

# Unsupervised Discovery of Implicit Gender Bias

Anjalie Field

Carnegie Mellon University  
anjalief@cs.cmu.edu

Yulia Tsvetkov

Carnegie Mellon University  
ytsvetko@cs.cmu.edu

## Abstract

Despite their prevalence in society, social biases are difficult to identify, primarily because human judgements in this domain can be unreliable. We take an *unsupervised* approach to identifying gender bias against women at a comment level and present a model that can surface text likely to contain bias. Our main challenge is forcing the model to focus on signs of implicit bias, rather than other artifacts in the data. Thus, our methodology involves reducing the influence of confounds through propensity matching and adversarial learning. Our analysis shows how biased comments directed towards female politicians contain mixed criticisms, while comments directed towards other female public figures focus on appearance and sexualization. Ultimately, our work offers a way to capture subtle biases in various domains without relying on subjective human judgements.<sup>1</sup>

## 1 Introduction

Despite widespread documentation of the negative impacts of bias, stereotypes, and prejudice (Krieger, 1990; Goldin, 1990; Steele and Aronson, 1995; Logel et al., 2009; Schluter, 2018), these concepts remain difficult to define and identify, especially for non-experts. Social biases appear to be a natural component of human cognition that allow people to make judgments efficiently (Kahneman et al., 1982). As a result, they are often *implicit*—people are unaware of their own biases (Blair, 2002; Bargh, 1999)—and manifest subtly, e.g., as microaggressions or condescension (Huckin, 2002; Sue, 2010).

Much NLP literature has examined biases in data, algorithms, or model performance, and the negative pipeline between them: models absorb and amplify data biases, which impacts performance (Sun et al., 2019). However, little work has looked

further up the pipeline and relied on the assumption that biases in data originate in human cognition.

In contrast, this assumption motivates our work: an *unsupervised approach to detecting implicit gender bias in text*. Text provides an ideal avenue for studying bias, because human cognition is closely tied to natural language. Psychology studies often examine human perceptions through word associations (Greenwald et al., 1998). However, the implicit nature of bias suggests that human annotations for bias detection may not be reliable, which motivates an unsupervised approach.

The goals of our work align with prior work in NLP that has examined biases in real-world data. However, prior work examines bias at a broad corpus level or relies on supervised models. While corpus-level analyses, e.g. associations between gendered words and stereotypes, can be insightful (Bolukbasi et al., 2016; Fast et al., 2016; Caliskan et al., 2017; Garg et al., 2018; Friedman et al., 2019; Chaloner and Maldonado, 2019), they are difficult to interpret over short text spans. They also often rely on human-defined “known” stereotypes, such as lists of traditionally male and female occupations obtained through crowd-sourcing, which restricts analysis to a narrow surface-level domain. Similarly, supervised approaches can provide insight into carefully defined types of bias (Wang and Potts, 2019; Breitfeller et al., 2019; Sap et al., 2020), but they rely on human annotations tasks, which are difficult to design or generalize to other domains, especially because social concepts differ across contexts and cultures (Dong et al., 2019).

Our work offers a new approach to surfacing gender bias that does not require direct supervision and is meaningful at a sentence or paragraph level. *We create a model that takes text in the 2<sup>nd</sup>-person perspective as input and predicts the gender of the person the text is addressed to.* If the classifier predicts the gender of the addressee with high confidence based only on the text directed to them, we

<sup>1</sup>Code and pre-trained models are available at [https://github.com/anjalief/unsupervised\\_gender\\_bias](https://github.com/anjalief/unsupervised_gender_bias)

