

# DoubleLingo: Causal Estimation with Large Language Models

Anonymous ACL submission

## Abstract

001 Estimating causal effects from non-randomized  
002 data requires assumptions about the underlying  
003 data-generating process. To achieve unbiased  
004 estimates of the causal effect of a treatment on  
005 an outcome, we must adjust for any confound-  
006 ing variables that influence both treatment and  
007 outcome. When such confounders include text  
008 data, existing causal inference methods strug-  
009 gle due to the high dimensionality of the text.  
010 The simple statistical models which have suffi-  
011 cient convergence criteria for causal estimation  
012 are not well-equipped to handle noisy unstruc-  
013 tured text, but flexible Large language models  
014 (LLMs) that excel at predictive tasks with text  
015 data do not meet the statistical assumptions  
016 necessary for causal estimation. Our method  
017 enables theoretically consistent estimation of  
018 causal effects using LLM-based nuisance mod-  
019 els by incorporating them within the framework  
020 of Double Machine Learning. On the best avail-  
021 able dataset for evaluating such methods, we  
022 obtain a 10.4% reduction in the relative abso-  
023 lute error for the estimated causal effect over  
024 existing methods.

## 025 1 Introduction

026 A common goal of scientific research is the analy-  
027 sis of causal relationships (Triantafillou et al., 2017;  
028 Sanna et al., 2019; Chang et al., 2022). Consider  
029 the following motivating example, where a phar-  
030 maceutical company wants to estimate the causal  
031 effect of the prescription of antibiotics (treatment)  
032 on the patient’s disease progression (outcome). The  
033 causal effect is defined as the expected change  
034 in disease progression across two *counterfactual*  
035 worlds which only differ in whether the patient  
036 is given antibiotics (Hernán, 2004). When ran-  
037 domization is impossible or unethical, we estimate  
038 causal effects from observational data using as-  
039 sumptions about the underlying data distribution.  
040 Confounders – variables affecting both the treat-  
041 ment and outcome – introduce potential bias that  
042 must be addressed.

043 When data is low-dimensional, confounding can  
044 be controlled for using various methods from the lit-  
045 erature (Pearl, 2009). However, several challenges  
046 arise in the case of high-dimensional confounders  
047 such as text. For example, assume the pharmaceu-  
048 tical company has free-text clinical notes that may  
049 include information about patients’ histories, diag-  
050 noses, or relationships with their doctors (Rajkomar  
051 et al., 2018). If these potential confounding vari-  
052 ables appear nowhere else in the patients’ records,  
053 then to account for confounding we must use text-  
054 based causal methods (Rosenbloom et al., 2011;  
055 Wu et al., 2013). Since text is high-dimensional,  
056 it requires sophisticated modeling that captures se-  
057 mantic meaning.

058 Existing models often utilize overly simplified  
059 representations of the text (Wood-Doughty et al.,  
060 2018; Keith et al., 2020), such as a *bag-of-words*  
061 (BoW) representation. While such representations  
062 combined with simple estimation models allow for  
063 consistent<sup>1</sup> estimation, they may fail to capture the  
064 true complexity of the text’s underlying relation-  
065 ships. The use of Large language models (LLMs)  
066 in causal estimation has only recently been stud-  
067 ied (Veitch et al., 2020), and many researchers  
068 suggest the need for more sophisticated natural  
069 language processing (NLP) techniques (Wood-  
070 Doughty et al., 2021; Feder et al., 2022; Keith et al.,  
071 2023). However, while LLMs excel at predictive  
072 tasks, they do not meet the necessary statistical  
073 assumptions for a consistent causal estimation.

074 We present **DoubleLingo**, combining Double  
075 Machine Learning with LLM-based nuisance mod-  
076 els to enable a theoretically consistent estimation of  
077 causal effects with text-based confounding. We test  
078 our model on a novel dataset (Keith et al., 2023),  
079 obtaining the best causal effect estimates reported  
080 thus far. In particular, our relative absolute error is  
081 over 10% lower than the best current models.

