ITE)

		Y ~~ ~
$\mathbf{\Lambda}$	Rheia	0
Va=0	Kronos	1
	Demeter	0
	Hades	0
	Hestia	0
vutcome nad	Poseidon	1
1 1 . 1	Hera	0
dividual been	Zeus	0
	Artemis	1
tron nlacoho	Apollo	1
iven placebo	Leto	0
	Ares	1
	Athena	1
	Hephaestus	0
	Aphrodite	0
va=1 $va=0$	Cyclope	0
	Persephone	1
	Hermes	1
	Hebe	1

Dionysus

1

ITEs Cannot Be "Identified"

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

"Cannot be measured from observable data"



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

Causal Estimand: ATE

Nindividuals i = 1, ..., N



Causal Estimand: ATE

 $Y_i^{A_i=1}$

Outcome had iOubeen vaccinatedbeen

Average Treatment Effect (ATE)

Is ATE = $E[Y_i | A_i = 1] - E[Y_i | A_i = 0]$?

 $Y_{i}^{A_{i}}=0$ Outcome had *i* been given placebo

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

	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



A = 1 (Patient Admitted to ICU)

A = 0 (Patient Not Admitted to ICU)

Identifying the ATE

Y = 1 (Patient Y = 0 (Patient Died) Survived)

100

25

10 1000

Does ICU admission cause death?

Identifying the ATE E[Death | Admitted] = $E[Y_i | A_i = 1]$ 75/100 = 75%

$\mathbf{E}[\mathbf{Death} \mid \mathbf{Not} \; \mathbf{Admitted}] = E[Y_i \mid A_i = o]$

$E[Y_i | A_i = 1] - E[Y_i | A_i = 0]$

100/1000 = 10%

65%

Is the ATE = $E[Y_i^{A_i=1}] - E[Y_i^{A_i=0}]$ $= E[Y_i | A_i = 1] - E[Y_i | A_i = 0]? - \text{not always}$

$E[Y_i=a | A_i=b]$ is an Association

- Association \neq Causation under "Confounding"
 - Confounders are factors that are common causes of both the treatment the the outcome
 - Previous example: patient's age
 - Confounders are often unobserved!

Simulation of Cofounding Bias

C jupyter	NLP+CSS 2	201 - Fa	ll 2021
File Edit	View Insert	Cell	Kernel
₽ + ≈ 4		Na Run	C
	Treatment Since the estimateffect will be bias	tion does	t Estir
In [76]:	<pre>%%time y = simulate X = sm.add_co model = sm.00 res = model.? res.summary()</pre>	d_data[; onstant) LS(endog fit(meth yname="}	:, 0] (simulat g=y, exc hod="pin ", xnam
Out[76]:	CPU times: us Wall time: 49	ser 428 5.8 ms	ms, sys
	OLS Regression Re Dep. Variable	sults	Y
	Model	:	OLS



Adj. R-squared: 0.071

The Magic of Randomized Experiments (RCTs)

$$E[Y_i^{A_i=a}] = E[Y_i | A_i=a]$$

Issues: Expensive, infeasible, unethical (eg. randomly send patients to ICU)

The only <u>perfect</u> solution to confounding: Randomly assign individuals to treatment actions

if $Y_i^a \perp A_i$ (randomization)

Dealing with Observational Data

For example, measure each patient's age, then compute ATE within each age bucket, average all the bucket-specific ATEs to get the overall ATE

Confounders are often unobserved!

Control for / condition on observed confounders

Dealing with Observational Data

Example: Effect of Airbnb certification on booking rates

Popular methods: Regression discontinuity designs, difference-in-differences, instrumental variables

Use controls + natural /quasi-experiment

Controls + Randomization

$E[Y_{i}^{A_{i}=a} | X_{i}] = E[Y_{i} | A_{i}=a, X_{i}] \text{ if } Y_{i}^{a} \perp A_{i} | X_{i}$

X_i is a scalar or vector of controls for individual i

Improve precision (smaller confidence intervals), deal with conditional randomization, endogenous selection bias, etc.

