

# Uninformative Input Features and Counterfactual Invariance: Two Perspectives on Spurious Correlations in Natural Language

Anonymous ACL submission

## Abstract

The natural language processing community has become increasingly interested in spurious correlations, and in methods for identifying and eliminating them. Gardner et al. (2021) argue that due to the compositional nature of language, *all* correlations between labels and individual input features are spurious. This paper analyzes this proposal in the context of a toy example, demonstrating three distinct conditions that can give rise to feature-label correlations through a simple PCFG. Linking the toy example to a structured causal model shows that (1) feature-label correlations can arise even when the label is invariant to interventions on the feature, and (2) feature-label correlations may be absent even when the label *is* sensitive to interventions on the feature. Because input features will be individually correlated with labels except in very rare circumstances, mitigation and stress tests should focus on those correlations that are counterfactually invariant under plausible causal models.

## 1 Introduction

Spurious correlations have increasingly preoccupied researchers in machine learning (Geirhos et al., 2020) and related fields, including natural language processing (Gururangan et al., 2018; McCoy et al., 2019, *inter alia*). However, the notion is frequently used without a formal definition. Gardner et al. (2021) propose a definition in terms of conditional probabilities: a feature  $X_i$  is spuriously correlated with the label  $Y$  unless  $P(Y | X_i)$  is uniform. The definition can be generalized from uniformity to independence ( $X_i \perp\!\!\!\perp Y$ ) without affecting the claims of the paper. They go on to argue that “in a language understanding problem, . . . *all* simple correlations between input features and output labels are spurious” (emphasis in the original). The property that individual input features should be independent of labels — which we will call *marginally uninformative input features (UIF)* — is treated as

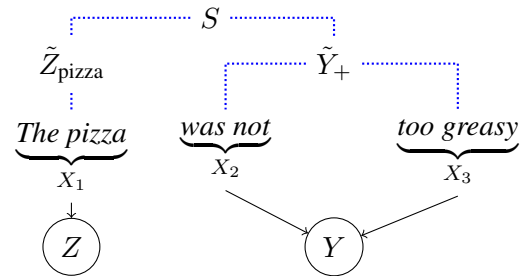


Figure 1: An instance from the toy model. The upper part of the figure corresponds to  $f_X$ , the function that generates the text via a PCFG (see fig. 2): nodes correspond to non-terminals in the grammar and edges represent context-free derivations. The lower part of the figure corresponds to the causal model of the sentiment  $Y$  and target  $Z$ . Here nodes correspond to variables and edges correspond to causal relationships.

an assumption about the nature of language processing and also as a desideratum that datasets should satisfy: if the label can be predicted from input features alone, then the dataset is too easy.<sup>1</sup>

The principle of UIF is based on the insight that linguistic context can invert the semantics of any subspan of a text (via, e.g., syntactic negation or discourse relations). Furthermore, the frequency of negation and other forms of semantic inversion may vary across datasets and deployment settings. A predictor that relies on, e.g., negation being rare, cannot be said to have truly achieved competence in the language processing task, and may perform poorly in domains in which these high-level distributional properties shift.

An especially provocative assertion of Gardner et al. is that all correlations between labels and individual input features have the same status. In the sentence *the pizza was amazing*, sup-

<sup>1</sup>To formalize the UIF assumption, it is necessary to clarify which features are “input features”: bytes, phonemes, word-pieces, words, phrases, or sentences? The selection of input features is a property of the model and not the dataset; one could use character-level features for natural language inference or sentence-level features for sentiment analysis.