<sup>1</sup>Defined in more detail in §3

## 2 Causal Inference Background

While causal inference is a broad and diverse field (Robins et al., 2000; Pearl, 2009), we provide a brief introduction here. For recent surveys of causal inference and natural language processing, see Keith et al. (2020) or Feder et al. (2022).

### 2.1 DAGs & Counterfactuals

The motivating example described above is illustrated by the directed acyclic graph (DAG) in Figure 1, where we use a binary random variable  $A$  to indicate whether the patient receives ( $A = 1$ ) antibiotics or not ( $A = 0$ ). We similarly use a binary  $Y$  to denote whether the disease progresses ( $Y = 1$ ) or not ( $Y = 0$ ). An arrow in the DAG such as  $A \rightarrow Y$  indicates that  $A$  has a potential causal effect on  $Y$ . Finally, we denote  $T$  as the patient medical records, and  $C$  as the set of all confounding variables contained in the records. Most importantly,  $C$  is unobserved — we don’t know the exact confounding variables, but we have access to the text  $T$  containing them. Hence, there exist some causal effects  $T \dashrightarrow A$  and  $T \dashrightarrow Y$  where the text  $T$  affects  $A$  and  $Y$  through the unobserved  $C$ . The *counterfactual* outcome  $Y^{a=1}$  represents the hypothetical disease progression had we intervened to assign  $A = 1$  (prescribe antibiotics), and  $Y^{a=0}$  is defined analogously. In causal inference, the most common estimand is the average treatment effect (*ATE*) of  $A$  on  $Y$ , computed as:

$$ATE = \mathbb{E}[Y^{a=1} - Y^{a=0}] \quad (1)$$

A fundamental problem is that we can never simultaneously observe both *counterfactuals*  $Y^{a=1}, Y^{a=0}$  (Holland, 1986), thus we need a way to compute the *ATE* only utilizing observed data.

### 2.2 Identification Assumptions

We proceed by assuming *consistency*, requiring that the outcome we observe for any possible treatment  $a$  is equal to the *counterfactual* outcome we would have observed had we intervened to assign  $A = a$ . Formally,

$$A = a \leftrightarrow Y^a = Y \quad (2)$$

We then assume *conditional exchangeability*, requiring the independence between our counterfactual  $Y^a$  and the observed treatment  $A$  conditioned on all confounders  $C$ , formalized as

$$Y^a \perp A \mid C \quad \forall a \in \{0, 1\} \quad (3)$$

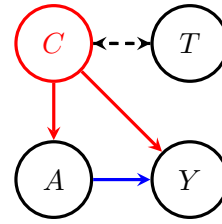


Figure 1: Textual Confounding DAG with Treatment  $A$ , Outcome  $Y$ , Confounders  $C$ , and Text  $T$ .

Using these assumptions, we may compute the counterfactual  $\mathbb{E}[Y^a]$  as follows

$$\mathbb{E}[Y^a] = \sum_C \mathbb{E}[Y^a \mid C] \mathbb{P}(C) \quad (4)$$

$$\stackrel{(3)}{=} \sum_C \mathbb{E}[Y^a \mid A = a, C] \mathbb{P}(C) \quad (5)$$

$$\stackrel{(2)}{=} \sum_C \mathbb{E}[Y \mid A = a, C] \mathbb{P}(C) \quad (6)$$

However, since  $C$  is unobserved, the main challenge is in modelling the text  $T$  to adjust for all of the confounding from  $C$ .

### 2.3 Causal Effect Estimation

In estimating the *ATE*, we thus require (a) a representation of the text and (b) an appropriate causal estimation method. As mentioned in §1, a BoW text representation is commonly used by existing text-based causal estimators. For (b), there are countless estimation methods, and we refer the reader to a much more exhaustive guide by Peters et al. (2017). One such commonly used method is the *Inverse Propensity of Treatment Weighting* (IPTW), where  $\mathbb{E}[Y^a]$  is calculated as follows for a dataset of size  $N$ .

$$\mathbb{E}[Y^a] = \frac{1}{N} \sum_{i \in [N]} Y_i \frac{\mathbb{1}(A_i = a)}{\mathbb{P}(A_i = a \mid T)} \quad (7)$$

Thus, combining (a) and (b), a common current method is to use IPTW and learn a *Logistic Regression* model  $\mathbb{P}(A \mid T)$  for the propensity of the treatment  $A$  given a BoW text representation  $T$ .