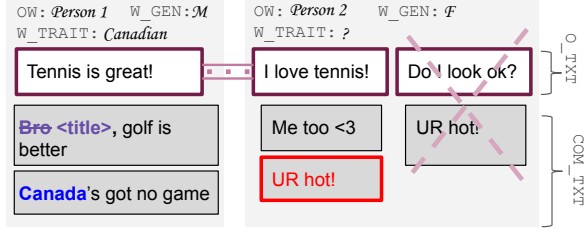


Figure 1: We train a classifier to predict the gender ( $W\_GEN$ ) of the person that text is addressed to ( $COM\_TXT$ ), while demoting features that are predictive of gender but not predictive of bias. Posts with similar content are matched through propensity scores; unmatched posts are discarded. Latent traits of the addressee (e.g., nationality) are demoted through an adversarial objective. Overtly gendered language (“Bro”) is substituted. Comments indicative of gender despite these restrictions are likely to contain bias.

hypothesize that the text is likely to contain bias. The main challenge is encouraging the model to focus on text features that are indicative of bias, rather than artifacts in data that correlate with the gender of the addressee but occur because of confounding variables (*confounds*). Thus, the core of our methodology focuses on reducing the influence of confounds. Our goal is not to improve accuracy of the gender-prediction task, but rather to validate that our methodology demotes confounds and surfaces comments likely to contain gender bias.

In §2, we define the problem and intuition behind our approach. We describe our methods for confound demotion in §3, and we evaluate them in §5. Our evaluation involves examining how confound control affects performance on in-domain and out-of-domain classification tasks, including detection of gender-based microaggressions. Our results suggest that our model successfully identifies text likely to contain bias against women, allowing us to analyze how this bias differs across domains (§6). To the best of our knowledge, this is the first work that aims to analyze bias in short text spans by learning implicit associations from data sets.

## 2 Problem Formulation

Our primary task is to detect gender bias in a communicative domain, specifically in texts targeting an addressee (i.e., 2<sup>nd</sup>-person) without relying on explicit bias annotations. Our goals align with a causality framework in that we seek to identify content that occurs because of the gender of the addressee rather than other factors. We can define a

counterfactual: *Would the addressee have received different text if their gender were different?*

While our framework is broadly applicable, in order to define consistent notation, we consider a setup where our primary text is a comment written in reply to text written by someone else. This includes domains like replies on social media posts, or comments on newspaper articles, and can be generalized other media, e.g., comments on YouTube videos. We identify the following variables:

- **OW**: “Original Writer”, the person who wrote the original text, e.g., the addressee
- **O\_TXT**: content of the original text
- **W\_GEN**: the gender (M, F) of OW. We use a binary variable because all of the individuals in our corpus identify as men or women, but our methodology is generalizable and can be used to examine bias against other genders.
- **W\_TRAITS**: any traits of OW other than gender, e.g., social role, age, nationality.
- **COM\_TXT**: comments replying to O\_TXT

Our goal is to detect bias in  $COM\_TXT$  values that occurs because of  $W\_GEN$ . A naive approach would train a classifier to predict  $W\_GEN$  from  $COM\_TXT$  and assume that any  $COM\_TXT$  values for which the classifier correctly predicts  $W\_GEN$  with high confidence contain bias. However,  $COM\_TXT$  may contain features that are predictive of  $W\_GEN$  but are not indicative of bias.

For example, in Figure 1, when the comment “UR hot!” ( $COM\_TXT$ ) is addressed to someone who said “I love tennis!” ( $O\_TXT$ ), it is an objectification and unsolicited reference to appearance, which could indicate bias. However, when it is addressed to someone who said “Do I look ok?”, it is likely not indicative of bias. If women ask “Do I look ok?” more frequently than men, this naive classifier would identify “UR hot!” as likely addressed towards a woman and identify it as biased. However, we only want the model to learn that references to appearance are indicative of gender if they occur in unsolicited contexts. Thus our model needs to account for the effects of  $O\_TXT$ : Because of correlations between  $W\_GEN$  and  $O\_TXT$ ,  $COM\_TXT$  values may contain features that are *predictive* of  $W\_GEN$ , but are *caused* by  $O\_TXT$ , rather than by  $W\_GEN$ . We face a similar problem with  $W\_TRAITS$ . From the synthetic example in Figure 1, if our data set contains more men from Canada than women, the model might learn that references to Canada indicate  $W\_GEN = M$ .

We provide additional empirical examples in §4.

