

# Conceptualizing Treatment Leakage in Text-based Causal Inference

Adel Daoud<sup>1,3,4</sup> Connor T. Jerzak<sup>1,3</sup> Richard Johansson<sup>2,4</sup>

<sup>1</sup>Linköping University, <sup>2</sup>University of Gothenburg,

<sup>3</sup>Harvard University, <sup>4</sup>Chalmers University of Technology

{adel.daoud, connor.jerzak}@liu.se, richard.johansson@gu.se

## Abstract

Causal inference methods that control for text-based confounders are becoming increasingly important in the social sciences and other disciplines where text is readily available. However, these methods rely on a critical assumption that there is no *treatment leakage*: that is, the text only contains information about the confounder and no information about treatment assignment. When this assumption does not hold, methods that control for text to adjust for confounders face the problem of post-treatment (collider) bias. However, the assumption that there is no treatment leakage may be unrealistic in real-world situations involving text, as human language is rich and flexible. Language appearing in a public policy document or health records may refer to the future and the past simultaneously, and thereby reveal information about the treatment assignment.

In this article, we define the treatment-leakage problem, and discuss the identification as well as the estimation challenges it raises. Second, we delineate the conditions under which leakage can be addressed by removing the treatment-related signal from the text in a pre-processing step we define as *text distillation*. Lastly, using simulation, we show how treatment leakage introduces a bias in estimates of the average treatment effect (ATE) and how text distillation can mitigate this bias.

## 1 Introduction

In observational settings, scholars need to collect information about potential confounders in order to estimate the causal effect ( $\tau$ ) of a treatment on an outcome (Daoud and Dubhashi, 2020). If we observed the set of confounders directly, we could condition on those quantities to recover unbiased causal effects. Yet, because some confounders  $U$  are difficult to measure directly, scholars are turning to alternative data sources, such as medical records, policy documents, or social media posts,

to indirectly measure (proxy) confounders (Kino et al., 2021). Recent methodological frameworks supply ways of integrating high-dimensional text data into causal estimation (Mozer et al., 2020; Roberts et al., 2020; Feder et al., 2021).

However, prior literature has primarily assumed that documents only contain information about the confounder, but not about the treatment—something we term the *no-treatment-leakage assumption*. Here, “contain information” means that the text is caused by the treatment (or the confounder) directly or indirectly. When treatment leakage occurs after treatment assignment, its bias is equivalent to a post-treatment bias (Pearl, 2015).

Treatment leakage leads to an identification challenge. The challenge is that  $W$  is both necessary for adjusting (as it is a proxy) yet it is also a post-treatment variable. Without treatment leakage,  $W$  would not be a post-treatment variable, as it does not harbour information about the treatment assignment. But because of leakage, scholars would have to accept bias arising from either adjusting on a post-treatment variable (arising from the part of  $W$  influenced by the treatment) or bias arising from not adjusting for unobserved confounding. Although several methodological studies develop and adapt causal-inference methods for text data (Keith et al., 2020), almost no studies examine the biasing influence of treatment leakage and how to counter this bias.

Our work investigates the treatment-leakage challenge. It shows that if  $W$  is the only available text representing  $U$  and there exists a distillation method,  $f$ , that has the ability to transform (e.g. partition)  $W$  into its post-treatment  $W_T$  and proxy textual-components  $W_U$ , then adjusting on  $W_U$  is the best one can do in identifying  $\tau$ . As  $W_U$  is not post-treatment, we can adjust for it to reduce the bias when estimating  $\tau$ . These  $f$  functions can represent a human annotator, identifying and removing parts of text (e.g., words, sentences) that

belong to  $\mathbf{W}_T$  and curating  $\mathbf{W}_U$ ; or, under additional assumptions,  $f$  can be based on supervised or unsupervised machine learning models that transform the text or its representation (Åkerström et al., 2019; Feder et al., 2021).

