Causal Representation Learning with Generative Artificial Intelligence: Application to Texts as Treatments

Kosuke Imai

Harvard University

Oxford Spring School in Advanced Research Methods University of Oxford April 3, 2025

Joint work with Kentaro Nakamura (Harvard Kennedy School)

Motivation

- Generative AI is transforming medicine, education, marketing, etc.
- Can methodologists get some help from generative AI too?

How can a methodologist use LLM to help improve research? Give me a short answer.

\$

A methodologist can use large language models (LLMs) to enhance research by automating literature reviews, generating hypotheses, designing study frameworks, and analyzing data. LLMs can also assist in drafting and refining research papers, coding qualitative data, and even suggesting relevant methodologies or statistical approaches, making the research process faster and more rigorous.

Using Generative AI to Improve Causal Inference

- LLM-assisted causal inference with unstructured data
 - use LLM to generate treatment texts
 - andomly assign generated texts to survey respondents and measure their reactions
 - extract the internal representation of generated texts from LLM
 - Machine learning using the extracted true vector representation
 - estimate a deconfounder that summarizes all confounding information
 - double machine learning for valid statistical inference
- Advantage: no need to estimate the representation
 - nonparametric modeling: no topic modeling
 - computational efficiency: no fine tuning
 - better empirical performance
- Extension: instrumental variable approach to perceived treatments

Motivating Application: Texts-as-Treatments

- Candidate Biography Experiment (Fong and Grimmer, 2016)
 - 1246 biographies of American politicians scraped from Wikipedia
 - 1,886 voters as respondents
 - randomly assign biographies to voters
 - feeling thermometer as the outcome
- Analysis
 - supervised topic model to discover 10 treatment features
 - estimate the average treatment effects of estimated topic proportions
- Existing approaches for texts-as-treatments:
 - Model-based approach (e.g., Egami et al. 2022; Fong and Grimmer, 2023)
 - causal representation learning based on embedding (e.g., Veitch et al. 2020; Pryzant et al. 2021; Gui and Veitch, 2023)

Using LLM to Generate Biographies (Llama 3)

Prompts	Generated candidate biography
(System): You are a text	Nathaniel C. Gilchrist was an American
generator who always pro-	politician who served as a member of the
duces a biography of the	United States House of Representatives from
instructed person.	Maryland's 4th congressional district from
(User): Create a biogra-	1895 to 1899. Born on November 24, 1847,
phy of an American politi-	in Baltimore, Maryland, Gilchrist was the son
cian named Nathaniel C.	of a prominent family of merchants and ship-
Gilchrist	builders. He was educated at the Baltimore
	City College and later attended the Univer-
	sity of Maryland, where he studied law. After
	graduating, he was admitted to the bar and
	began practicing law in Baltimore.

• Requirements:

- open-source LLM: GPT, PaLM, Llama, OPT, etc.
- **2** assume LLM is compliant with prompt or use human coder

Text Reuse: Using the Existing Texts as Treatments

Prompts	Generated candidate biography		
(System): You are a text generator	Higgins was born in		
who just repeats the input text.	Red Lion Hundred in New Castle		
(User): Higgins was born in	County, Delaware. He attended		
Red Lion Hundred in New Castle	Newark Academy and Delaware Col-		
County, Delaware. He attended	lege, and graduated from Yale Col-		
Newark Academy and Delaware Col-	lege in 1861, where he was a mem-		
lege, and graduated from Yale Col-	ber of Skull and Bones. After study-		
lege in 1861, where he was a mem-	ing law at the Harvard Law School,		
ber of Skull and Bones. After study-	he was admitted to the bar in 1864		
ing law at the Harvard Law School,	and began practice in Wilmington,		
he was admitted to the bar in 1864	Delaware. He also served for a time		
and began practice in Wilmington,	in the United States Army in 1864.		
Delaware. He also served for a time			
in the United States Army in 1864.			