We refer to factors that might influence COM.TXT as *confounding variables* and the artifacts that they produce in COM.TXT as *confounds*. We distinguish two types: *observed* and *latent*. **Latent confounding variables cannot be controlled if they are entirely unknown; instead, we assume there are observed signals that can be used to infer them, but the values themselves are difficult to explicitly enumerate.** In addition to confounds introduced by O.TXT and W.TRAITS, COM.TXT may also contain overt signals, e.g. titles like **“Ma’am” or “Sir”**, that are predictive of gender, but not indicative of bias. We thus identify 3 factors to account for: O.TXT, W.TRAITS, and overt signals.

### 3 Methodology

Our overall methodology centers on creating a classifier that predicts gender of the addressee while controlling for the effects of observed confounding variables (O.TXT), latent confounding variables (W.TRAITS), and overt signals. The input to the prediction model is COM.TXT, while the output is W\_GEN, and we aim to identify bias in COM.TXT.

#### 3.1 Controlling Observed Confounding Variables through Propensity Matching

Our primary method for controlling for O.TXT is *propensity matching*. Propensity matching was developed to replicate the conditions of randomized trials in causal inference studies (Rosenbaum and Rubin, 1983, 1985). In this step, we discard any COM.TXT training samples whose associated O.TXT is heavily affiliated with only one gender. In Figure 1, if we assume that only women post “Do I look ok?”, we would discard all comments posted in reply to the O.TXT “Do I look ok?”. We ultimately seek to balance our data set, so that the set of all COM.TXT where W\_GEN = M has similar associated O.TXT as the set of all COM.TXT where W\_GEN = F. Thus, we match each O.TXT where W\_GEN = F with a similar O.TXT where W\_GEN = M and discard all unmatched data.

Ideally, we would match O.TXT values written by men with identical O.TXT values written by women, but this is infeasible in practice. Instead, the key insight behind propensity matching is that it is sufficient to match data points based on the probability of the target variable, e.g., the probability that W\_GEN = F (Rosenbaum and Rubin, 1983, 1985). Thus, the propensity score  $e_i$  for a COM.TXT<sub>*i*</sub> is

defined as the probability that W\_GEN = F, given the confounding variable, O.TXT<sub>*i*</sub>:

$$e_i(\text{COM.TXT}_i) = P(W\_GEN_i = F | O\_TXT_i)$$

To balance our data set, we need to ensure that the set of COM.TXT where W\_GEN = M has a similar propensity score distribution as the set of COM.TXT where W\_GEN = F. Because propensity scores are dependent on O.TXT, all COM.TXT replied to the same O.TXT have the same propensity score. We can then equate  $e_i(\text{COM.TXT}_i) = e_i(\text{O.TXT}_i)$ , and focus estimating O.TXT scores.

Propensity scores can be estimated by using a classification model that is trained to predict the target attribute W\_GEN<sub>*i*</sub> = F from the observed confounding variable O.TXT<sub>*i*</sub> (Westreich D, 2010; Lee et al., 2010). We use a bidirectional LSTM encoder followed by two feedforward layers with a tanh activation function and a softmax in the final layer. Then, we use greedy matching to match each O.TXT<sub>*i*</sub> where the true value of W\_GEN<sub>*i*</sub> is F with O.TXT<sub>*j*</sub> where the true value of W\_GEN<sub>*j*</sub> is M and  $|e_i(\text{O.TXT}_i) - e_j(\text{O.TXT}_j)|$  is minimal (Gu and Rosenbaum, 1993).

We institute a threshold  $c$  (Stuart, 2010), where we discard O.TXT<sub>*i*</sub> if we cannot find a O.TXT<sub>*j*</sub> such that  $|e_i(\text{O.TXT}_i) - e_j(\text{O.TXT}_j)| \leq c$ . Thus, for example, we would match a post written by a woman that is “stereotypically female” (e.g.,  $e_i$  is large) with a post written by a man that is also “stereotypically female” (e.g.,  $e_j$  is also large). In Figure 1, we match “Tennis is great” with “I love tennis”, and we discard “Do I look ok?” as unable to be matched. However, using propensity matching rather than direct matching allows us to match O.TXT values that are about different topics, as long as they are equally likely to have been written by a woman.