## 3 Model

Any estimator  $\hat{\theta}$  of the true *ATE* estimate  $\theta$  must be both unbiased and consistent such that

$$\mathbb{E}[\hat{\theta}] = \theta \quad \text{and} \quad \hat{\theta} \xrightarrow{P} \theta \quad (8)$$

While LLMs have drastically changed the field of NLP (Vaswani et al., 2017; Min et al., 2023),

they are not consistent estimators of causal parameters due to both explicit and implicit regularization (Neyshtabur, 2017; Chernozhukov et al., 2018). Thus, a naive approach of using an LLM such as BERT (Devlin et al., 2019) to learn the propensity  $\mathbb{P}(A | T)$  in Equation (7) would be biased.

### 3.1 Double Machine Learning

We thus turn to Double Machine Learning (DML), which has never previously been used in the context of LLMs. As proven by Chernozhukov et al. (2018), regularization bias in complex ML models can be overcome by utilizing orthogonalization. In particular, we partial out this bias by learning classifiers for both the treatment  $\mathbb{E}[A | T]$  and outcome  $\mathbb{E}[Y | T]$ . Accordingly, we obtain a consistent estimate of the *ATE* by regressing the residuals

$$Y - \mathbb{E}[Y | T] \sim A - \mathbb{E}[A | T] \quad (9)$$

Additionally, as we fit both  $\mathbb{E}[A | T]$  and  $\mathbb{E}[Y | T]$ , the estimation is doubly robust such that only one of the two models need to be correctly specified to obtain an unbiased *ATE* (Funk et al., 2011). Finally, we utilize sample splitting (Stone, 1974), where we train on half of the data, using the other half for estimation, preventing any estimation bias induced by overfitting. A potential concern is that DML requires our nuisance models to converge at  $N^{-1/4}$  rates such that the overall estimator is  $\sqrt{N}$ -consistent<sup>2</sup>, that is

$$\hat{\theta} - \theta = O_p(N^{-1/2}) \quad (10)$$

While there is research on the rate of convergence of misclassification probability (Gurevych et al., 2022) for encoder-based transformer classifiers such as BERT, its convergence rate for semiparametric inference is unknown.

### 3.2 Faster Converging Model Variations

Since fully fine-tuning BERT classifiers within the DML framework may not be appropriate, we present **DoubleLingo**, utilizing two faster converging model variations.

**BERT+Adapter.** Our first configuration utilizes parameter efficient transfer learning in the form of adapters (Houlsby et al., 2019). Thus, instead of fine-tuning all of BERT, we only fine-tune the adapter layers. While there are no theoretical

bounds for the convergence of adapters, they empirically demonstrate a much quicker convergence compared to fine-tuning the full network.

**Embedding+FFN.** Fully-connected feedforward neural networks (FFNs) with the ReLU activation function have been proven to converge at  $N^{-1/4}$  rates for their use in semiparametric inference (Farrell et al., 2021). Thus, instead of fine-tuning BERT at all, a potential approach is to fine-tune a feedforward layer on top of BERT’s pre-trained embeddings. Since BERT doesn’t learn independent sentence embeddings, we could instead use the  $[CLS]$  encoding or pool the sequence of hidden states for the whole input. Instead, we utilize embeddings from pre-trained sentence transformers (Reimers and Gurevych, 2019), which are much more semantically meaningful. To our knowledge, sentence transformer embeddings have never been utilized in the context of causal inference estimation, thus we further contribute to the literature by analyzing their causal estimation capabilities compared to simpler text representations.

## 4 Causal Dataset & Experiment

Unlike supervised learning models, which can be evaluated on held-out test sets with ground-truth labels, causal estimation methods require evaluations with *counterfactual* ground-truth, which is impossible to measure from observed data (Holland, 1986). Researchers often turn to (semi-)synthetic data, for which there is a tension between generating realistic text and maintaining full knowledge of the underlying data-generating process (DGP) (Wood-Doughty et al., 2021). Most current datasets fail to accomplish both, either fully specifying the DGP but with unrealistic text (Johansson et al., 2016; Yao et al., 2019), or using real-world text inside a semi-synthetic DGP (Veitch et al., 2020).