• Deep Generative Model

 $\mathbb{P}(\boldsymbol{X}_i \mid \boldsymbol{h}_{\boldsymbol{\gamma}}(\boldsymbol{R}_i)), \ \mathbb{P}(\boldsymbol{R}_i \mid \boldsymbol{P}_i).$

- **P**_i: prompt
- X_i: treatment object (generated texts)
- **R**_i: hidden states or internal representations
- $h_{\gamma}(R_i)$: deterministic function from hidden states to last layer
- Other variables
 - $Y_i(\mathbf{x})$: Potential outcome when exposed to treatment object \mathbf{x}
 - Y_i: Outcome (collected from the survey respondents)
 - *T_i*: Binary treatment feature (e.g., military experiences)
 - **U**_i: Confounding features (e.g., college education)

Assumptions



• Separability:

$$Y_i(\boldsymbol{X}_i) = Y_i(g_T(\boldsymbol{X}_i), g_{\boldsymbol{U}}(\boldsymbol{X}_i)) = Y_i(T_i, \boldsymbol{U}_i)$$

• Lemma: separability implies overlap

$$\mathbb{P}(T_i = t \mid \boldsymbol{U}_i = \boldsymbol{u}) > 0.$$

• Deterministic decording: $\mathbb{P}(X_i | h_{\gamma}(R_i))$ is degenerate

Nonparametric Identification

• Average treatment effect (ATE):

$$au := \mathbb{E}[Y_i(1, \boldsymbol{U}_i) - Y_i(0, \boldsymbol{U}_i)]$$

• Under these assumptions, there exists a Deconfounder $f : \mathbb{R}^r \to \mathbb{R}^q$ with $q \leq r$ such that

$$Y_i \perp \mathbf{R}_i \mid T_i = t, f(\mathbf{R}_i), \quad t \in \{0, 1\}$$

- Deconfounder does not have to be unique
- Example: Confounding Features **U**_i (deterministic function of **R**_i)
- By adjusting for this Deconfounder, we can identify the ATE
- Direct adjustment for R_i leads to the lack of overlap

Estimation and Inference



Estimate the outcome models and deconfounder via TarNet (Shalit et al. 2017):

$$\{\hat{\boldsymbol{\lambda}}, \hat{\boldsymbol{\theta}}_0, \hat{\boldsymbol{\theta}}_1\} = \operatorname*{argmin}_{\boldsymbol{\lambda}, \boldsymbol{\theta}_0, \boldsymbol{\theta}_1} \frac{1}{n} \sum_{i=1}^n \{Y_i - \mu_{T_i}(\boldsymbol{f}(\boldsymbol{R}_i; \boldsymbol{\lambda}); \boldsymbol{\theta}_{T_i})\}^2$$

estimate the propensity score using the estimated Deconfounder

$$\pi(\boldsymbol{f}(\boldsymbol{R}_i, \hat{\boldsymbol{\lambda}})) = \mathbb{P}(T_i = 1 \mid \boldsymbol{f}(\boldsymbol{R}_i, \hat{\boldsymbol{\lambda}}))$$

Popular DragonNet (Shi et al. 2019) jointly estimates the outcome models, propensity score, and deconfounder, leading to the lack of overlap

Double Machine Learning (Chernozhukov et al. 2018)

- Cross-fitting:
 - **1** randomly divide the data into *K* folds
 - 2 for each k = 1, ..., K, use the *k*th fold as the test set and the remaining k 1 folds as the training set
 - randomly split the training set further into two subsets
 - ② use the first subset to estimate outcome models and deconfounder
 - estimate by the second subset to estimate propensity score given the estimated deconfounder
 - Ompute the ATE estimator as:

$$\begin{aligned} \hat{\tau} &= \frac{1}{nK} \sum_{k=1}^{K} \sum_{i:l(i)=k} \hat{\mu}_{1}^{(-k)}(\hat{f}^{(-k)}(R_{i})) - \hat{\mu}_{0}^{(-k)}(\hat{f}^{(-k)}(R_{i})) \\ &+ \frac{T_{i}\{Y_{i} - \hat{\mu}_{1}^{(-k)}(\hat{f}^{(-k)}(R_{i}))\}}{\hat{\pi}^{(-k)}(\hat{f}^{(-k)}(R_{i}))} - \frac{(1 - T_{i})\{Y_{i} - \hat{\mu}_{0}^{(-k)}(\hat{f}^{(-k)}(R_{i}))\}}{1 - \hat{\pi}^{(-k)}(\hat{f}^{(-k)}(R_{i}))} \end{aligned}$$