Finally, our actual model input consists of COM.TXT, not of O.TXT. Once we have matched pairs of O.TXT values, we need to ensure that we have an equal number of COM.TXT values for each O.TXT in the pair in order to have a balanced data set. Then, for each matched [O.TXT<sub>*i*</sub>, O.TXT<sub>*j*</sub>], we randomly downsample to have an equal number of COM.TXT values for each O.TXT in the pair. In this way, we balance the training set of COM.TXT in terms of how predictive the confounding variable O.TXT is of the target attribute W\_GEN.

### 3.2 Controlling Latent Confounding Variables through Adversarial Training

While propensity matching is a desirable way to control for confounding variables because of established literature, matching is only possible for observed variables (Gu and Rosenbaum, 1993; Rosenbaum, 1988). In our data, while O\_TXT is observed, W\_TRAITS is not possible to match on (further discussion in §4). Instead, we use an adversarial objective drawn from Kumar et al. (2019) to encourage the model to ignore W\_TRAITS.

**Confound representation** While we cannot explicitly enumerate W\_TRAITS, we know that they are associated with the identity of OW, and we can infer them from COM\_TXT addressed to OW. We use associations between OW and COM\_TXT to derive a feature vector for each COM\_TXT<sub>i</sub> that reflects W\_TRAITS<sub>i</sub>. The latent confounds to demote are represented as multinomial distributions, derived from log-odds scores (Monroe et al., 2008).

For each label OW =  $k$  and each word type  $w$  in all COM\_TXT, we calculate the log-odds score  $lo(w, k) \in \mathbf{R}$ , where higher scores indicate stronger associations between  $k$  and the word. In Figure 1,  $lo(\text{Canada}, \text{Person } 1)$  would be high, as COM\_TXT values addressed to Person 1 often contain the word Canada. Then, following Kumar et al. (2019), we define a distribution: for all  $k \in \text{OW}$  and an input COM\_TXT<sub>i</sub>, =  $\langle w_1, \dots, w_n \rangle$ :

$$p(k|\text{COM\_TXT}_i) \propto p(k) \prod_{i=1}^n p(w_i|k)$$

$p(k)$  is estimated from the distribution of  $k$  in the training data, i.e., the proportion of COM\_TXT values addressed to OW =  $k$ .  $p(w_i|k)$  is proportional to  $\sigma(lo(w, k))$ , where we use the sigmoid function ( $\sigma$ ) to map log-odds scores to the range [0,1] and then normalize them over the vocabulary to obtain valid probabilities. For each input COM\_TXT<sub>i</sub>, we then obtain a vector whose elements are  $p(k|\text{COM\_TXT}_i)$  and whose dimensionality is the number of OW individuals in the training set. We normalize these vectors to obtain multinomial probability distributions which reflect COM\_TXT<sub>i</sub>’s association with each OW individual. Thus, when we demote this vector during training, we force the classifier to learn features that are indicative of the group W\_GEN and not features that are indicative of individual members of this group

(e.g., some group members are from Canada). We refer to the confound vector as  $t_i$ . Justification for the log-odds representation as opposed to alternatives is presented in Kumar et al. (2019).

**Training Procedure** Our goal is to obtain a model that can predict W\_GEN, but cannot predict the latent confounds represented by  $t_i$ . To achieve this, the model is trained in an alternate GAN-like procedure (Goodfellow et al., 2014).

First, the input  $x \in \text{COM\_TXT}$  is encoded using an encoder neural network  $h(x; \theta_h)$  to obtain a hidden representation  $\mathbf{h}_x$ . This representation is then passed through two feedforward networks: (1)  $c(h(x); \theta_c)$  to predict the label  $y \in \{\mathbf{M}, \mathbf{F}\}$ ; and (2) an adversary network  $\text{adv}(h(x); \theta_a)$  to predict the vector representation of the latent confounds.