pose that both *pizza* and *amazing* are correlated with positive sentiment because the reviewers like pizza. There is an intuitive difference between these two correlations, because the modified sentence *the movie was amazing* should have the same label as the original, while *the pizza was greasy* should not. This intuition can be formalized using the framework of causality, which has generally treated spurious correlations as those that arise without a direct causal explanation (Simon, 1954). Given a causal model of the data generating process, we can compute the *interventional* distribution  $P(Y \mid \text{do}(X_1 := x_1), X_2, X_3)$ , which corresponds to the distribution over  $Y$  in a data generating process in which the value of  $X_1$  is surgically set to  $x_1$  (Pearl, 1995; Peters et al., 2017; Feder et al., 2021).<sup>2</sup> When such interventions do not affect  $Y$  for any given example, we say that  $Y$  and  $X_1$  are *counterfactually invariant* (Veitch et al., 2021). Violations of UIF are particularly troubling when they are accompanied by counterfactual invariance, because non-causal correlations often do not transfer to other domains (Schölkopf et al., 2012; Bühlmann, 2020).

This paper uses a toy example to relate the UIF property to (1) the production probabilities in probabilistic context-free grammars (PCFGs), and (2) counterfactual invariance in structured causal models. The connection to PCFGs provides additional motivation for the UIF criterion from the perspective of domain generalization, while clarifying the scenarios that can give rise to violations of UIF, which Gardner et al. attribute too narrowly to “bias and priming effects” in annotators. The connection to counterfactual invariance highlights the ways in which these concepts do and do not align. Efforts to remove artifacts from the training and evaluation of NLP systems will be most productive when focused at the intersection of these two views of spurious correlations: violations of UIF for input features to which the label is counterfactually invariant according to a causal model of the data.

## 2 Toy Example

Consider a simplified targeted sentiment analysis task (Mitchell et al., 2013), in which the sentiment is  $Y$ , the target is  $Z$ , and the sentences are all of the form  $(X_1, X_2, X_3)$ , with  $X_1$  specifying a target noun phrase,  $X_2$  a copula-like expression, and

<sup>2</sup>Space does not permit a discussion of the distinction between interventions and counterfactuals (see Pearl, 2009).

$$\begin{aligned}
 U &:= N_U & (1) \\
 (X_1, X_2, X_3) &:= f_X(U, N_X) & (2) \\
 Z &:= f_Z(X_1, N_Z) & (3) \\
 Y &:= f_Y(X_2, X_3, N_Y). & (4)
 \end{aligned}$$

Figure 2: Causal model for the toy example shown in fig. 1.  $N_U, N_X, N_Y, N_Z$  indicate independent noise variables, and  $f_X, f_Y, f_Z$  indicate deterministic functions that map from causes to effects (for more details on the notation, see Peters et al., 2017).

$X_3$  a predicative adjectival phrase. For example,  $Y = \text{POS}$ ,  $Z = \text{PIZZA}$ ,  $X_1 = \text{the pizza}$ ,  $X_2 = \text{turned out to be}$ ,  $X_3 = \text{crispy and delicious}$ . We will treat this data as generated from the causal model shown in fig. 2. This causal model can be summarized by two assertions: (1) the target  $Z$  is a direct effect of only the span  $X_1$ ; (2) the sentiment label  $Y$  is a direct effect of only the spans  $X_2$  and  $X_3$ . The function  $f_X$  can represent any generative model of text: an n-gram model, a grammar-based formalism, a deep autoregressive network, etc.

**Aside on the direction of causation.** We treat the text as the cause of the labels, rather than the converse. This distinction is somewhat vexed (Schölkopf et al., 2012; Jin et al., 2021). In some cases the direction of causation is clear from the task (e.g., table-to-text generation, summarization, and translation), but often the problem could be framed in either direction: perhaps the writer had the label in mind when producing the text, and thus the text is an effect of the label; or perhaps it is better to think of the annotator, who must read the text to arrive at the label, regardless of the writer’s intentions. When the labels cause the text, the notion of counterfactual invariance can be restated in terms of the invariance of text features to perturbations on labels, e.g.  $P(X_1 \mid \text{do}(Y := y), Z)$ . As the toy example is meant to serve only an expository purpose, we leave elaboration of the relationship of UIF to such models for future work.

