

# Identifying Spurious Correlations for Robust Text Classification

**Zhao Wang**

Illinois Institute of Technology  
zwang185@hawk.iit.edu

**Aron Culotta**

Tulane University  
aculotta@tulane.edu

## Abstract

The predictions of text classifiers are often driven by spurious correlations – e.g., the term *Spielberg* correlates with positively reviewed movies, even though the term itself does not semantically convey a positive sentiment. In this paper, we propose a method to distinguish spurious and genuine correlations in text classification. We treat this as a supervised classification problem, using features derived from treatment effect estimators to distinguish spurious correlations from “genuine” ones. Due to the generic nature of these features, and their small dimensionality, we find that the approach works well even with limited training examples, and that it is possible to transport the word classifier to new domains. Experiments on four datasets (sentiment classification and toxicity detection) suggest that using this approach to inform feature selection also leads to more robust classification, as measured by improved worst-case accuracy on the samples affected by spurious correlations.

## 1 Introduction

Text classifiers often rely on spurious correlations. For example, consider sentiment classification of movie reviews. The term *Spielberg* may be correlated with the positive class because many of director Steven Spielberg’s movies have positive reviews. However, the term itself does not indicate a positive review. In other words, the term *Spielberg* does not *cause* the review to be positive. Similarly, consider the problem of toxicity classification of online comments. Terms indicative of certain ethnic groups may be associated with the toxic class because those groups are often victims of harassment, not because those terms are toxic themselves.

Oftentimes, such spurious correlations do not harm prediction accuracy because the same correlations exist in both training and testing data

(under the common assumption of i.i.d. sampling). However, they can still be problematic for several reasons. For example, under *dataset shift* (Quionero-Candela et al., 2009), the testing distribution differs from the training distribution. E.g., if Steven Spielberg makes a new, bad movie, the sentiment classifier may incorrectly classify the reviews as positive because they contain the term *Spielberg*. Additionally, if the spurious correlations indicate demographic attributes, then the classifier may suffer from issues of *algorithmic fairness* (Kleinberg et al., 2018). For example, the toxicity classifier may unfairly over-predict the toxic class for comments discussing certain demographic groups. Finally, in settings where classifiers must explain their decisions to humans, such spurious correlations can reduce trust in autonomous systems (Guidotti et al., 2018).

In this paper, we propose a method to distinguish *spurious* correlations, like *Spielberg*, from *genuine* correlations, like *wonderful*, which more reliably indicate the class label. Our approach is to treat this as a separate classification task, using features drawn from treatment effect estimation approaches that isolate the impact each word has on the class label, while controlling for the context in which it appears.

We conduct classification experiments with four datasets and two tasks (sentiment classification and toxicity detection), focusing on the problem of *short text* classification (i.e., single sentences or tweets). We find that with a small number of labeled word examples (200-300), we can fit a classifier to distinguish spurious and genuine correlations with moderate to high accuracy (.66-.82 area under the ROC curve), even when tested on the terms most strongly correlated with the class label. In addition, due to the generic nature of the features, we find that we can train a word classifier on one domain and transfer it to another domain without

much loss in accuracy.

Finally, we apply the word classifier to inform feature selection for the original classification task (e.g., sentiment classification and toxicity detection). Following recent work on distributionally robust classification (Sagawa et al., 2020a), we measure worst-case accuracy by considering samples of data most affected by spurious correlations. We find that removing terms in order of their predicted probability of being spurious correlations can result in more robust classification with respect to this worst-case accuracy.

## 2 Problem and Motivation

We consider binary classification of short documents, e.g., sentences or tweets. Each sentence is a sequence of words  $s = \langle w_1 \dots w_k \rangle$  with a corresponding binary label  $y \in \{-1, 1\}$ . To classify a sentence  $s$ , it is first transformed into a feature vector  $x$  via a feature function  $g : s \mapsto x$ . Then, the feature vector is assigned a label by a classification function  $f : (x; \theta) \mapsto \{-1, 1\}$ , with model parameters  $\theta$ . Parameters  $\theta$  are typically estimated from a set of i.i.d. labeled examples  $\mathcal{D} = \{(s_1, y_1) \dots (s_n, y_n)\}$  by minimizing some loss function  $\mathcal{L} : \theta^* \leftarrow \arg \min_{\theta} \mathcal{L}(\mathcal{D}, \theta)$ .

To illustrate the problem addressed in this paper, we will first consider the simple approach of a bag-of-words logistic regression classifier. In this setting, the feature function  $g(s)$  simply maps a document to a word count vector  $x = \{x_1 \dots x_V\}$ , for vocabulary size  $V$ , and the classification function is the logistic function  $f(x; \theta) = \frac{1}{1+e^{-(x, \theta)}}$ . After estimating parameters  $\theta$  on labeled data  $\mathcal{D}$ , we can then examine the coefficients corresponding to each word in the vocabulary to see which words are most important in the model.