### 4.1 Dataset and Baselines

A recent novel dataset employs a randomized controlled trial (RCT) rejection sampling algorithm to create text-based datasets that both contain real text and are based on a realistic DGP (Keith et al., 2023). In particular, the authors fix  $C$  to be a single binary confounding variable contained in the text and choose RCT’s with an existing  $C \rightarrow Y$  relationship. They then sample the dataset to artificially create a  $C \rightarrow A$  relationship and evaluate 8 different models over 100 sampled dataset subsets.

<sup>2</sup>As  $N \rightarrow \infty$  estimation error goes to 0 at a rate of  $\sqrt{N}$

Unadjusted	Oracle (C)	<i>TF-IDF+FFN</i>	LR <sub>Q</sub>	LR <sub>IPTW</sub>	LR <sub>AIPTW</sub>	LR <sub>DML</sub>
0.214 (0.08)	0.115 (0.09)	0.118 (0.09)	1.408 (1.00)	0.470 (0.16)	1.579 (0.66)	1.899 (0.91)
<b>BERT+Adapter</b>	<b>SPECTER+FFN</b>	<b>MPNetV2+FFN</b>	CB <sub>Q</sub>	CB <sub>IPTW</sub>	CB <sub>AIPTW</sub>	CB <sub>DML</sub>
0.104 (0.08)	0.104 (0.08)	0.103 (0.08)	0.237 (0.10)	0.141 (0.11)	0.115 (0.10)	0.128 (0.10)

Table 1: Relative Absolute Error mean (variance) for all methods over 100 subsets. §4.2 describes our **DoubleLingo** methods and *TF-IDF+FFN* baseline. Logistic Regression (LR), CatBoost (CB), Oracle, and Unadjusted baselines all use code from Keith et al. (2023). Our methods achieve the best (lowest) error and variance.

They train *Logistic Regression* and *CatBoost* nuisance models based on a BoW representation for the text, combining both with 4 different causal estimation techniques, including IPTW, Augmented-IPTW (AIPTW), Outcome Regression (Q), and DML. They finally evaluate an *Oracle* with full access to the unobserved  $C$  value.

## 4.2 DoubleLingo Experiments

We now describe our methods that use LLMs inside the DML framework. Our **BERT+Adapter** method fine-tunes adapters within BERT classifiers for both  $A$  and  $Y$  (Houlsby et al., 2019). Our **Embedding+FFN** configuration uses two sentence transformers. First, *all-mpnet-base-v2*<sup>3</sup>, based on *MPNet* (Song et al., 2020) and fine-tuned on 1 million sentence pairs. Second, *SPECTER* (Cohan et al., 2020), pre-trained on a dataset of scientific paper titles and abstracts which matches the format of Keith et al. (2023). For both **Embedding+FFN** methods, we use a single hidden layer, ReLU activation functions, and the AdamW optimizer (Loshchilov and Hutter, 2018) to obtain  $N^{-1/4}$  convergence (Farrell et al., 2021). Finally, we implement a *TF-IDF+FFN* baseline, following Manzoor et al. (2023), which uses DML with FFNs with batch normalization (Ioffe and Szegedy, 2015) and a TF-IDF text representation. A more detailed implementation, including specific hyperparameters and RCT parameterization choices are provided in Appendix A.

## 5 Results and Conclusions

Table 1 shows that our three **DoubleLingo** estimators obtain the lowest *ATE* relative absolute error (0.103), a 10.4% decrease from the prior best (0.115). These results provide strong empirical evidence that the DML framework successfully enables the use of LLMs in causal estimation. Notably, the prior best was achieved by both a BoW

model (CB<sub>AIPTW</sub>) and the *Oracle* estimator which calculates the estimates using the unobserved  $C$  values. If  $C$  contained all causes of  $A$  and  $Y$ , it would be the theoretically-optimal efficient adjustment set (Rotnitzky and Snucler, 2020) and the Oracle should – asymptotically – be impossible to outperform. However, while the  $C \rightarrow A$  relationship is artificially induced by the sampling procedure of Keith et al. (2023), the  $C \rightarrow Y$  correlation is confirmed to exist<sup>4</sup>; we hypothesize that the underlying complexity of the  $T \rightarrow Y$  relationship is not fully captured by the binary topic  $C$  and thus modeling  $T$  allows for more efficient estimation.