In this paper, we define key assumptions and demonstrate the mechanics of text distillation in a simulated experiment. Using a language model, we generate synthetic documents  $\mathbf{W}$  so that they contain information about the treatment assignment,  $T$ , and the unobserved confounding,  $U$ , imprinted paragraph by paragraph. Because we control which paragraph is affected by  $T$  (injecting post-treatment bias) or by  $U$  (infusing knowledge about the confounder), we have an oracle distillation function,  $f$ , that mimics human coding. This oracle method perfectly distills  $\mathbf{W}$ , and supplies  $\mathbf{W}_U$ . Then, when using  $\mathbf{W}_U$  in our causal model, we reduce bias of  $\hat{\tau}$  markedly. Although our oracle is idealized, it deepens intuition, and in future work, we will investigate the conditions under which automated methods can be applied to obtain  $f$ .

By conceptualizing the problem of treatment leakage in text data and investigating its impact, scholars developing causal methods can be better positioned to tailor their frameworks to reduce bias; domain scholars can better calibrate their data collection procedure to account for this leakage.

## 2 Treatment Leakage in Text Data

While the literature on dealing with confounding in observational studies is established (Rubin, 1974), recent advances have been made in the analysis of text-based causal inference. Indeed, text  $\mathbf{W}$  is widely available in the health and social sciences (Gentzkow et al., 2019; Kino et al., 2021), and can be used to proxy for some confounders,  $U$ , that would otherwise remain unobserved (Keith et al., 2020). If the text only contains information about  $U$  and no other factors, then  $\mathbf{W}$  is a faithful representation of  $U$  and we denote it as  $\mathbf{W}_U$ . However, text, by its nature as a medium of creativity, rarely has fixed boundaries, and can contain information not only about confounders, but also leak information about the treatment assignment and its effects.

The future- and backward-looking nature of text can exacerbate treatment leakage. Documents that often contain backward looking temporally (e.g. in much of journalism) or has an unknown production date, will like contain information about the treatment and its effects. Using these docu-

ments directly for causal inference would inject post-treatment bias. Conversely, documents that reference the future (e.g., many public-policy documents in the economy and polity) may also lead to unfavorable RMSE if they predict the future well (see §2.1.2). As a result, a substantial amount of real-world text containing rich information about confounding factors might be affected by that language can reference the future, post-treatment state.

### 2.1 Characterizing Treatment Leakage

We define *treatment leakage* as when the text,  $\mathbf{W}$ , is affected by treatment status,  $T$ : that is,  $\mathbf{W}$  is conditionally dependent on  $T$  given  $U$ .

$$\textit{Treatment leakage: } \mathbf{W} \not\perp T | U$$

The treatment leakage can take different forms. In the most straightforward case, we can assume that a portion  $\mathbf{W}_T$  is affected by  $T$  while another portion  $\mathbf{W}_U$  is affected by  $U$ . However, in the general case it may be difficult to partition the document into treatment- and confounder-related passages, and we should see  $T$  and  $U$  as latent factors controlling the data-generating process. For instance,  $T$  may affect the overall tone or sentiment of a document.

We can quantify the degree of treatment leakage in different ways. If the text can be partitioned into treatment- and confounder-related passages  $\mathbf{W}_T$  and  $\mathbf{W}_U$  as described above, we can consider the fraction  $\frac{|\mathbf{W}_T|}{|\mathbf{W}|}$  to be a measure of the degree of treatment leakage; this also assumes that each partition carries strength equal to the number of its elements (e.g., words) and each element has the same strengths. In the general case, we may turn to information-theoretical quantities, for instance the conditional mutual information between  $\mathbf{W}$  and  $T$  given  $U$ .

In the following, we discuss a number of situations in which treatment leakage can occur.

#### 2.1.1 Case 1: Text is Post-treatment

In one form of this phenomenon, there is a causal relationship between the treatment status  $T$  and the text. Figure 1, panel *a.*, shows a directed acyclic graph (DAG) representing this scenario where the text affected by the treatment status. This sort of treatment leakage induces post-treatment bias: when the text is affected by the treatment, conditioning on the text (which is a collider) opens the path from  $T$  to  $Y$  through  $\mathbf{W}$  and  $U$ , will in general yield biased estimates (in the notation of Pearl (2015),  $(Y \not\perp T | \mathbf{W})_{G_T}$ ).

Identification assumptions may also be hard to maintain, with the treated/control units having distinct text features (e.g. if all treated units have associated texts referring to the treatment). This lack of overlap would violate the identification assumptions of causal estimators such as Inverse Propensity Score Weighting (IPW) (Heinrich et al., 2010), and could lead to extreme estimated probabilities, something we see empirically in Figure 3.