In Figure 1, we show eight words with high magnitude coefficients for a classifier fit on a dataset of movie reviews (Pang and Lee, 2005), where class 1 means positive sentiment and  $-1$  means negative sentiment. We will return shortly to the meaning of the  $x$ -axis; for now, let us consider the  $y$ -axis, which is the estimated coefficient  $\theta_w$  for each word. Of the four words strongly correlated with the positive class ( $\theta_w > 0$ ), two seem genuine (*enjoyable*, *masterpiece*), while two seem spurious (*animated*, *spielberg*). (Steven Spielberg is very successful American director and producer.) Similarly, of the words correlated with the negative class, two seem genuine (*boring*, *failure*) and two seem spurious

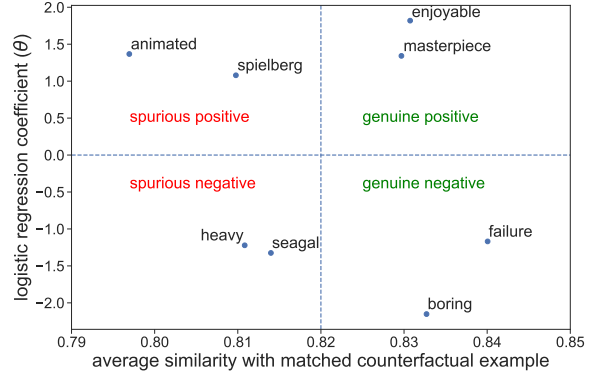


Figure 1: Motivating example of spurious and genuine correlations in a sentiment classification task.

(*heavy*, *seagal*). (Steven Seagal is an American actor mostly known for martial-arts movies.) Furthermore, in some cases the spurious term actually has a *larger* coefficient than the genuine term (e.g., *seagal* versus *failure*).

Our goal in this paper is to distinguish between spurious and genuine correlations. Without wading into long-standing debates over the nature of causality (Aldrich et al., 1995), we simplify the distinction between genuine and spurious correlations as a dichotomous decision: the discovered relationship between word  $w$  and label  $y$  is genuine if, all else being equal, one would expect  $w$  to be a determining factor in assigning a label to a sentence. We use human annotators to make this distinction for training and evaluating models.

In this light, our problem is related to prior work on active learning with rationales (Zaidan et al., 2007; Sharma et al., 2015) and interactive feature selection (Raghavan et al., 2005). However, our goal is not solely to improve prediction accuracy, but also to improve robustness across different groups affected by these spurious correlations.

## 3 Methods

Our definition of *genuine* correlation given above fits well within the counterfactual framework of causal inference (Winship and Morgan, 1999). If word  $w$  in  $s$  were replaced with some other word  $w'$ , how likely is it that the label  $y$  would change? Since conducting randomized control trials to answer this counterfactual for many terms and sentences is infeasible, we instead resort to matching methods, commonly used to estimate average treatment effects from observational data (Imbens, 2004; King and Nielsen, 2019). The intuition is as follows: if  $w$  is a reliable piece of evidence to

determine the label of  $s$ , we should be able to find a very similar sentence  $s'$  that (i) does not contain  $w$ , and (ii) has the opposite label of  $s$ .

Paul (2017) proposed a similar formulation, using propensity score matching to estimate the treatment effect for each term, then performing feature selection based on these estimates. Beyond recent critiques of propensity scores (King and Nielsen, 2019), any matching approach will create matches of varying quality, making it difficult to distinguish between spurious and genuine correlations. Returning to Figure 1, the  $x$ -axis shows the average quality of the counterfactual match for each term, where a larger value means that the linguistic context of the counterfactual sentence is very similar to the original sentence. (These are computed by cosine similarity of sentence embeddings, described in §3.2.) Even though these terms have very similar average treatment effect estimates, the quality of the match seems to be a viable signal of whether the term is spurious or genuine.

More generally, building on prior work that treats causal inference as a classification problem (Lopez-Paz et al., 2015), we can derive a number of features from the components of the treatment effect estimates (enumerated in §3.3), and from these fit a classification model to determine whether a word should be labeled as spurious or genuine. This *word classifier* can then be used in a number of ways to improve the document classifier (e.g., to inform feature selection, to place priors on word coefficients, etc.).