Our results specifically support the hypothesis that the text representation itself matters to causal estimation. Among all DML methods with feed-forward classifiers, our **Embedding+FFN** methods’ outperformance of our *TF-IDF+FFN* baseline shows that better representations can enable lower estimation error. Appendix B also shows our models’ slightly better classification accuracy than the *TF-IDF+FFN* baseline during estimation.

Between our three proposed methods, we see no large differences in performance. This suggests that while the incorporation of LLMs into the estimators is essential, the specific architecture and training setup matters less. However, **BERT+Adapter** trains two to three times slower than **Embedding+FFN**. We also see little difference between the two pre-trained embeddings, despite the similarity of the *SPECTER* embedding’s dataset to that of our evaluation data.

This work proposes **DoubleLingo**, a theoretically consistent causal estimator that uses LLM nuisance models inside the DML framework. We show that both adapters and sentence transformers can achieve the lowest estimation error on the best available dataset for evaluating methods that account for text confounding. We include our code as an appendix that reproduces our provided results.

<sup>3</sup><https://hf.co/sentence-transformers/all-mpnet-base-v2>

<sup>4</sup>Authors verify that  $C \not\perp Y$  with an odds ratio test

## 329 Limitations

330 The main limitation of our estimation procedure  
331 is compute time – training the **BERT+Adapter**  
332 configuration on 100 sampled dataset subsets takes  
333 10 hours parallelized across 2 RTX 8000’s, signifi-  
334 cantly longer than the baseline *Linear Regression*  
335 or *CatBoost* models. In particular, our model’s re-  
336 liance on sample-splitting and double robustness to  
337 obtain a consistent final estimate requires training  
338 4 times as many models per each dataset subset.  
339 However, it’s important to note that the **Embed-**  
340 **ding+FFN** configurations only take a third of the  
341 time, yet achieve identical results.

342 Additionally, our work only focuses on causal  
343 estimation with text-based confounding. In partic-  
344 ular, dealing with textual treatments or outcomes  
345 is still an open problem in the field (Feder et al.,  
346 2022). Finally, we only train on a single English-  
347 language dataset, and resultingly encourage future  
348 work to expand on this by testing other types of  
349 text-based RCT’s.

## 350 References

- 351 Chun-Wei Chang, Stephan B Munch, and Chih-hao  
352 Hsieh. 2022. Comments on identifying causal rela-  
353 tionships in nonlinear dynamical systems via empir-  
354 ical mode decomposition. *Nature communications*,  
355 13(1):2860.
- 356 Victor Chernozhukov, Denis Chetverikov, Mert Demirer,  
357 Esther Duflo, Christian Hansen, Whitney Newey, and  
358 James Robins. 2018. Double/debiased machine learn-  
359 ing for treatment and structural parameters: Double/  
360 debiased machine learning. *The Econometrics*  
361 *Journal*, 21(1).
- 362 Arman Cohan, Sergey Feldman, Iz Beltagy, Doug  
363 Downey, and Daniel Weld. 2020. **SPECTER:**  
364 **Document-level representation learning using**  
365 **citation-informed transformers.** In *Proceedings*  
366 *of the 58th Annual Meeting of the Association*  
367 *for Computational Linguistics*, pages 2270–2282,  
368 Online. Association for Computational Linguistics.
- 369 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and  
370 Kristina Toutanova. 2019. **BERT: Pre-training of**  
371 **deep bidirectional transformers for language**  
372 **understanding.** In *Proceedings of the 2019 Conference*  
373 *of the North American Chapter of the Association for*  
374 *Computational Linguistics: Human Language Tech-*  
375 *nologies, Volume 1 (Long and Short Papers)*, pages  
376 4171–4186, Minneapolis, Minnesota. Association for  
377 Computational Linguistics.
- 378 Max H Farrell, Tengyuan Liang, and Sanjog Misra.  
379 2021. Deep neural networks for estimation and infer-  
380 ence. *Econometrica*, 89(1):181–213.