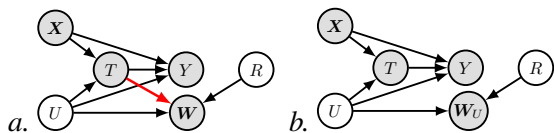


Figure 1: A causal model consisting of observed variables (shaded): confounders ( $X$ ), treatment ( $T$ ), outcome ( $Y$ ), document ( $W$ ), and unobserved variables (unshaded): confounder ( $U$ ) and residual factors ( $R$ ). The red-colored edge in *a.* represents the treatment leakage. In *b.*, A distillation function  $f$  has removed the treatment information in the text, leaving only information from the confounder. A perfect intervention of  $f$  is equivalent with deleting the red arrow; a less than perfect intervention reduces at least its dependence.

### 2.1.2 Other Cases

Figure 1 shows a case when text is post-treatment, but in other cases the precise DAG structure may not be known. For example, text may represent a mediator if the document includes post-treatment information and also affected the outcome (if, for example, the text is congressional speech and the outcome is a roll call vote). If the proxy text is pre-treatment and directly affects the treatment, conditioning on the treatment-related portion of the text could increase the variance of estimation, leading to unfavorable RMSE (Myers et al., 2011).

## 3 Text Distillation as Preprocessing

*Text distillation* is a form of text preprocessing. It has to target any text (e.g., tone, words, sentences) that belongs to  $W_T$ , and remove it from  $W$ . Thus, distillation ensures that the treatment signal is negated. As Figure 1, panel *b.* shows, if distillation is perfectly successful, it results in cutting the red arrow (from  $T$  to  $W$ ). The arrow is cut, because the distillation function has removed  $W_T$  from  $W$ , supplying  $W_U$  for causal analysis.

### 3.1 Assumptions for Valid Distillation

Depending on how the treatment leakage is manifested in  $W$ , we need to introduce assumptions

to make distillation feasible. As already discussed in §2.1, in some cases we may assume that  $W$  contains treatment-related passages  $W_T$  and confounder-related passages  $W_U$ . We may further assume that the text is *separable*: that is,  $W_T$  and  $W_U$  do not overlap.

$$\text{Separability Assumption: } W_U \cap W_T = \emptyset$$

Assuming separability, a perfect distillator will produce  $W^* = f(W)$  that is equivalent to the confounder-related portion of the text,  $W_U$ . *Perfect* distillation means that the distillator  $f$  identified text that contains the same information about  $U$  as  $W_U$  has. Thus, if  $W_U$  is a valid adjustment set, then  $W^*$  is that as well. The separability assumption is appealing because it implies that researchers only need to find a valid partition of the text (and do not need to consider all possible text transformations).

This separability assumption is particularly plausible for text data, which by its nature consists of a sequence of linguistic signifiers which can be decomposed into smaller units (e.g. paragraphs).

While plausible for many circumstances, in some cases separability may not hold, as when the entire tone of the text is affected by the treatment. In this more complicated setting, we need a more general assumption, that the transformed text,  $W^*$ , is conditionally independent of  $T$  given  $U$ . That is, the conditional mutual information between  $W^*$  and  $T$  given  $U$  is zero, while information about  $U$  in  $W^*$  is maintained. Despite the benefits of this more general framing, because  $U$  is unobserved, it may be difficult for investigators to assess whether the assumption is satisfied or whether ethically problematic information has been included in the  $f$  function (e.g., race; Menon and Williamson (2018)). Unlike numerical data, as text data is readable, scholars can examine and validate whether  $W^*$  still contains information about  $T$ .

## 4 Experimental Setup

We use simulation to illustrate the dynamics of text distillation and build on the framework for evaluating text-based causal inference methods introduced by Wood-Doughty et al. (2021). We generate numerical covariates from the model in Figure 1; the general procedure is described in §A, with implementation details in §B. Parameters are selected so that ATE estimates  $\hat{\tau}$  are biased if the estimator does not account for the unobserved confounder  $U$ .