To build the word classifier, we assume a human has annotated a small number of terms as spurious or genuine, which we can use as training data. While this places an additional cost on annotation, the nature of the features reduces this burden — there are not very many features in the word classifier, and they are mostly generic / domain independent. As a result, in experiments below we find that useful word classifiers can be built from a small number of labeled terms (200-300). Furthermore, and perhaps more importantly, we find that the word classifier can be transported to new domains with little loss in accuracy. This suggests that one can label words once in one domain, fit a word classifier, and apply it in new domains without annotating additional words.

### 3.1 Overview of approach

The main stages of our approach are as follows:

1. Given training data  $\mathcal{D} = \{(s_1, y_1) \dots (s_n, y_n)\}$  for the primary classification task, fit an initial classifier  $f(x; \theta)$ .
2. Extract from  $f(x; \theta)$  the words  $\mathcal{W} = \{w_1 \dots w_m\}$  that are most strongly associated with each class according to the initial classifier. E.g., for logistic regression, we may extract the words with the highest magnitude coefficients for each class. For more complex models, other transparency algorithms may be used (Martens and Provost, 2014).
3. For each word, compute features that indicate its likelihood to be spurious or genuine (§3.3).
4. Fit a word classifier  $h(w; \lambda)$  on a human-annotated subset of  $\mathcal{W}$ .
5. Apply  $h(w; \lambda)$  on remaining words to estimate the probability that they are spurious correlations.

After the final step, one may use the posterior probabilities in a number of ways to improve classification: e.g., to sort terms for feature selection, to place priors on word coefficients, to set attention weights in neural networks, etc. In this paper, we focus on feature selection, leaving other options for future work.

Additionally, we experiment with domain adaptation settings, where  $h(w; \lambda)$  is fit on one domain and applied to another domain for feature selection, without requiring additional labeled words from that domain.

### 3.2 Matching

Most of the features for the word classifier are inspired by matching approaches from causal inference (Stuart, 2010). The idea is to match sentences containing different words in similar contexts so that we can isolate the effect that one particular word choice has on the class label.

For a word  $w$  and a sentence  $s$  containing this word, we let  $s[\hat{w}]$  be the sentence  $s$  with word  $w$  removed. The goal of matching is to find some other context  $s'[\hat{w}']$  such that  $w \notin s'$  and  $s[\hat{w}]$  is semantically similar to  $s'[\hat{w}']$ . We use a *best match* approach, finding the closest match  $s^* \leftarrow \arg \max_{s'} \text{sim}(s[\hat{w}], s'[\hat{w}'])$ . With this best match, we can compute measures such as the average treatment effect (ATE) of word  $w$  in  $N$  sentences:

$$\tau_w = \frac{1}{N} \sum_{\{s|w \in s\}} y_s - y_{s^*}$$

it's <b>refreshing</b> to see a movie that (1)
it's <b>rare</b> to see a movie that (-1)
cast has a lot of <b>fun</b> with the material (1)
comedy with a lot of <b>unfunny</b> (-1)
smoothly under the direction of <b>spielberg</b> (1)
it works under the direction of <b>kevin</b> (1)
refreshingly different slice of asian <b>cinema</b> (1)
an interesting slice of <b>history</b> (1)
charting the rise of hip-hop <b>culture</b> in general (1)
hip-hop has a <b>history</b> , and it's a metaphor (1)

Table 1: Examples of matched contexts from IMDB dataset; word substitutions are shown in bold.

Thus, a term  $w$  will have a large value of  $\tau_w$  if (i) it often appears in positive sentences, and (ii) very similar sentences where  $w$  is swapped with  $w'$  have negative labels.

In our experiments, to improve the quality of matches, we limit contexts to the five previous and five subsequent words to  $w$ , then represent the context by concatenating the last four layers of a pre-trained BERT model (Devlin et al., 2018). We use the cosine similarity of context embeddings as a measure of semantic similarity.

Take one example from Table 1: “it’s *refreshing* to see a movie that (1)” is matched with “it’s *rare* to see a movie that (-1)”. Words *refreshing* and *rare* appear in similar contexts, but adding *refreshing* to this context makes the sentence positive, while adding *rare* to this context makes it negative. If most of the pairwise matches show that adding *refreshing* is more positive than adding other substitution words, then *refreshing* is very likely to be a genuine positive word.

On the contrary, if adding other substitution words for similar contexts does not change the label, then  $w$  is likely to be a spuriously correlated word. Take another example from Table 1, “smoothly under the direction of *spielberg* (1)” is matched with “it works under the direction of *kevin* (1)”, *spielberg* and *kevin* appear in similar contexts, and substituting *spielberg* with *kevin* does not make any difference in the label. If most pairwise matches show that substituting *spielberg* to other words does not change the label, then *spielberg* is very likely to be a spurious positive word.