- Amir Feder, Katherine A. Keith, Emaad Manzoor, Reid  
Pryzant, Dhanya Sridhar, Zach Wood-Doughty, Ja-  
cob Eisenstein, Justin Grimmer, Roi Reichart, Mar-  
garet E. Roberts, Brandon M. Stewart, Victor Veitch,  
and Diyi Yang. 2022. **Causal inference in natural lan-**  
**guage processing: Estimation, prediction, interpreta-**  
**tion and beyond.** *Transactions of the Association for*  
*Computational Linguistics*, 10:1138–1158. 381  
382  
383  
384  
385  
386  
387  
388
- 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. 389  
390  
391  
392
- Iryna Gurevych, Michael Kohler, and Gözde Gül Şahin.  
2022. On the rate of convergence of a classifier based  
on a transformer encoder. *IEEE Transactions on*  
*Information Theory*, 68(12):8139–8155. 393  
394  
395  
396
- Miguel Angel Hernán. 2004. A definition of causal  
effect for epidemiological research. *Journal of Epi-*  
*demiology & Community Health*, 58(4):265–271. 397  
398  
399
- Paul W. Holland. 1986. **Statistics and causal infer-**  
**ence.** *Journal of the American Statistical Association*,  
81(396):945–960. 400  
401  
402
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski,  
Bruna Morrone, Quentin De Laroussilhe, Andrea  
Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019.  
Parameter-efficient transfer learning for nlp. In *In-*  
*ternational Conference on Machine Learning*, pages  
2790–2799. PMLR. 403  
404  
405  
406  
407  
408
- Sergey Ioffe and Christian Szegedy. 2015. **Batch nor-**  
**malization: Accelerating deep network training by re-**  
**ducing internal covariate shift.** In *Proceedings of the*  
*32nd International Conference on Machine Learn-*  
*ing*, volume 37 of *Proceedings of Machine Learning*  
*Research*, pages 448–456, Lille, France. PMLR. 409  
410  
411  
412  
413  
414
- Fredrik Johansson, Uri Shalit, and David Sontag. 2016.  
**Learning representations for counterfactual inference.**  
In *Proceedings of The 33rd International Conference*  
*on Machine Learning*, volume 48 of *Proceedings of*  
*Machine Learning Research*, pages 3020–3029, New  
York, New York, USA. PMLR. 415  
416  
417  
418  
419  
420
- 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 As-*  
*sociation for Computational Linguistics*, pages 5332–  
5344, Online. Association for Computational Lin-  
guistics. 421  
422  
423  
424  
425  
426  
427
- Katherine A Keith, Sergey Feldman, David Jurgens,  
Jonathan Bragg, and Rohit Bhattacharya. 2023. Rct  
rejection sampling for causal estimation evaluation.  
*arXiv preprint arXiv:2307.15176*. 428  
429  
430  
431
- Ilya Loshchilov and Frank Hutter. 2018. Decoupled  
weight decay regularization. In *International Confer-*  
*ence on Learning Representations*. 432  
433  
434