Following Wood-Doughty et al. (2021), we generate documents,  $\mathbf{W}$ , by sampling from an English-language GPT-2 model (Radford et al., 2019). In contrast to their approach, text generation is conditioned not only on  $U$  but also on  $T$ . As described in detail in §A, we define paragraph-level topics, where some topics are associated with  $U$ , some with  $T$ , and some with a residual topic related only to other background variables ( $R$  in Figure 1). For a given paragraph topic, we define a number of prompts and a distribution shift that increases the probability of generating topic-related keywords.

As we simulate and record which paragraphs are affected by  $T$  and by  $U$ , our distillator  $f$  has oracle properties. We can then use  $f$  to investigate three idealized distillation scenarios. The first is when a distillator was not applied or the distillator failed to do any distillation  $f(\mathbf{W}) = \mathbf{W}$ . It outputs the same corpus. The second is when it perfectly distills  $\mathbf{W}$ , excluding all paragraphs affected by  $T$ . That is, apply  $f(\mathbf{W}) = \mathbf{W}^*$  such that  $\mathbf{W}^* = \mathbf{W}_U$ . The third scenario is when  $f$  was overly aggressive and accidentally removed not only paragraphs related to  $T$  but also those related to  $U$ , resulting in  $\mathbf{W}^{**}$ . This corpus violates the proxy-faithfulness assumption that  $\mathbf{W}^{**}$  fully measures  $U$ . Then, we use the three corpora, one at a time, for causal inference. We use an Inverse Propensity Weighting (IPW) estimator, fully described in Appendix C.

## 5 Experiments and Results

Based on the setting described in §4, our analysis produces six estimates, three based on distillation and three based on facts about the data-generating process. Figure 2 shows all estimates.

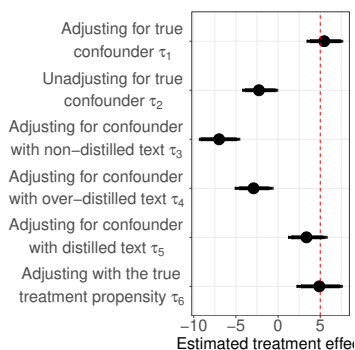


Figure 2: Estimates under different distillation regimes.

The *first estimate*,  $\hat{\tau}_1 = 5.5$ , is the baseline where all information is known to the outcome model, including  $U$ . Because this linear model adjusting for  $U$  and  $\mathbf{X}$  is equivalent to the data-

generating model, and the estimated effect would be equal to the true value of 5 without sampling noise. The bootstrapped 95% confidence interval (CI) is 3.4 to 7.6. The *second estimate*,  $\hat{\tau}_2 = -2.3$ , is obtained when  $U$  is omitted from the model to induce omitted variable bias (CI: -4.2, -0.1).

The *third estimate*,  $\hat{\tau}_3$ , uses IPW to estimate the ATE (see §C). Here, we use the non-distilled documents,  $\mathbf{W}$ , to estimate propensities. As Figure 2 shows, in the absence of distillation, the bias *increases* compared to conditioning on  $\mathbf{X}$  alone, producing  $\hat{\tau}_3 = -7.0$  (CI: -9.4, -4.6). The *fourth estimate*,  $\hat{\tau}_4$ , applies overly aggressive distillation. This approach gives a result similar to the unadjusted estimate:  $\hat{\tau}_4 = -2.9$  (CI: -5.1, -0.6).

The *fifth estimate*,  $\hat{\tau}_5$ , applies oracle distillation by removing the paragraphs we know were affected by  $T$ . Using  $\mathbf{W}^*$ , the bias is reduced substantially, yielding an estimate  $\hat{\tau}_5 = 3.5$  (CI: 1.2, 5.8). As the CI of this  $\hat{\tau}$  includes the true  $\tau = 5$ , we conclude that distillation successfully recovers  $\tau$ . However, we note that this recovery is not perfect and will be affected by sampling and modeling parameters.

The *sixth estimate*,  $\hat{\tau}_6$ , demonstrates the impact of model selection for the propensity estimator. Using the *true* (simulated) propensity, the IPW estimate is  $\hat{\tau}_6 = 4.9$  (CI: 2.2, 7.6). This result shows that further gains could be made by careful model selection (Chernozhukov et al., 2018).