### 3.3 Features for Word Classification

While the matching approach above is a traditional way to estimate the causal effect of a word  $w$  given observational data, there are many well known limitations to matching approaches. A primary difficulty is that high quality matches may not exist

in the data, leading to biased estimates. Inspired by supervised learning approaches to causal inference (Lopez-Paz et al., 2015), rather than directly use the ATE to distinguish between spurious and genuine correlations, we instead compute a number of features to summarize information about the matching process. In addition to the ATE itself, we calculate the following features:

- The average context similarity of every match for word  $w$ .
- The context similarity of the top-5 most closest matches.
- The maximum and standard deviation of the similarity score.
- The context similarity of closest positive and negative sentences.
- The weighted average treatment effect, where Eq. 3.2 is weighted by the similarity between  $s$  and  $s^*$ .
- The ATE restricted to the top-5 most similar matches for sentences containing  $w$ .
- The word’s coefficient in the initial sentence classifier.
- Finally, to capture subtle semantic differences between the original and matched sentences, we compute features such as the average embedding difference from all matches, the top-3 most different dimensions from the average embedding, and the maximum value along each dimension.

### 3.4 Measuring the Impact of Spurious Correlations on Classification

After we train the word classifier to identify spurious and genuine words, we are further interested in exploring how spurious correlations affect classification performance on test data. As discussed in §1, measuring robustness can be difficult when data are sampled i.i.d. because the same spurious correlations exist in the training and testing data. Thus, we would not expect accuracy to necessarily improve on a random sample when spurious words are removed. Instead, we are interested in measuring the *robustness* of the classifier, where robustness is with respect to which subgroup of data is being considered.

Motivated by (Sagawa et al., 2020a), we divide the test data into two groups and explore the model performance on each. The first group, called the *minority group*, contains sentences in which the spurious correlation is expected to mislead the clas-



sifier. From our running example, that would be a *negative* sentiment sentence containing *spielberg*, or a *positive* sentiment sentence containing *seagal*. Analogously, the *majority group* contains examples in which the spurious correlation helps the classifier (e.g., positive sentiment documents containing *spielberg*). In §4.4, we conduct experiments to see how removing terms that are predicted to be spurious affects accuracies on majority and minority groups.

## 4 Experiments

### 4.1 Data

We experiment with four datasets for two binary classification tasks: sentiment classification and toxicity detection.<sup>1</sup>

- *IMDB movie reviews*: movie review sentences labeled with their overall sentiment polarity (positive or negative) (Pang and Lee, 2005) (version 1.0).
- *Kindle reviews*: product reviews from Amazon Kindle Store with ratings range from 1-5 (He and McAuley, 2016). We first fit a binary classification model on this dataset to identify keywords, and then split each review into single sentences and assign each sentence the same rating as the original review. We select sentences that contain the keywords and then remove sentences that have fewer than 5 or more than 40 words, and finally label the remaining sentences rated {4,5} as positive and sentences rated {1,2} as negative.
- *Toxic comment*: a dataset of comments from Wikipedia’s talk page.<sup>2</sup> Comments are labeled by human raters for toxic behavior (e.g., comments that are rude, disrespectful, offensive or otherwise likely to make someone leave a discussion). Each comment was shown up to 10 annotators and the fraction of human raters who believed the comment is toxic serve as the final toxic score that ranges from 0.0 to 1.0. We follow the same processing steps in *Kindle reviews dataset*: split comments into sentences, select sentences using keywords, and limit sentence length. We label sentences with toxicity score  $\geq 0.7$  as toxic and  $\leq 0.5$  as non-toxic.

<sup>1</sup>All code and data will be shared publicly upon publication.

<sup>2</sup><https://www.kaggle.com/c/jigsaw-unintended-bias-in-toxicity-classification>

	#docs	#top words	#matched sentences
IMDB	10,662	366	8,882
Kindle	20,232	270	24,882
Toxic comment	15,216	329	84,14
Toxic tweet	6,774	341	9,224

Table 2: Corpus summary

- *Toxic tweet*: tweets collected through Twitter Streaming API by matching toxic keywords from HateBase and labeled as toxic or non-toxic by human raters (Bahar et al., 2020).

All datasets are sampled to have an equal class balance. The basic dataset information is summarized in Table 2.

### 4.2 Creating Matched Sentences

We first get pairwise matched sentences for words of interest. In this work, we focus on words that have relatively strong correlations with each class. So we fit a logistic regression classifier for each dataset and select the top features by placing a threshold on coefficient magnitude (i.e., words with high positive or negative coefficients). For IMDB movie reviews, Kindle reviews, and Toxic comments, we use a coefficient threshold 1.0; and for Toxic tweet, we use threshold 0.7 (to generate comparable number of candidate words).