### 2.1 Counterfactual invariance $\nRightarrow$ UIF

The causal model implies several counterfactual invariance properties: intervention on  $X_1$  will not affect  $Y$ , nor will intervention on  $X_2$  or  $X_3$  affect  $Z$ . This is because  $X_1$  blocks the influence of  $X_2$  and  $X_3$  on  $Z$ , and vice versa for  $Y$ . Conversely,  $(X_3, Y)$  are not counterfactually invariant in gen-

eral because  $X_3$  is an ancestor of  $Y$  in the causal graph, and similarly for  $(X_2, Y)$  and  $(X_1, Z)$ .

Counterfactual invariance does not imply that the associated input features are marginally uninformative of the label. Consider a classical spurious correlation in which pizza tends to receive positive sentiment and sushi receives negative sentiment. This correlation is produced when  $f_X$  encodes a PCFG with the top-level production:

$$\begin{aligned}
 S \rightarrow \tilde{Z}_{\text{pizza}} \tilde{Y}_+ & (1 + \alpha)/4 \\
 & \tilde{Z}_{\text{sushi}} \tilde{Y}_- & (1 + \alpha)/4 \\
 & \tilde{Z}_{\text{pizza}} \tilde{Y}_- & (1 - \alpha)/4 \\
 & \tilde{Z}_{\text{sushi}} \tilde{Y}_+ & (1 - \alpha)/4,
 \end{aligned} \tag{5}$$

with the right column indicating the probability of each rule expansion and  $\alpha \in [-1, 1]$ .<sup>3</sup> The non-terminal symbols  $\tilde{Z}_{\text{pizza}}$ ,  $\tilde{Z}_{\text{sushi}}$ ,  $\tilde{Y}_+$ ,  $\tilde{Y}_-$  are intentionally chosen to correspond to the labels  $Z$  and  $Y$ . Subsequent rules in the grammar can then be designed to ensure that  $\tilde{Z}_{\text{pizza}}$  usually produces values of  $X_1$  that make  $Z = \text{PIZZA}$  likely, and analogously for the other non-terminals and associated labels. The unification of PCFGs and structured causal models is shown in fig. 1.

When  $\alpha \neq 0$ , there may be an association between  $X_1$  and  $(X_2, X_3)$ . As a result, there exist  $(x_1, x'_1)$  such that,

$$\begin{aligned}
 P(Y|X_1 = x_1) & \\
 &= \sum_{X_2, X_3} P(Y | X_2, X_3) P(X_2, X_3 | X_1 = x_1) \\
 &\neq \sum_{X_2, X_3} P(Y | X_2, X_3) P(X_2, X_3 | X_1 = x'_1) \\
 &= P(Y|X_1 = x'_1),
 \end{aligned} \tag{6}$$

creating a violation of UIF. The same argument can be applied to  $P(Z | X_2)$  and  $P(Z | X_3)$ . UIF is also violated in  $P(Z | X_1)$ ,  $P(Y | X_2)$ , and  $P(Y | X_3)$ , but for a different reason: these distributions are conditioned on the direct causal parents of the labels in  $f_Y$  and  $f_Z$ . Manipulation of the data distribution to ensure that  $\alpha = 0$  (deconfounding  $\tilde{Y}$  and  $\tilde{Z}$ ) can remove only the violations of UIF induced by  $f_X$ , but not those induced by the direct causal relationships encoded in  $f_Y$  and  $f_Z$ .

<sup>3</sup>The stochasticity of the grammar is encoded in the deterministic function  $f_X$  through the noise variable  $N_X$ . Let  $N_X \sim \text{Uniform}(0, 1)$ , and choose the first rule expansion of  $S$  when  $N_X < (1 + \alpha)/4$ , the second rule expansion when  $(1 + \alpha)/4 \leq N_X < (1 + \alpha)/2$ , and so on.