Figure 3 shows distributions of propensity values for  $\hat{\tau}_3$ ,  $\hat{\tau}_5$ , and  $\hat{\tau}_6$ . Without distillation (red), the estimated propensities cluster near 0 and 1.  $T$  is predicted almost perfectly, as mentioned in §2.1.1, causing the IPW estimate to be similar to the unweighted one. Conversely, with distillation, the predicted probabilities are now similar to the data-generating propensities, and thereby, the resulting causal estimate is improved.

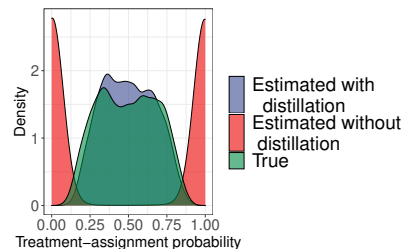


Figure 3: Estimated and true assignment probabilities.

## 6 Discussion

This paper shows the critical role of the no-treatment-leakage assumption when using text for

causal inference. While text is becoming an established data source, it may harbour valuable information about a confounder but also contaminating information about post-treatment effects. This issue has seen little discussion in text-based causal inference literature (Mozer et al., 2020; Roberts et al., 2020; Feder et al., 2021; Daoud and Dubhashi, 2020), but has the potential to severely bias causal estimates, potentially leading to false discoveries or invalid policy recommendations in social and health settings (Kino et al., 2021; Daoud et al., 2017; Balgi et al., 2022).

Before discussing the implication of treatment leakage, three limitations should be considered. First, more work is required to show how the no-treatment-leakage assumption operates under different covariance structures (i.e., different data-generating processes). Second, a larger simulation framework is needed to decompose estimator bias and variance. Third, all results are based on simulated data, and more research is needed to generalize our insights to real data. Although simulate data are idealized, they provide a benefit by allowing us to analyze the mechanics of treatment leakage and text distillation in a controlled environment. Based on our simulated data, our analysis shows that when the no-treatment-leakage assumption is violated, effect estimates will be severely biased. In the presence of treatment leakage, scholars may be better off abstaining from using a non-distilled text to adjust for confounding. Although, in theory, the best solution is to use a text distillation that removes all treatment leakage, in practice, using distillation can be difficult to achieve.

Therefore, one critical extension of our work is to develop methods that estimates the amount of treatment leakage in text. This estimate will enable applied researchers to make an informed decision about whether to adjusting for text-based confounding or abstain from it when leakage is high, and when text distillation is not an option.

A second extensions is to develop a generalized framework that accounts for when the adjusted text represents multiple nodes in a DAG (in combination with the confounding and the treatment or without them). While our article focuses on treatment leakage, there are other types of leakage when a single document is a function of combinations of DAG nodes such as the outcome, confounder, treatment, mediator, or instrument. Thus, a generalization of the no-treatment-leakage assumption is the

*no-node-leakage assumption*. Such methods will benefit from insights established in the literature on causal inference with proxies (Peña, 2020; VanderWeele, 2019; Miao et al., 2018; Rissanen and Marttinen, 2021). A third extension is to develop a variety of text distillation methods, suitable for different application settings. Researchers need alternative frameworks when human partitioning of text is not possible to achieve manually, because of corpus size or language complexities. *Automatic distillation* could be attempted with additional assumptions, perhaps building from the literature on removing sensitive information in text representations (Bolukbasi et al., 2016; Ravfogel et al., 2020).

## Acknowledgements

Richard Johansson was supported by the projects *Interpreting and Grounding Pre-trained Representations for NLP* and *Representation Learning for Conversational AI*, both funded by Wallenberg AI, Autonomous Systems and Software Program (WASP) funded by the Knut and Alice Wallenberg Foundation. Adel Daoud would like to acknowledge a grant from The Royal Swedish Academy of Letters, History and Antiquities.