We find matched sentences for each word following the method in §3.2. Table 1 shows five examples of pairwise matches. The total number of matched sentences are shown in Table 2.

### 4.3 Word Classification

The goal of word classification is to distinguish between spurious words and genuine words. We first manually label a small set of words as spurious or genuine (Table 3). For sentiment classification, we consider both positive and negative words. For toxicity classification, we only consider toxic words. We had two annotators annotate each term; agreement was generally high for this task (e.g., 96% raw agreement), with the main discrepancies arising from knowledge of slang and abbreviations.

We represent each word with the numerical features calculated from matched sentences (§3.3), standardized to have zero mean and unit variance. Finally, we apply a logistic regression model for the word classifier. We explore the word classifier performance for same domain and domain adaptation.

**Same domain:** We apply 10-fold cross validation to estimate the word classifier’s accuracy within the same domain. The idea is that in practice one would label a set of words, fit a classifier, then apply to the remaining words.

**Domain adaptation:** To reduce the word annotation burden, we are interested in understanding whether a word classifier trained in one domain can be applied in another. Thus, we measure cross-domain accuracy, e.g., by fitting the word classifier on IMDB data and evaluating on Kindle data.

#### 4.4 Feature Selection Based on Spurious Correlation

We compare several strategies to do feature selection for the initial document classification tasks.

Using the word classifier, each word is assigned a probability of being spurious, which we use to sort terms for feature selection. That is, words deemed most likely to be spurious are removed first. As comparison, we experiment with four different strategies to rank words in the order to be removed.

**Oracle** This is the gold standard. We treat the manually labeled spurious words as equally important and sort them in random order. This gold standard ensures that the removed features are definitely spurious.

**Random** This is a baseline method that sorts the top words in random order, where top words could be spurious or genuine, and the words are removed in random order.

**Same domain prediction** We sort words in descending order of the probability of being spurious, according to the word classifier trained on the same domain (using cross-validation).

**Domain adaptation prediction** This is a similar sorting process with the previous strategy except that the probability is from domain adaptation, where the word classifier is trained on a different dataset. We consider domain transfer between (IMDB  $\leftrightarrow$  Kindle) and (Toxic comment  $\leftrightarrow$  Toxic tweet).

In the document classification task, we sample majority and minority groups by selecting equal number of sentences for each top word to ensure a fair comparison during feature selection. We check feature selection performance for each group by gradually removing spurious words following the order of each strategy described above. As a final comparison, we also implement the method sug-

gested in Sagawa et al. (2020b), which reduces the effect of spurious correlation from training data. To do so, we sample the majority and minority group from training data, and down-sample the majority group to have equal size with minority group. We then fit the document classifier on the new training data and evaluate its performance on the test set. Note that this method assumes knowledge of which terms are spurious. Our approach can be seen as a way to first estimate which words are spurious and then adjust the classifier accordingly.

## 5 Results and Discussion

In this section, we show results for identifying spurious correlations and then analyze the effect of removing spurious correlations in different cases.

### 5.1 Word Classification

Table 3 shows the ROC AUC scores for classifier performance. To place these numbers in context, recall that the words being classified were specifically selected because of their strong correlation with the class labels. For example, some spurious positive words appear in 20 positive documents and only a few negative documents. Despite the challenging nature of this task, Table 3 shows that word classifier performs well at classifying spurious and genuine words with AUC scores range from 0.657 to 0.823. Furthermore, the domain adaptation results indicate limited degradation in accuracy, and occasionally improvements in accuracy. The exception is toxic comment, where accuracy is 10% worse for domain adaptation. We suspect that this is caused by the low quality sentences in the toxic tweet dataset (this is the only dataset that is tweets instead of sentences).

Fig 2 shows an example of the domain adaptation results. This word classifier is trained on words from Kindle reviews and predicted on words from IMDB movie reviews. We observe that *culture*, *spielberg*, *russian*, *cinema* are correctly predicted to have high probabilities of being spurious, while *refreshing*, *heartbreaking*, *wonderful*, *fun* are correctly predicted to have relatively lower probabilities of being spurious. We also observe that the predictions for *unique* and *ages* do not agree with human labels. Additionally, we show top-5 spurious and genuine words predicted for each dataset in Table 4. Error analysis suggests that misclassifications are often due to small sample sizes – there are some genuine words that simply do not appear

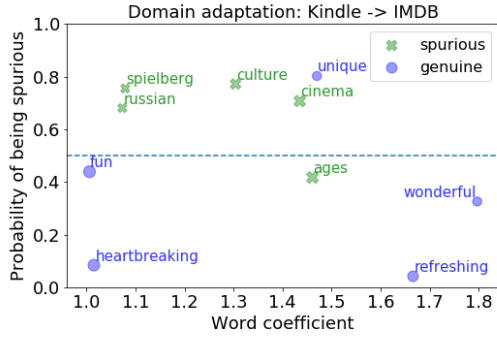


Figure 2: Example of spurious and genuine words predicted from word classifier trained on words from Kindle reviews and applied to words from IMDB reviews.