435	Emaad Manzoor, George H Chen, Dokyun Lee, and Michael D Smith. 2023. Influence via ethos: On the persuasive power of reputation in deliberation online. <i>Management Science</i> .	<i>Neural Information Processing Systems</i> , 33:16857–16867.	490
436			491
437		M. Stone. 1974. Cross-validatory choice and assessment of statistical predictions. <i>Journal of the Royal Statistical Society. Series B (Methodological)</i> , 36(2):111–147.	492
438			493
439	Bonan Min, Hayley Ross, Elior Sulem, Amir Poursan Ben Veyseh, Thien Huu Nguyen, Oscar Sainz, Eneko Agirre, Ilana Heintz, and Dan Roth. 2023. Recent advances in natural language processing via large pre-trained language models: A survey. <i>ACM Computing Surveys</i> , 56(2):1–40.	Sofia Triantafillou, Vincenzo Lagani, Christina Heinze-Deml, Angelika Schmidt, Jesper Tegner, and Ioannis Tsamardinos. 2017. Predicting causal relationships from biological data: Applying automated causal discovery on mass cytometry data of human immune cells. <i>Scientific reports</i> , 7(1):12724.	494
440			495
441			496
442			497
443			498
444			499
445	Behnam Neyshabur. 2017. Implicit regularization in deep learning. <i>arXiv preprint arXiv:1709.01953</i> .		500
446			501
447	Judea Pearl. 2009. <i>Causality</i> . Cambridge university press.	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	502
448			503
449	Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. 2017. <i>Elements of Causal Inference: Foundations and Learning Algorithms</i> . The MIT Press.		504
450			505
451			506
452	Alvin Rajkomar, Eyal Oren, Kai Chen, Andrew M Dai, Nissan Hajaj, Michaela Hardt, Peter J Liu, Xiaobing Liu, Jake Marcus, Mimi Sun, et al. 2018. Scalable and accurate deep learning with electronic health records. <i>NPJ digital medicine</i> , 1(1):18.	Victor Veitch, Dhanya Sridhar, and David Blei. 2020. Adapting text embeddings for causal inference. In <i>Conference on Uncertainty in Artificial Intelligence</i> , pages 919–928. PMLR.	507
453			508
454			509
455			510
456			511
457	Nils Reimers and Iryna Gurevych. 2019. <i>Sentence-BERT: Sentence embeddings using Siamese BERT-networks</i> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.	Zach Wood-Doughty, Ilya Shpitser, and Mark Dredze. 2018. <i>Challenges of using text classifiers for causal inference</i> . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 4586–4598, Brussels, Belgium. Association for Computational Linguistics.	512
458			513
459			514
460			515
461			516
462			517
463			518
464			519
465	James M Robins, Miguel Angel Hernan, and Babette Brumback. 2000. Marginal structural models and causal inference in epidemiology. <i>Epidemiology</i> , pages 550–560.	Zach Wood-Doughty, Ilya Shpitser, and Mark Dredze. 2021. Generating synthetic text data to evaluate causal inference methods. <i>arXiv preprint arXiv:2102.05638</i> .	520
466			521
467			522
468			523
469	S Trent Rosenbloom, Joshua C Denny, Hua Xu, Nancy Lorenzi, William W Stead, and Kevin B Johnson. 2011. Data from clinical notes: a perspective on the tension between structure and flexible documentation. <i>Journal of the American Medical Informatics Association</i> , 18(2):181–186.	Chia-Yi Wu, Chin-Kuo Chang, Debbie Robson, Richard Jackson, Shaw-Ji Chen, Richard D Hayes, and Robert Stewart. 2013. Evaluation of smoking status identification using electronic health records and open-text information in a large mental health case register. <i>PloS one</i> , 8(9):e74262.	524
470			525
471			526
472			527
473			528
474			529
475	Andrea Rotnitzky and Ezequiel Smucler. 2020. Efficient adjustment sets for population average causal treatment effect estimation in graphical models. <i>The Journal of Machine Learning Research</i> , 21(1):7642–7727.	Liuyi Yao, Sheng Li, Yaliang Li, Hongfei Xue, Jing Gao, and Aidong Zhang. 2019. <i>On the estimation of treatment effect with text covariates</i> . In <i>Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19</i> , pages 4106–4113. International Joint Conferences on Artificial Intelligence Organization.	530
476			531
477			532
478			533
479			
480	Serena Sanna, Natalie R van Zuydam, Anubha Mahajan, Alexander Kurilshikov, Arnau Vich Vila, Urmo Vösa, Zlatan Mujagic, Ad AM Masclee, Daisy MAE Jonkers, Marije Oosting, et al. 2019. Causal relationships among the gut microbiome, short-chain fatty acids and metabolic diseases. <i>Nature genetics</i> , 51(4):600–605.		
481			
482			
483			
484			
485			
486			
487	Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2020. Mpnet: Masked and permuted pre-training for language understanding. <i>Advances in</i>		
488			
489			