## References

- Joakim Åkerström, Adel Daoud, and Richard Johansson. 2019. [Natural language processing in policy evaluation: Extracting policy conditions from IMF loan agreements](#). In *Proceedings of the 22nd Nordic Conference on Computational Linguistics*, pages 316–320, Turku, Finland. Linköping University Electronic Press.
- Sourabh Balgi, Jose M. Peña, and Adel Daoud. 2022. [Personalized Public Policy Analysis in Social Sciences using Causal-Graphical Normalizing Flows](#). *Association for the Advancement of Artificial Intelligence: AI for Social Impact track*. ArXiv: 2202.03281.
- Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. [Man is to computer programmer as woman is to homemaker? Debiasing word embeddings](#). In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc.
- Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James Robins. 2018. [Double/debiased machine learning for treatment and structural parameters](#). *The Econometrics Journal*, 21(1):C1–C68.
- Adel Daoud and Devdatt Dubhashi. 2020. [Statistical modeling: The three cultures](#). *arXiv:2012.04570*.

- Adel Daoud, Elias Nosrati, Bernhard Reinsberg, Alexander E. Kentikelenis, Thomas H. Stubbs, and Lawrence P. King. 2017. [Impact of International Monetary Fund programs on child health](#). *Proceedings of the National Academy of Sciences*, 114(25):6492–6497.
- Adel Daoud, Bernhard Reinsberg, Alexander E. Kentikelenis, Thomas H. Stubbs, and Lawrence P. King. 2019. [The International Monetary Fund’s interventions in food and agriculture: An analysis of loans and conditions](#). *Food Policy*, 83:204–218.
- Amir Feder, Katherine A. Keith, Emaad Manzoor, Reid Pryzant, Dhanya Sridhar, Zach Wood-Doughty, Jacob Eisenstein, Justin Grimm, Roi Reichart, Margaret E. Roberts, Brandon M. Stewart, Victor Veitch, and Diyi Yang. 2021. [Causal inference in natural language processing: Estimation, prediction, interpretation and beyond](#). *arXiv preprint arXiv:2109.00725*.
- Michele Jonsson Funk, Daniel Westreich, Chris Wiesen, Til Stürmer, M Alan Brookhart, and Marie Davidian. 2011. [Doubly robust estimation of causal effects](#). *American journal of epidemiology*, 173(7):761–767.
- Matthew Gentzkow, Bryan Kelly, and Matt Taddy. 2019. [Text as data](#). *Journal of Economic Literature*, 57(3):535–74.
- Carolyn Heinrich, Alessandro Maffioli, and Gonzalo Vazquez. 2010. [A primer for applying propensity-score matching](#). *Inter-American Development Bank*.
- Katherine Keith, David Jensen, and Brendan O’Connor. 2020. [Text and causal inference: A review of using text to remove confounding from causal estimates](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5332–5344, Online. Association for Computational Linguistics.
- Shiho Kino, Yu-Tien Hsu, Koichiro Shiba, Yung-Shin Chien, Carol Mita, Ichiro Kawachi, and Adel Daoud. 2021. [A scoping review on the use of machine learning in research on social determinants of health: Trends and research prospects](#). *SSM - Population Health*, 15:100836.
- Aditya Krishna Menon and Robert C Williamson. 2018. [The cost of fairness in binary classification](#). In *Conference on Fairness, Accountability and Transparency*, pages 107–118. PMLR.
- Wang Miao, Zhi Geng, and Eric J. Tchetgen Tchetgen. 2018. [Identifying causal effects with proxy variables of an unmeasured confounder](#). *Biometrika*, 105(4):987–993.
- Reagan Mozer, Luke Miratrix, Aaron Russell Kaufman, and L. Jason Anastasopoulos. 2020. [Matching with text data: An experimental evaluation of methods for matching documents and of measuring match quality](#). *Political Analysis*, 28(4):445–468.
- Jessica A. Myers, Jeremy A. Rassen, Joshua J. Gagne, Krista F. Huybrechts, Sebastian Schneeweiss, Kenneth J. Rothman, Marshall M. Joffe, and Robert J. Glynn. 2011. [Effects of adjusting for instrumental variables on bias and precision of effect estimates](#). *American journal of epidemiology*, 174(11):1213–1222.
- Judea Pearl. 2015. [Conditioning on post-treatment variables](#). *Journal of Causal Inference*, 3(1):131–137.
- Jose M. Peña. 2020. [On the Monotonicity of a Nondifferentially Mismeasured Binary Confounder](#). *Journal of Causal Inference*, 8(1):150–163. Publisher: De Gruyter Section: Journal of Causal Inference.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. [Language models are unsupervised multitask learners](#). *OpenAI blog*, 1(8):9.
- Shauli Ravfogel, Yanai Elazar, Hila Gonen, Michael Twiton, and Yoav Goldberg. 2020. [Null it out: Guarding protected attributes by iterative nullspace projection](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7237–7256. Association for Computational Linguistics.
- Severi Rissanen and Pekka Marttinen. 2021. [A critical look at the consistency of causal estimation with deep latent variable models](#). In *Advances in Neural Information Processing Systems*, volume 34, pages 4207–4217. Curran Associates, Inc.
- Margaret E. Roberts, Brandon M. Stewart, and Richard A. Nielsen. 2020. [Adjusting for confounding with text matching](#). *American Journal of Political Science*, 64(4):887–903.
- Paul R. Rosenbaum and Donald B. Rubin. 1983. [The central role of the propensity score in observational studies for causal effects](#). *Biometrika*, 70(1):41–55.
- Donald B. Rubin. 1974. [Estimating causal effects of treatments in randomized and nonrandomized studies](#). *Journal of Educational Psychology*, 66(5):688–701.
- Tyler J. VanderWeele. 2019. [Principles of confounder selection](#). *European Journal of Epidemiology*, 34:211–219.
- Zach Wood-Doughty, Ilya Shpitser, and Mark Dredze. 2021. [Generating synthetic text data to evaluate causal inference methods](#). *arXiv preprint arXiv:2102.05638*.