	IMDB reviews	Kindle reviews	Toxic comment	Toxic tweet
#spurious	90	119	40	72
#genuine	174	100	73	45
same domain	0.776	0.657	0.823	0.686
domain adaptation	0.741	0.699	0.726	0.744

Table 3: Word classifier performance (AUC score)

enough to find good matches. In future work we will investigate how data size influences accuracy.

Examining the top coefficients in the word classifier, we find that features over the match similarity scores tend to be highly correlated with genuine words, often more so than the ATE scores themselves. For example, in IMDB the average match similarity score has a coefficient of 1.3, while the ATE has a coefficient of 0.8. These results suggest that combining traditional ATE estimates with other features derived from the matching procedure can provide stronger signals for distinguishing spurious and genuine correlations.

## 5.2 Feature Selection by Removing Spurious Correlations

We apply different feature selection strategies introduced in §4.4 and test their performance on majority and minority groups. Additionally, we show results for “All”, which is the union of majority and minority groups.

Fig 3 shows feature selection performance on IMDB movie reviews dataset and Kindle reviews dataset. The starting point in each plot shows the performance of not removing any feature. The horizontal line shows the performance of the method suggested in Sagawa et al. (2020b) (which does not do feature selection).

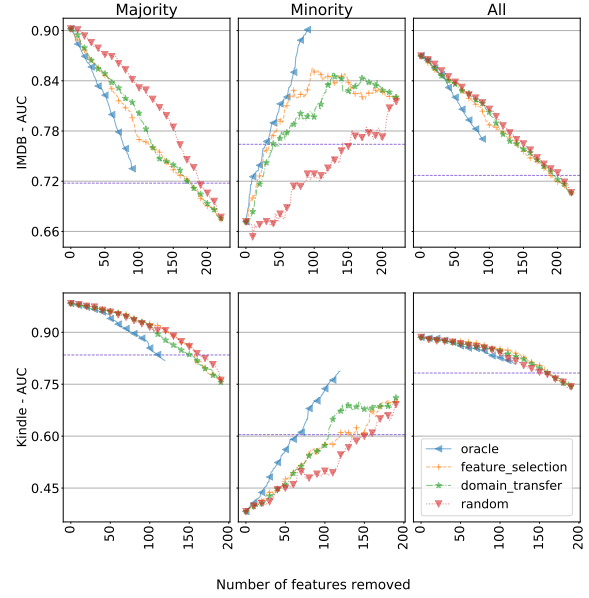


Figure 3: Feature selection for sentiment classification.

For the majority group, because spurious correlations learned during model training agree with class labels, the model will perform well on this group, and removing spurious features will hurt performance (i.e., about 20% drop of AUC score in both datasets). On the contrary, the spurious correlations do not hold in minority group. Thus, the model does not perform well at the starting point when we are not removing any spurious feature, and the performance increases when we gradually remove spurious features. After removing enough spurious features, the model performance stabilizes.

For IMDB reviews, removing spurious features improves performance by up to 20% AUC for minority group, and feature selection using the word classifier outperforms random ordering substantially. For Kindle, removing spurious features improves accuracy by up to 30% AUC for minority group. Interestingly, domain adaptation actually appears to outperform the within-domain results, which is in line with word classifier performance shown in Table 3 (i.e., domain adaptation outperforms within domain AUC by 4.2% for Kindle). The “All” set is the test set that combines the majority and minority groups. It shows the trade-off between performance on majority group and minority group. If removing spurious features hurts more on the majority group than it helps minority group, then the performance on test set would decrease, and vice versa. In our experiment, because the majority group has more samples than the minority group, the final performance on test set generally

IMDB		Kindle		Toxic comment		Toxic tweet	
spurious	genuine	spurious	genuine	spurious	genuine	spurious	genuine
unintentional	refreshing	boy	omg	intelligence	idiot	edkrassen	cunt
russian	horrible	issues	definitely	parasites	stupid	hi	twat
benigni	uninspired	benefits	draw	sucking	idiots	pathetic	retard
animated	strength	teaches	returned	mongering	stupidity	side	pussy
pulls	exhilarating	girl	halfway	lifetime	moron	example	ass
visceral	refreshing	finds	omg	mongering	stupid	aint	cunt
mike	rare	mother	highly	lunatics	idiot	between	twat
unintentional	horrible	girl	returned	slaughter	idiots	wet	retard
strange	ingenious	us	down	narrative	idiotic	side	faggot
intelligent	sly	humans	enjoyed	brothers	stupidity	rather	pussy

Table 4: Top 5 spurious and genuine words predicted by the in-domain word classifier (first five rows) and cross-domain classifier (last five rows).

decreases when removing spurious features. (Additional results for toxic datasets can be found in the Appendix.)