We train the encoder, so that  $\mathbf{h}_x$  does not contain any information predictive of the confound vector, but does contain information predictive of the target attribute. Thus our primary training objective is:

$$\min_{c,h} \frac{1}{N} \sum_{i=1}^N \text{CE}(c(h_{x_i}), y_i) + \text{KL}(\text{adv}(h_{x_i}), \mathbb{U}_K)$$

where  $\mathbb{U}$  represents a uniform distribution, CE represents cross-entropy loss, and KL represents KL-divergence. We refer to Kumar et al. (2019) for the training procedure that alternates minimizing this objective and training the adversary.

### 3.3 Overt Signals

Finally, we control for overt signals using word substitutions that replace gendered terms with more neutral language, for example woman  $\rightarrow$   $\langle \text{person} \rangle$  and man  $\rightarrow$   $\langle \text{person} \rangle$ . We create a 66-term list of substitutions from existing resources (Zhao et al., 2018; Bolukbasi et al., 2016) as well as our observations of the data. We also use substitutions to remove the names of addressees from comment, replacing OW’s “Firstname” and “Lastname” with “ $\langle \text{name} \rangle$ ” in COM\_TXT. We do not attempt to identify nicknames, as the confound demotion method described in §3.2 should already mitigate the influence of individual names, and we perform the substitution as merely an extra precaution.

## 4 Experimental Setup

Our primary data is the Facebook subsection of the RiGender corpus (Voigt et al., 2018). The data contains two subsections: *Politicians* (400K posts and 13.9M replies addressed to 412 then-current U.S.



members of Congress), and *Public Figures* (118K posts and 10.7M replies addressed to 105 famous people such as actresses and tennis players).

We can show that O\_TXT is a confounding variable by computing log-odds scores between the words in O\_TXT and W\_GEN (Monroe et al., 2008). In the Politicians data, the most female-associated words are *women*, *Congresswoman*, *sexual*, and *assault*. The most male-associated words are *Obamacare*, *Iran*, *EPA*, and *spending*. It is evident that male and female politicians post about different topics, e.g., female politicians likely post more about sexual assault. A naive model may predict that comments using sexual language are addressed towards women, but increased sexual language may occur because of O\_TXT, rather than gender bias.

A similar problem occurs with W\_TRAITS, e.g., the corpus has more comments addressed to female tennis players (9 players; 184K comments) than male players (1 player; 29K comments). The model can obtain high accuracy by predicting  $W\_GEN = F$  for COM\_TXT with the word “tennis”. Unlike O\_TXT, which is observable from the data, we have no way of enumerating every possible value in W\_TRAITS. Even if we could enumerate them, we do not expect propensity matching over W\_TRAITS to work, because we cannot find reasonable matches, e.g., there is only one senior senator from Massachusetts. Additionally, W\_TRAITS can be as fine-grained as names: we cannot find a male senator whom commenters call “Liz Warren”.

We divide each data set into train, dev, and test sets, enforcing no OW overlap between subsets. We perform propensity matching and derive the confound vectors to demote using only the training data. We apply word substitutions to all subsets.<sup>2</sup>

## 5 Evaluation

We train our model to predict W\_GEN from COM\_TXT, employing propensity matching over O\_TXT, word substitutions over COM\_TXT, and W\_TRAITS demotion. We focus on evaluating how well our model controls for confounds and whether or not it captures gendered language. Successful demotion of confounds would suggest that our model learns to identify text indicative of gender bias.