Model	Accuracy	
	$\mathbb{E}[A   T]$	$\mathbb{E}[Y   T]$
Logistic Regression	75.5	82.8
CatBoost	80.3	95.5
<i>TF-IDF+FFN</i>	80.6	95.3
<b>SPECTER+FFN</b>	82.8	95.7
<b>MPNetV2+FFN</b>	83.2	95.7
<b>BERT+Adapter</b>	83.2	95.7

Table 2: Average Predictive Accuracy over 100 dataset subsets

## A Implementation

This section gives a more detailed overview of our implementation, including specific hyper-parameter values for both model configurations and parameterization choices of  $\mathbb{P}(A | C)$  required by the RCT rejection sampling algorithm.

**BERT+Adapters.** For our BERT adapter configuration, we use a batch size of 128, the maximum that can fit parallelized across two RTX 8000’s. We use default values for beta and weight decay, setting  $B_1 = 0.9$ ,  $B_2 = 0.999$ ,  $\lambda = 0$ . We manually optimize for the learning rate and number of epochs based on validation accuracy on a small subset of the 100 datasets, resulting in a learning rate of  $3e-4$  over 5 epochs. Our estimation takes around 10 hours to complete. For the estimation of a single dataset, we suggest practitioners perform a larger search over hyper-parameters, however the use of sample-splitting and doubly-robust estimation requires training 4 times the number of models. Thus, a simple grid-search over just 10 hyper-parameter combinations with 4-fold cross-validation over 100 dataset seeds would require the training of 16,000 models. Finally, we use  $\text{BERT}_{\text{BASE}}$  which has 109,482,240 parameters, however the use of adapters allows us to only fine-tune 894,528 parameters.

**Embedding+FFN.** For all of our FFN configurations, we use the same batch size of 128 and the same default beta and weight decay values. We use a single hidden layer with the same number of nodes as the input layer, equal to 768 for both sentence transformers. Since these FFNs are much quicker to train, we perform a search over the learning rates,  $\{1e-5, 1e-4, 1e-3, 1e-2\}$ , combined with early-stopping for each one of the 100 dataset subsets.

**RCT parameterization.** The RCT rejection sampling algorithm requires practitioners to specify  $\mathbb{P}(A | C)$ . In particular, the authors choose  $C$  to be a binary random variable representing the specific text topic. We accordingly utilize the default provided RCT using medicine ( $C = 0$ ) and physics ( $C = 1$ ) articles. Authors then define  $\mathbb{P}(A | C)$  as follows

$$\mathbb{P}(A = 1 | C) = \begin{cases} \zeta_0 & \text{if } C = 0 \\ \zeta_1 & \text{if } C = 1 \end{cases}$$

which is used in sampling the RCT to create an artificial  $C \rightarrow A$  effect. We utilize the default choices of  $\zeta_0 = 0.85$  and  $\zeta_1 = 0.15$  which induce the highest amount of confounding. For a much more thorough explanation, we direct readers to [Keith et al. \(2023\)](#).

## B Nuisance Model Predictive Accuracy

Specific values for the average predictive accuracy during estimation of all tested nuisance models are provided in Table 2. A similar trend appears compared to causal estimation results in Table 1, where the largest improvement occurs from simply switching to non-linear nuisance models (*CatBoost* vs. *LogisticRegression*).

While our three **DoubleLingo** model configurations achieve the best predictive accuracies (83.2%, 95.7%), the values are only slightly higher than those for the *TF-IDF+FFN* implementation. Here, it’s important to note that predictive accuracy alone does not directly contribute to a more accurate estimation ([Wood-Doughty et al., 2018](#)).

## C Use of Scientific Artifacts & Licensing

Our work uses the RCT rejection sampling dataset by [Keith et al. \(2023\)](#). In particular, the dataset is fully in English, containing publicly available paper titles and abstracts. The authors remove any potentially personally identifiable information from the dataset (author names, user ids, user IP addresses, or session ids). The dataset is made publicly available for research purposes (apache-2.0).

Finally, **DoubleLingo** uses the Hugging Face implementations for *bert-base-uncased*, *allenai/specter*, and *all-mpnet-base-v2*, all made publicly available for research purposes (apache-2.0).