## 6 Related Work

Wood-Doughty et al. (2018) and Keith et al. (2020) provide good overviews of the growing line of research combining causal inference and text classification. Two of the most closely related works, mentioned previously, are Sagawa et al. (2020b) and Paul (2017).

Sagawa et al. (2020b) investigate how spurious correlations arise in classifiers due to overparameterization. They compare overparameterized models with underparameterized models and show that overparameterization hurts worst-group error, where the spurious correlation does not hold. They do simulation experiments with core features encoding actual label and spurious features encoding spurious attributes. Results show that the relative size of the majority group and minority group as well as the informativeness of spurious features modulates the effect of overparameterization. While Sagawa et al. (2020b) assume it is known ahead of time which features are spurious, here we instead try to predict that in a supervised learning setting.

Paul (2017) proposes to do feature selection for text classification by causal inference. He adapts the idea of propensity score matching to document classification and identifies causal features from matched samples. Results show meaningful word features and interpretable causal associations. Our primary contributions beyond this prior work is (a) to use features of the matching process to better identify spurious terms using supervised learning, and (b) to analyze effects in terms of majority and minority groups. Indeed, we find that using the

treatment effect estimates alone for the word classifier results in worse accuracy than combining it with the additional features.

Recently, Kaushik et al. (2020) show the prevalence of spurious correlations in machine learning by having humans make minimal edits to change the class label of a document. Doing so reveals large drops in accuracy due to the model’s overdependence on spurious correlations.

Another line of work investigates how confounds can lead to spurious correlations in text classification (Elazar and Goldberg, 2018; Landeiro and Culotta, 2018; Pryzant et al., 2018; Garg et al., 2019). These methods typically require the confounding variables to be identified beforehand (though Kumar et al. (2019) is an exception).

A final line of work views spurious correlations as a result of an adversarial, data poisoning attack (Chen et al., 2017; Dai et al., 2019). The idea is that an attacker injects spurious correlations into the training data, so as to control the model’s predictions on new data. While most of this research focuses on the nature of the attack models, future work may be able to combine the approaches in this paper to defend against such attacks.

## 7 Conclusion

We have proposed a supervised classification method to distinguish spurious and genuine correlations in text classification. Using features derived from matched samples, we find that this classifier achieves moderate to high accuracy even tested on strongly correlated terms. Additionally, due to the generic nature of the features, we find that this classifier does not suffer much degradation in accuracy when trained on one dataset and applied to another dataset. Finally, we use this classifier to inform feature selection for the original classi-



fication tasks. Results show that removing words in the order of their predicted probability of being spurious results in more robust classification with respect to worst-case accuracy.