• Double robustness, asymptotic normality

Practical Implementation Details

• Internal representation extracted from LLM is still high-dimensional:

dim(\mathbf{R}) = number of tokens × 4096 for Llama 3 (8 billion parameters)

- Pooling strategies depend on deep generative models
 - BERT: the first special classification token [CLS]
 - Llama 3: the hidden states of the last token
- TarNet requires hyperparameter tuning
 - size and depth of layers
 - learning rate
 - maximum epoch size
- Use of automatic hyperparameter optimization methods (e.g., Optuna)

Simulation Study Setup

- A simulation based on the candidate biography experiment
 - Create 4,000 sets of the first, middle, and last names of political candidates via randomly sampling from the Fong and Grimmer data
 - Use Llama 3 to generate a biography for each US political candidate's
 - Instruct LLM to repeat the same texts for reuse
- The data generating process:

$$egin{aligned} Y_i &= lpha_1 \mathcal{T}_i + lpha_2 \mathcal{T}_i h_1(oldsymbol{X}_i) - lpha_3 h_1(oldsymbol{X}_i) - lpha_4 h_2(oldsymbol{X}_i) + \epsilon_i \ \epsilon_i &\sim \mathcal{N}(\mu_i, 1) \end{aligned}$$

where

- T_i: military background (binary)
- $h_1(X_i)$: topic-model based confounder
- $h_2(\mathbf{X}_i)$: sentiment-analysis based confounder
- $2 \times 3 = 6$ scenarios:
 - separability holds or does not hold (separate or overlapping topics)
 weak, medium, or strong confounding

Simulation Results



- Proposed (new) - Proposed (reuse) - T-Learner with BERT · · DR-Learner with BERT

Distribution of Estimated Propensity Score



- Treatment Group

Performance across Different Sample Sizes



Empirical Analysis I: Candidate Biography Experiment

• Analyze the original survey by Fong and Grimmer (2016)

- 1,246 Congressional candidate biographies from Wikipedia
- 1,886 survey participants with a total of 5,291 observations
- evaluate a biography using the feeling thermometer [0, 100]
- Keyword-based treatment coding: "military", "war", "veteran", or "army"
- use text-reuse approach with Llama 3

Methods	ATE	95% Conf. Int.	Runtime (sec.)
Proposed method (reuse)	5.462	[2.790, 8.135]	28.9
T-learner with BERT	-2.557	[-2.608, -2.505]	6139.7
DR-learner with BERT	-67.777	[-109.967, -25.587]	6210.3

Empirical Analysis II: Chinese Internet Sensorship

- Texts as confounders (Roberts *et al.* 2020)
- Question: Are Chinese social media users who had their post censored more likely to be censored for subsequent posts?
 - Treatment: whether or not a post was censored (at least once in the first half of 2012)
 - Outcomes: censorship during four weeks after a censored post
 - number of posts
 - 2 proportion of censored posts
 - I proportion of removed accounts
 - text-as-confounder: contents of posts
- Original analysis: Matching (CEM) with topic proportions (STM) and propensity score (inverse regression)
- Our reanalysis:
 - Text reuse with Llama 3
 - Apply the deconfounder method:
 - entire sample (4155 users; 75324 posts)
 - 2 matched samlpe (628 users; 879 posts)

ATE Estimates



• Proportion of top 30 censorship related keywords from Fu et al. (2013)

Outcome	Deconfouder	TIRM (Roberts et al.)
Number of Posts	-0.010	0.005
	[-0.076, 0.056]	[-0.072, 0.060]
Rate of censorship	-0.023	0.126
	[-0.089, 0.043]	[0.060, 0.190]
Rate of missing posts	-0.022	0.024
	[-0.089, 0.043]	[-0.042, 0.090]

Concluding Remarks

• Generative AI can be used to improve causal inference

- generate treatments at scale
- enables the extraction of true internal representation
- better causal representation learning
- Open-source software GPI (GenAl Powered Inference) available at https://gpi-pack.github.io/

Further extensions

- images and videos
- tabular data too
- interpretation of estimated deconfounder
- discovery of treatment concepts
- policy learning with unstructured treatments