**Discussion.** The example shows how violations to UIF can emerge via confounding, creating classical spurious correlations in the sense of Simon (1954): informativeness despite counterfactual invariance. Such correlations are unlikely to be robust because it is not difficult to imagine a domain in which the sign of  $\alpha$  changes, impairing the performance of predictors that have learned the spurious correlation. In contrast, feature-label correlations that arise directly from the causal model, such as  $(Z, X_1)$ , are only damaging under more extreme forms of concept shift, in which the meanings of the features themselves change.<sup>4</sup>

## 2.2 UIF $\nRightarrow$ Counterfactual Invariance

Violations of counterfactual invariance can occur even when UIF is satisfied. To show this, we supply two more productions for the grammar:

$$\begin{aligned}
 \tilde{Y}_+ \rightarrow \text{COP}_+ \text{ADJP}_+ & \beta_+ \\
 & \text{COP}_- \text{ADJP}_- & 1 - \beta_+
 \end{aligned} \tag{7}$$

$$\begin{aligned}
 \tilde{Y}_- \rightarrow \text{COP}_+ \text{ADJP}_- & \beta_- \\
 & \text{COP}_- \text{ADJP}_+ & 1 - \beta_-
 \end{aligned} \tag{8}$$

Here the non-terminal  $\text{COP}_+$  produces a ‘‘positive’’ copula in  $X_2$  (*is, was, is universally agreed to be*),  $\text{COP}_-$  produces a negated copula in  $X_2$  (*isn’t, wasn’t, was the furthest possible thing from*),  $\text{ADJP}_+$  produces positive-sentiment adjectival phrases in  $X_3$  (*great, delicious*), and  $\text{ADJP}_-$  produces negative-sentiment adjectival phrases in  $X_3$  (*disappointing, totally unappetizing*). There are two special cases of interest:

- When  $\beta_+ = \beta_-$ , the probability of using a negated copula is independent of  $Y$ , so  $X_2$  satisfies UIF with regard to  $Y$ , while  $X_3$  generally does not.
- When  $\beta_+ = 1 - \beta_-$ , the use of negation is balanced to make the distribution over sentiment terms independent of  $Y$ , so  $X_3$  satisfies UIF with  $Y$ , while  $X_2$  generally does not.

<sup>4</sup>This basic intuition is sometimes formalized as the *principle of sparse mechanism shift*, which states that complex causal systems are usually composed of smaller independent parts, and that domain shifts typically affect only a few components (Schölkopf et al., 2021). A related principle arises in the context of natural language: distributional frequencies are more likely to change across domains than categorical facts about language. Biber (1991), for example, makes this argument explicitly in the analysis of register. In our model, the implication is that the probabilistic rule expansions in  $f_X$  are more likely to change than the basic properties of the lexicon, which govern which terminal symbols can be emitted by each non-terminal.

216 Combining these cases, both  $X_2$  and  $X_3$  satisfy  
217 UIF with  $Y$  when  $\beta_+ = \beta_- = \frac{1}{2}$ , meaning that  
218 negated and non-negated copula are equally likely  
219 and are independent of  $Y$ .

220 **Discussion.** UIF is violated not only by con-  
221 founding, as discussed in the previous section, but  
222 also in mild settings that do not meet any reason-  
223 able definition of bias: unless  $\beta_+ = \beta_- = 1/2$   
224 then at least one of  $X_2$  and  $X_3$  is marginally infor-  
225 mative of  $Y$ . Furthermore, UIF has no impact on  
226 the counterfactual invariance of  $X_2$  and  $X_3$  on  $Y$ .  
227 Neither is counterfactually invariant even when the  
228 generative model is parametrized to make UIF hold  
229 for all input features (see also Pearl, 2009, page  
230 185). This is because the overall sentiment can  
231 be directly affected by adding or removing nega-  
232 tion and by flipping the polarity of the sentiment-  
233 carrying adjective.