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*

 $\mathbf{ATE}_{i} = E[r_{::}^{A_{j}=1}] - E[r_{::}^{A_{j}=0}]$ IJ

- 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

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

Would having a theorem improve a paper's rating? For a given reviewer, if I <u>fix</u>

the research topic, 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

- More examples:
- Does having higher reputation on a debating website make you more persuasive? (https://arxiv.org/abs/ 2006.00707) Instrumental variable (quasi-experiment) for reputation — need to control for argument text

• Text and Causal Inference: A Review of Using Text to Remove Confounding from Causal Estimates (ACL 2020)

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

First Attempt: Control for Words

Controlling for Text (autosaved)		e e	Logout
Widgets Help	Trusted	P	ython 3 O
Markdown			

Treatment Effect Estimation 3: Regress Y_i on Z_i and X_i

```
# X = sm.add constant(np.hstack((simulated data[:, 1:2], simulated data[:,
```


Second Attempt: Control for Topics

C jupyter	NLP+CSS 201 - Fall 2021 - Controlling for Text (autosaved)	
File Edit	View Insert Cell Kernel Widgets Help Trusted Python 3 O	
E + % (A IN Run ■ C IN Markdown	
	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 ef print_top_words(model, feature_names, n_top_words): for topic_idx, topic in enumerate(model.components_): message = "Topic #%d: " % topic_idx message += " ".join([feature_names[i]</pre>	
	<pre>umtopics = 50 mf = NMF(n_components=numtopics).fit(tfidf_vectors) fidf_feature_names = vectorizer.get_feature_names() rint_top_words(nmf, tfidf_feature_names, 10)</pre>	

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?

- General "recipe" to perform statistically valid estimation/inference after incorporating ML models
- Neutralizes "regularization bias" and "overfitting bias" that arise from ML model estimation
- Fast $O(\sqrt{n})$ convergence rates despite slowlyconverging nonparametric ML models

Partially Linear Regression Model

$r_{ii} = \theta_0 + \theta_1 A_j + f(\mathbf{text}_j) + \epsilon_i$ function

Special moment condition combines correlations between text, treatment, and outcome

$\begin{array}{ll} \underbrace{\mathsf{Outcome}}_{i} & \operatorname{Treatment}_{i} \\ \mathbb{E}[(r_{ij} - \mathbb{E}[r_{ij} | \operatorname{text}_{j}] - \beta_{1}(A_{j} - \mathbb{E}[A_{j} | \operatorname{text}_{j}])) \times \\ & (A_{j} - \mathbb{E}[A_{j} | \operatorname{text}_{j}])] = 0 \end{array}$

- Empirically solve for the coefficients 2.

Double ML Procedure

1. Construct a Neyman-orthogonal moment condition for the regression equation

In practice: Can be done using a sequence of ordinary-least-squares regressions

- Compute the prediction errors of treatment from text, outcome from text
- 2. Regress the outcome prediction error on the treatment prediction error
- Text prediction models must be trained on a <u>held out</u> subset of the data (called sample-splitting)

Double ML Theoretical Details

Check out Chris Felton's slides: <u>https://</u> <u>scholar.princeton.edu/sites/default/files/bstewart/files/</u> <u>chern.handout.pdf</u>

Double ML + Text Examples

💭 Jupyter	NLP+CSS 201 - Fall 2021 - Controlling for Text (autosaved)
File Edit	View Insert Cell Kernel Widgets Help Trusted Python 3 O
₽ + ≈ 2	A ↓ NRun ■ C ▶ Markdown ↓ ■
	Treatment Effect Estimation 5: Regress Y_i on Z_i and X_i
	using Double Machine Learning
	Play around with the type of ML models used to predict the treatment and outcome.
	Note that some models take really long to train (eg. Random Forests).
	Using the EconML package
In [8]:	<pre>%%time W = simulated details 01 wavel() # out some</pre>
	Y = Simulated_data[:, 0].ravel() # Outcome
	W = simulated data[:, 3:] # text
	<pre>dml_mult = LinearDMLCateEstimator(model_y=LassoCV(cv=2, n_alphas=1, verbos)</pre>
	model t=LassoCV(cv=2, n alphas=1, verbos

linear_first_stages=True, n_splits=2)

What We've Learned Today

- 1. ITEs, ATEs, randomized experiments
- 2. Observational data and confounding bias
- 3. Double machine learning

4. Controlling for text with double machine learning

Alternative Approaches

- 1. Causal Forests: Restricted to tree-structured
- 2. Causal BERT, DragonNet, etc.: Do not have
- colleagues): Approach worth exploring

models, double ML permits using <u>neural networks</u>

consistency guarantees, ways to do inference

3. Targeted learning / TMLE (van der Laan and

Next Steps

- and Mostly Harmless Econometrics

• 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

• Preprint using double ML to control for text: On the Persuasive Power of Reputation in Deliberation Online

• Recommended Books: Causal Inference: What If (free)