**Observed Confounding Variable Demotion** In Figure 2, we show log-odds scores, measuring association between O\_TXT and W\_GEN in the train-

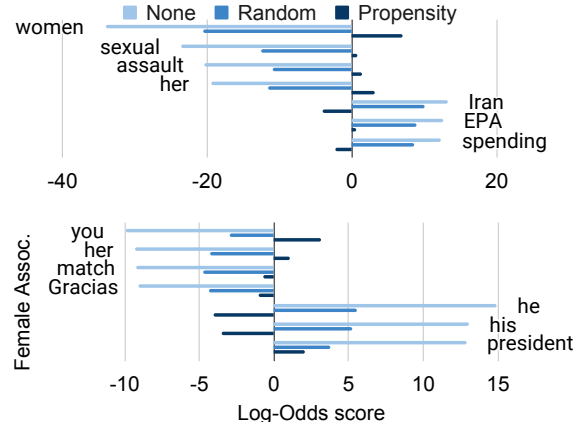


Figure 2: Log-odds scores for most polar words in Politicians (top) and Public Figures (bottom) data, with different matching methods. Propensity matching best reduces polarity.

ing set before and after propensity matching. For comparison, we also show scores for a randomly matched data set, in which we balance O\_TXT to have an equal proportion of F and M labels by random sampling (constructed to be the same size as the propensity matched set). In the Politicians and Public Figures data, propensity matching reduces the magnitude of the most polar words: log-odds scores for the matched data are closer to zero than for the non-matched or randomly matched data.<sup>3</sup> Further, propensity matching can even cause the polarity to change direction: words that were originally female-associated (e.g. “her”) become slightly male-associated. These figures suggest that propensity matching effectively reduces the confounding influence of O\_TXT.

**Latent Confounding Variable Demotion** We evaluate how well our model demotes the influence of latent confounding variables over the held-out test sets (Table 1). We created data splits so that there is no overlap in OW values between the train and test sets. While there may still be overlap in some latent W\_TRAITS, we expect there to be less overlap in W\_TRAITS between the train and test set than within the train set. Thus, improved performance over the held-out test set would suggest that demotion effectively reduces the influence of the latent confounding variables—the model learns characteristics of comments addressed to women generally rather than characteristics specific to the

<sup>2</sup>We provide additional details in Appendix A.

<sup>3</sup>Polarities were reduced without producing new ones: in the Politicians data, the magnitude of the 2 most polar words decreased from -34.0 and 17.9 to -7.68 and 8.52, and in the Public Figures data, from -45.5 and 39.3 to -5.29 and 9.43.

	Public Figs		Politicians	
	F1	Acc.	F1	Acc.
base	74.9	63.8	23.2	73.2
+demotion	<b>76.1</b>	<b>65.1</b>	17.4	<b>77.1</b>
+match	65.4	56.0	28.5	46.7
+match+dem.	68.2	59.7	<b>28.8</b>	51.4

Table 1: Evaluation over held-out test sets, where  $W\_GEN = F$  is considered the positive class. Latent confound demotion improves performance.

individual people in the training set. We do not necessarily expect propensity matching to improve performance, as this method reduces the influence of confounding variables that have high overlap between the train and test sets.

Because the data set is imbalanced (the Politicians test set is 82%M and the Public Figures test set is 35.9%M), we report F1 and accuracy scores in Table 1, where  $W\_GEN = F$  is considered the positive class. As expected, models with demotion perform best on all metrics, with the exception of recall in the Politicians data.<sup>4</sup> We note that in general lack of performance improvement on the test set does not necessarily mean the model is not working, and it could indicate that there is not biased language in the data set. However in this case, since we do observe biased comments in this data (e.g. Table 3), and we do observe a performance increase, the performance increase suggests that confound demotion improves the model’s ability to generalize beyond the individuals in the training set and capture characteristics of language addressed to women in general.

**Detection of Sexist Comments** Finally, we evaluate if our model captures gender-biased language by using it to identify gender-based microaggressions, i.e., “you’re too pretty to be a computer scientist!”. This task is notoriously difficult because words like “pretty” often register as positive content (Breitfeller et al., 2019; Jurgens et al., 2019). Our goal is not to maximize accuracy over microaggression classification, but rather to assess whether or not our model has encoded any indicators of gender bias from the RtGender data set, which would be indicated by better than random performance.

<sup>4</sup>Appendix B reports precision and recall. The discrepancies between F1 and Accuracy are explained by the imbalance in the data set, particularly in the Politicians data set, which is imbalanced in favor of M while we report metrics assuming F is the positive class.