### 234 3 Conclusions

235 In the toy example, violations of UIF arise from  
236 three distinct phenomena: confounding between  
237 the sentiment and the target ( $\alpha \neq 0$ , leading to  
238  $X_1 \not\perp Y$ ); confounding between the sentiment  
239 and the use of negation ( $\beta_+ \neq \beta_-$ , leading to  
240  $X_2 \not\perp Y$ ); and lack of a perfect balance in the prob-  
241 ability of negation between positive- and negative-  
242 sentiment examples ( $\beta_+ \neq 1 - \beta_-$ , leading to  
243  $X_3 \not\perp Y$ .) The conditions required to satisfy UIF  
244 are thus progressively less plausible as we move  
245 from  $X_1$  to  $X_3$ , and full UIF is achieved only in the  
246 perfectly balanced case of  $\alpha = 0, \beta_+ = \beta_- = \frac{1}{2}$ .  
247 The number of such constraints will increase with  
248 the size of the grammar, making UIF vanishingly  
249 rare in more general settings. Note that this general  
250 conclusion follows from the PCFG analysis, and  
251 can be derived without reference to causality.

252 The toy example also demonstrates the discon-  
253 nect between the UIF view of spurious correlations  
254 and the causal view: counterfactual invariance does  
255 not imply UIF because  $X_1$  can be marginally infor-  
256 mative of  $Y$  even when  $X_1$  and  $Y$  are counter-  
257 factually invariant (these are the artifacts that we  
258 want to remove); UIF does not imply counterfactual  
259 invariance because both  $X_2$  and  $X_3$  can be  
260 uninformative of  $Y$  even when  $Y$  is sensitive to  
261 interventions on both features. From a theoretical  
262 perspective, it is unsurprising that these two views  
263 diverge, because UIF is a purely observational cri-  
264 terion while counterfactual invariance requires an  
265 explicit causal model. Indeed, this relationship is

266 discussed in depth by Pearl (2009, §6.3), albeit out-  
267 side the context of language. The two perspectives  
268 can be seen as complementary, in that violation of  
269 UIF is a necessary but insufficient condition for a  
270 spurious correlation in the causal sense.

271 It is of course possible to quibble with the causal  
272 model presented here, and in real applications it is  
273 likely impractical to construct full causal models  
274 of language. How then can we use causal insights  
275 to go beyond sensitivity analysis to design better  
276 benchmarks and more robust language understand-  
277 ing systems? In some cases it is possible to elabo-  
278 rate partial causal models of a task, with associated  
279 invariance properties: for example, the sentiment of  
280 a movie review should be invariant to (though not  
281 independent of) the identities of the actors in the  
282 movie. Several existing approaches can be viewed  
283 as instantiations of partial causal models: for ex-  
284 ample, data augmentation, causally-motivated reg-  
285 ularizers, stress tests, and “worst-subgroup” per-  
286 formance metrics (and associated robust optimizers)  
287 can be seen as enforcing or testing task-specific in-  
288 variance properties that provide robustness against  
289 known distributional shifts (e.g., Lu et al., 2020;  
290 Ribeiro et al., 2020; Kaushik et al., 2021; Koh et al.,  
291 2021; Veitch et al., 2021). Such approaches gener-  
292 ally require domain knowledge about the linguistic  
293 and causal properties of the task at hand — or to  
294 put it more positively, they make it possible for  
295 such domain knowledge to be brought to bear.

296 A final observation, pertaining to both UIF and  
297 counterfactual invariance, is the parallel treatment  
298 of  $X_2$  (the copula) and  $X_3$  (the adjectival phrase).  
299 From a lexical semantic perspective, only  $X_3$  is  
300 directly associated with the sentiment, while  $X_2$   
301 plays a functional role by potentially reversing  $X_3$ .  
302 It may therefore seem undesirable to learn a corre-  
303 lation between  $X_2$  and  $Y$ , and preferable to attach  
304 that relationship exclusively to  $X_3$ . Yet neither UIF  
305 nor counterfactual invariance is capable of making  
306 such a distinction. While it is possible to enforce  
307 uninformative on  $X_2$  heuristically, e.g. by sam-  
308 pling or augmenting the data to ensure  $\beta_+ = \beta_-$ ,  
309 those same heuristics could be applied to enforce  
310 uninformative on  $X_3$  by making  $\beta_+ = 1 - \beta_-$ .  
311 Singling out  $X_2$  requires additional justification.  
312 Such a principle might be found in the multitask  
313 setting, in which we prefer feature-label informa-  
314 tiveness to be sparse, with each feature directly  
315 informing only a few labels.