## References

- John Aldrich et al. 1995. Correlations genuine and spurious in pearson and yule. *Statistical science*, 10(4):364–376.
- Radfar Bahar, Shivaram Karthik, and Aron Culotta. 2020. Characterizing variation in toxic language by social context. *ICWSM-2020*.
- Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. 2017. Targeted backdoor attacks on deep learning systems using data poisoning. *arXiv preprint arXiv:1712.05526*.
- Jiazhu Dai, Chuanshuai Chen, and Yike Guo. 2019. A backdoor attack against lstm-based text classification systems. *CoRR*, abs/1905.12457.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805.
- Yanai Elazar and Yoav Goldberg. 2018. Adversarial removal of demographic attributes from text data. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 11–21.
- Sahaj Garg, Vincent Perot, Nicole Limtiaco, Ankur Taly, Ed H Chi, and Alex Beutel. 2019. Counterfactual fairness in text classification through robustness. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, pages 219–226.
- Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Franco Turini, Fosca Giannotti, and Dino Pedreschi. 2018. A survey of methods for explaining black box models. *ACM computing surveys (CSUR)*, 51(5):1–42.
- Ruining He and Julian McAuley. 2016. Ups and downs. *Proceedings of the 25th International Conference on World Wide Web - WWW '16*.
- Guido W Imbens. 2004. Nonparametric estimation of average treatment effects under exogeneity: A review. *Review of Economics and statistics*, 86(1):4–29.
- Divyansh Kaushik, Eduard Hovy, and Zachary C Lipton. 2020. Learning the difference that makes a difference with counterfactually-augmented data. In *ICLR*.
- Katherine A Keith, David Jensen, and Brendan O'Connor. 2020. Text and causal inference: A review of using text to remove confounding from causal estimates. *arXiv preprint arXiv:2005.00649*.
- Gary King and Richard Nielsen. 2019. Why propensity scores should not be used for matching. *Political Analysis*, 27(4):435–454.
- Jon Kleinberg, Jens Ludwig, Sendhil Mullainathan, and Ashesh Rambachan. 2018. Algorithmic fairness. In *Aea papers and proceedings*, volume 108, pages 22–27.
- Sachin Kumar, Shuly Wintner, Noah A Smith, and Yulia Tsvetkov. 2019. Topics to avoid: Demoting latent confounds in text classification. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4144–4154.
- Virgile Landeiro and Aron Culotta. 2018. Robust text classification under confounding shift. *Journal of Artificial Intelligence Research*, 63:391–419.
- David Lopez-Paz, Krikamol Muandet, Bernhard Schölkopf, and Iliya Tolstikhin. 2015. Towards a learning theory of cause-effect inference. In *International Conference on Machine Learning*, pages 1452–1461.
- David Martens and Foster Provost. 2014. Explaining data-driven document classifications. *Mis Quarterly*, 38(1):73–100.
- Bo Pang and Lillian Lee. 2005. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In *Proceedings of the 43rd annual meeting on association for computational linguistics*, pages 115–124. Association for Computational Linguistics.
- Michael J Paul. 2017. Feature selection as causal inference: Experiments with text classification. In *Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017)*, pages 163–172.
- Reid Pryzant, Kelly Shen, Dan Jurafsky, and Stefan Wagner. 2018. Deconfounded lexicon induction for interpretable social science. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1615–1625.
- Joaquin Quionero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D Lawrence. 2009. *Dataset shift in machine learning*. The MIT Press.
- Hema Raghavan, Omid Madani, and Rosie Jones. 2005. Interactive feature selection. In *IJCAI*, volume 5, pages 841–846.
- Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. 2020a. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. *ICLR*.

- Shiori Sagawa, Aditi Raghunathan, Pang Wei Koh, and Percy Liang. 2020b. [An investigation of why over-parameterization exacerbates spurious correlations](#).
- Manali Sharma, Di Zhuang, and Mustafa Bilgic. 2015. Active learning with rationales for text classification. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 441–451.
- Elizabeth A Stuart. 2010. Matching methods for causal inference: A review and a look forward. *Statistical science: a review journal of the Institute of Mathematical Statistics*, 25(1):1.
- Christopher Winship and Stephen L Morgan. 1999. The estimation of causal effects from observational data. *Annual review of sociology*, 25(1):659–706.
- Zach Wood-Doughty, Ilya Shpitser, and Mark Dredze. 2018. Challenges of using text classifiers for causal inference. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, volume 2018, page 4586. NIH Public Access.
- Omar Zaidan, Jason Eisner, and Christine Piatko. 2007. Using “annotator rationales” to improve machine learning for text categorization. In *Human language technologies 2007: The conference of the North American chapter of the association for computational linguistics; proceedings of the main conference*, pages 260–267.

950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999

## A Appendices

We also perform feature selection on toxic comment and toxic tweet datasets, where we only focus on toxic features. As we observe in Fig 4, for toxic comment data, removing spurious features improves performance by up to 20% AUC for minority group, and 30% AUC for the toxic tweet dataset. Compared with sentiment datasets, toxic datasets have fewer spurious words because the words to be removed are only from toxic class, while in sentiment datasets, the words to be removed are from both positive and negative classes. Besides that, toxic tweet dataset is noisy with many short and incomplete sentences. So the feature selection methods perform slightly different on toxic datasets compared with sentiment datasets.

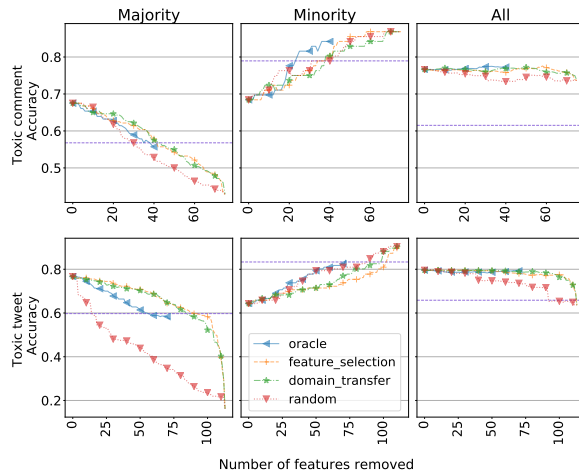


Figure 4: Feature selection for toxic datasets. As we only focus on toxic words, these groups only have sentences from one class, so we show accuracy score on y-axis.