We use a corpus of self-reported microaggressions.<sup>5</sup> In the absence of negative examples that contain no microaggressions, we focus on distinguishing gender-tagged microaggressions (704 posts) from other forms of microaggressions, e.g., racism-tagged (900 posts). We train our model on either the Politicians or Public Figures training data sets, and then we test our model on the microaggressions data set. Because most gender-related microaggressions target women, if our model predicts that the reported microaggression was addressed to a woman (e.g.  $W\_GEN = F$ ), we assume that the post is a gender-tagged microaggression. Thus, our models are *not trained at all* for identifying gender-tagged microaggressions.

Table 2 shows results from our models and two random baselines. “Random” guesses gender-tagged or not with equal probability. “Class Random” guesses gender-tagged or not according to true test distributions (56.1% gender-tagged). All models outperform “Class random”, and all models with demotion also outperform “Random”.

Propensity matching improves F1 when training on the Politicians data, but not Public Figures. Several differences could explain this: the Public Figures set is smaller, so propensity matching causes a more substantial size reduction. Also, the Politicians data is more heavily imbalanced, though notably, it is imbalanced in the same direction as the microaggressions data, while the Public Figures data is imbalanced oppositely. Finally, many microaggressions contain references to appearance, which are also common in the Public Figures data. Many comments to people like actresses focus on their looks, especially because they often post photos. However, by controlling for O\_TXT, propensity matching discards many of these comments. Thus, by demoting a confounding variable, we make the prediction task more difficult. Our goal in confound demotion is not to improve accuracy, but to increase confidence in model outputs.

Nevertheless, the general better-than-random performance of all models is striking, as it suggests strong bias in the underlying training data, which is encoded by our models.

## 6 Analysis of Encoded Bias

Finally, we analyze what type of bias our model learns: (1) we identify words that most impact model confidence; (2) we compare posts surfaced

<sup>5</sup>Details in Appendix C

	Public Figs		Politicians	
	F1	Acc.	F1	Acc.
base	61.3	57.3	48.1	64.2
+demotion	<b>62.2</b>	57.9	53.7	61.5
+match	38.9	55.9	46.9	50.7
+match+dem.	50.9	57.0	<b>56.9</b>	49.9
Random	46.0	49.8	-	-
Class Random	42.1	48.3	-	-

Table 2: Evaluation over the microaggressions data set. Despite not being trained for this task, our models achieve better-than-random performance.

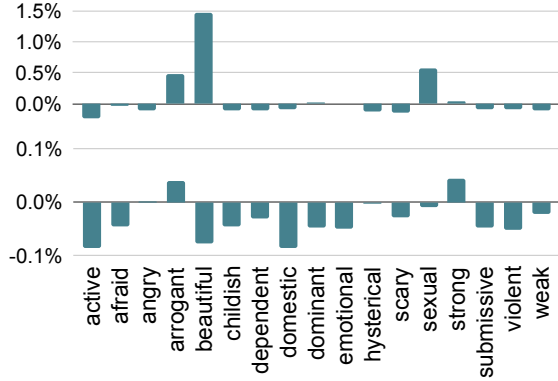


Figure 3: Lexicon differentials between comments with a high likelihood of bias and random samples with  $W\_GEN = \mathbf{F}$  for Public Figures (top) and Politicians (bottom) data. In the Public Figures data, high-likelihood comments are more focused on appearance.

by our model with prior work on stereotypes; (3) we show example posts surfaced by our model. Throughout this section, we use *prediction score* to refer to the output of the final softmax layer of the prediction model, which we take as an estimate of model confidence. We generally focus on COM\_TXT for which our model predicts  $W\_GEN = \mathbf{F}$  with a high prediction score. These are the posts our model identifies as likely to contain bias against women: despite the matching and demotion methods, the model still predicts  $W\_GEN = \mathbf{F}$  with high confidence.