316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371

## References

Douglas Biber. 1991. *Variation across speech and writing*. Cambridge University Press.

Peter Bühlmann. 2020. Invariance, causality and robustness. *Statistical Science*, 35(3):404–426.

Amir Feder, Katherine A. Keith, Emaad Manzoor, Reid Pryzant, Dhanya Sridhar, Zach Wood-Doughty, Jacob Eisenstein, Justin Grimmer, 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*.

Matt Gardner, William Merrill, Jesse Dodge, Matthew Peters, Alexis Ross, Sameer Singh, and Noah A. Smith. 2021. [Competency problems: On finding and removing artifacts in language data](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1801–1813, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. 2020. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673.

Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and Noah A Smith. 2018. Annotation artifacts in natural language inference data. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 107–112.

Zhijing Jin, Julius von Kügelgen, Jingwei Ni, Tejas Vaidhya, Ayush Kaushal, Mrinmaya Sachan, and Bernhard Schölkopf. 2021. [Causal direction of data collection matters: Implications of causal and anticausal learning for NLP](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 9499–9513, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Divyansh Kaushik, Amrith Setlur, Eduard Hovy, and Zachary C Lipton. 2021. Explaining the efficacy of counterfactually augmented data. In *International Conference on Learning Representations*.

Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al. 2021. Wilds: A benchmark of in-the-wild distribution shifts. In *International Conference on Machine Learning*, pages 5637–5664. PMLR.

Kaiji Lu, Piotr Mardziel, Fangjing Wu, Preetam Amancharla, and Anupam Datta. 2020. Gender bias in

neural natural language processing. In *Logic, Language, and Security*, pages 189–202. Springer. 372  
373

Tom McCoy, Ellie Pavlick, and Tal Linzen. 2019. Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3428–3448. 374  
375  
376  
377  
378

Margaret Mitchell, Jacqui Aguilar, Theresa Wilson, and Benjamin Van Durme. 2013. Open domain targeted sentiment. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1643–1654. 379  
380  
381  
382  
383

Judea Pearl. 1995. Causal diagrams for empirical research. *Biometrika*, 82(4):669–688. 384  
385

Judea Pearl. 2009. *Causality*. Cambridge university press. 386  
387

Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. 2017. *Elements of causal inference: foundations and learning algorithms*. The MIT Press. 388  
389  
390  
391

Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. [Beyond accuracy: Behavioral testing of NLP models with CheckList](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4902–4912, Online. Association for Computational Linguistics. 392  
393  
394  
395  
396  
397  
398

Bernhard Schölkopf, Dominik Janzing, Jonas Peters, Eleni Sgouritsa, Kun Zhang, and Joris Mooij. 2012. On causal and anticausal learning. In *Proceedings of the 29th International Conference on Machine Learning*, pages 459–466. 399  
400  
401  
402  
403

Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner, Anirudh Goyal, and Yoshua Bengio. 2021. Toward causal representation learning. *Proceedings of the IEEE*, 109(5):612–634. 404  
405  
406  
407  
408

Herbert A Simon. 1954. Spurious correlation: A causal interpretation. *Journal of the American statistical Association*, 49(267):467–479. 409  
410  
411

Victor Veitch, Alexander D’Amour, Steve Yadlowsky, and Jacob Eisenstein. 2021. Counterfactual invariance to spurious correlations in text classification. *Advances in Neural Information Processing Systems*, 34. 412  
413  
414  
415  
416