## A Synthetic Data Generation

We first summarize the general approach in this section and provide details for the simulation in §5 in the next section.

For each document  $i$ , we first draw observed and unobserved confounders  $X_i$  and  $U_i$ , and then the

treatment  $T_i$ . For each paragraph  $j$  in the document, we draw a paragraph topic  $Z_{ij}$ , depending on the values of  $U_i$  and  $T_i$ , and then a prompt  $W_{ij}^0$  depending on the value of  $Z_{ij}$ . Finally, we sample from the GPT-2 language model<sup>1</sup> to generate the paragraph text  $W_{ij}$ , starting from the prompt  $W_{ij}^0$  and with a vocabulary distribution shift defined by  $Z_{ij}$ . Algorithm 1 shows the pseudocode.

---

**Algorithm 1** Generation of synthetic data.

---

```

for  $i \in 1, \dots, N$ 
   $\mathbf{X}_i \sim f_X$ 
   $U_i \sim f_U$ 
   $T_i \sim \text{Bernoulli}(\text{sigmoid}(f_T(\mathbf{X}_i, U_i)))$ 
   $Y_i \sim f_Y(\mathbf{X}_i, U_i, T_i)$ 
  for  $j \in 1, \dots, K$ 
     $Z_{ij} \sim \text{Categorical}(f_Z(U_i, T_i))$ 
     $W_{ij}^0 \sim \text{Categorical}(f_{W^0}(Z_{ij}))$ 
     $W_{ij} \sim \text{LM}(W_{ij}^0, Z_{ij})$ 

```

---

In the pseudocode above, the functions  $f_X$ ,  $f_U$ ,  $f_T$ , and  $f_Y$  define the distributions of the observed confounders, unobserved confounder, treatment and outcome, respectively. On the paragraph level, the function  $f_Z$  defines a categorical distribution over paragraph topics, and  $f_{W^0}$  a categorical distribution over prompts.

Similarly to Wood-Doughty et al. (2021), we use two mechanisms to condition the generation of a paragraph on a topic  $Z$ : a prompt and a vocabulary distribution shift. The distribution shift is designed to promote a set of *keywords* related to the topic and we implement it by multiplying the language model probabilities by a topic-specific vector  $\theta_Z$  of scale factors:

$$P'(w|\text{context}, Z) \propto P_{\text{LM}}(w|\text{context}) \cdot \theta_Z(w)$$

## B Parameterization Used in §5

In §5, we generated  $N = 10,000$  instances, each consisting of numerical values and a document. We used the following distributions to generate the document-level variables:  $f_X$  was a 3-dimensional isotropic Gaussian;  $f_U$  was an even coin toss;  $f_T$  was linear in  $\mathbf{X}_i$  and  $U_i$ ;  $f_Y$  was Gaussian with a mean defined by a linear function of  $\mathbf{X}_i$ ,  $U_i$ , and  $T_i$  and a fixed standard deviation.

Each document consisted of  $K = 20$  paragraphs. For the paragraph generation, we defined five dif-

ferent topics: two corresponding to positive and negative treatment values; two corresponding to positive and negative values of the unobserved confounder; one general background topic that was unrelated to  $U$  or  $T$  (but conceptually thought of as controlled by other “residual” variables  $R$ ). For a document with given values of  $U$  and  $T$ , we set the topic distribution  $f_Z$  to select the  $U$  topic with a probability of 0.2, the  $T$  topic with a probability of 0.2, and the general topic with a probability of 0.6.

The generated texts were designed to simulate a hypothetical use case where the researchers want to investigate the effect of IMF programs on some country-level indicator (cf. Daoud et al., 2019). The treatment variable  $T$  represents the presence or absence of an IMF program; the unseen confounder  $U$  represents the political situation of the country with respect to the IMF. For each topic except the general topic, we define four different prompts: for instance, for a positive treatment value, one of the prompts was *The International Monetary Fund mandates the deregulation of [COUNTRY]’s labor market*. In the analysis, “[COUNTRY]” is substituted by randomly sampled country names.

All topics except the general topic defined a distribution shift used when generating from the language model. We used 8 topic keywords for each of these topics. For these keywords, the corresponding entries in the vocabulary distribution shift vector  $\log \theta_Z$  were set to a value that defines the strength of the effect of  $T$  on  $\mathbf{W}$ ; for all other words except these keywords,  $\log \theta_Z$  was 0. Since our focus in this paper is on a clear-cut use case where the effects are strong, we set the strength parameter to a value of 4, which gives a noticeable effect on the generated texts.

The text generation model was run on a single GPU (NVIDIA GeForce GTX TITAN X). Generating the 10,000 documents took around 10 hours. The generation of random text is within the intended use of the GPT-2 model.

The implementation of the algorithm to generate the synthetic data is available in our repository.<sup>2</sup>

## C IPW Details

### C.1 Background

The ATE is defined as  $\tau = \mathbb{E}[Y_i(1) - Y_i(0)]$ , where  $Y_i(t)$  is the potential outcome for unit  $i$  under

<sup>1</sup>We used the implementation from the HuggingFace repository, <https://huggingface.co/gpt2>.

<sup>2</sup><https://github.com/adeldaoud/AIforTextandCausalInference>

treatment  $t$ . It can be identified in randomized experiments (Rubin, 1974). However, the situation is more complicated in the observational setting, where the treatment is not randomized to units but could be correlated with confounders,  $\mathbf{X}_i$ , that are associated with the treatment and the outcome. In that setting, we can, with additional assumptions, still recover the ATE using Inverse Propensity Weighting (IPW) or related robust methods (Funk et al., 2011), where observations are weighted by the inverse of their estimated treatment probabilities  $\hat{\pi}(\mathbf{X}_i) = \widehat{\Pr}(T_i = 1 | \mathbf{X}_i)$  (Rosenbaum and Rubin, 1983):  $\hat{\tau} = \frac{1}{n} \sum_{i=1}^n \left\{ \frac{T_i Y_i}{\hat{\pi}(\mathbf{X}_i)} - \frac{(1-T_i) Y_i}{1-\hat{\pi}(\mathbf{X}_i)} \right\}$ .

## C.2 Estimation

ATE estimates based on Inverse Propensity Weighting (see §C.1) require the estimation of the propensity scores,  $\widehat{\Pr}(T | \mathbf{X}, \mathbf{W})$ . To estimate these scores, we applied a  $L_1$ -regularized logistic regression model using the `glmnet` package in R. The regularization strength ( $\lambda$ ) was set automatically via 10-fold cross-validation. When estimating propensities, we represented the (non-distilled or distilled) document as an  $L_2$ -normalized TF-IDF vector using the 256 most frequent terms in the vocabulary, while the numerical covariates  $\mathbf{X}$  were standardized.