**Influential words** We identify words that strongly influence the model’s decisions by masking out words from comments in the test set and examining the impact on prediction score. For each data set, we take the 500 comments from the test set for which the model predicts  $W\_GEN = \mathbf{F}$  with maximal prediction scores. We then generate masked posts: for every word  $w$  in the post, we

generate a version of the post that omits  $w$ . We run these masked posts through our gender-prediction model and compare the prediction scores where  $w$  is omitted and where  $w$  is not omitted, averaging across all occurrences of  $w$  in the 500 posts. We then examine the set of  $w$  words with the highest differential in prediction score - these are words that, when omitted, cause the model to less associate  $W\_GEN$  with  $\mathbf{F}$ .

In the Public Figures data, the most influential words are appearance-driven and sexualized: *beautiful*, *bellissima*, *amore*, *amo*, *love*, *linda*, *sexo*. In contrast, influential words in the Politicians data are more mixed. Words include references to strength and competence, e.g., *force*, *situation*, as well as traditionally domestic terms, e.g., *spouse*<sup>6</sup>, *family*, *love*. When we repeat this process using the 500 highest-confidence posts from the training set instead of the test set, we find similar results. Influential words in the Public Figures training data primarily refer to appearance, while ones in the Politicians training data include terms like *DINO*.<sup>7</sup> However, influential words from the training data also includes some correlative terms, like names of states, that we would expect the latent confound demotion to de-emphasize. While §5 suggests that our model successfully reduces the influence of confounding variables, more work is needed to eliminate them entirely.

**Comparison to stereotype lexicons** In order to better understand these trends, we draw from prior work on stereotype detection (Fast et al., 2016). We take the set of test comments for which our model predicts  $W\_GEN = \mathbf{F}$  with a high prediction score ( $\geq 0.99$  for Public Figures;  $\geq 0.95$  for Politicians). Then, we compute the difference in frequency of words from a stereotype lexicon (Fast et al., 2016) in this high-confidence prediction set with their frequency in a random sample of the same number of comments where the true value of  $W\_GEN = \mathbf{F}$ .<sup>8</sup>

Figure 3 reports results, which reflect the same trends observed in the influential words. In the Public Figures data, the lexicons that overlap the most with the high-bias posts are “Beautiful”, “Arrogant”, and “Sexual”, which suggests that bias

<sup>6</sup>“<” indicate overt terms substituted out. “<spouse>” replaced “husband”, “husbands”, “wife”, and “wives”.

<sup>7</sup>“Democrat in Name Only” a political insult

<sup>8</sup>We ignore non-English comments and lemmatize the comment text and lexicons. We randomly sample twice and average frequencies between samples. Lexicon counts are normalized by total number of words in the sample.

O_TXT	From reintroducing my legislation to curb sexual assault on college campuses to...
COM_TXT	DINO I hope another real Democrat challenges you next election
O_TXT	Donald Trump is the President, not our ruler...Speak up! Call the White House...
COM_TXT	<name> Shea-Porter, I did not vote for you and have no clue why anyone should have. You do not belong in politics
O_TXT	I am wondering about the guy who actually cried over spilt milk? He must have had...
COM_TXT	Total tangent I know but, you're gorgeous.
O_TXT	Bob and I join Bill Hemmer on America's Newsroom to discuss whether or not...
COM_TXT	I like Bob, but you're hot, so kick <theirs> butt.

Table 3: Example comments surfaced by our model from Politicians (top) and Public Figures (bottom) data sets.

in these comments focuses on appearance and sexualization. In contrast, bias in comments directed towards politicians are less focused, and differences between the high-confidence prediction posts and the random sample are smaller. The two most prominent lexicons are “Arrogant” (primarily driven by lexicon words *special*, *proud*) and “Strong”. Notably, we do not account for negation of lexicon words. A narrative of power is reflected in comments surfaced by our model: “you & Nikki Haley lost my vote on the flag issue *your both weak*”. We provide more examples in Table 3.

Because the stereotype lexicons are small and scores can be dominated by a few words, we also compare LIWC scores (Pennebaker et al., 2001). While most LIWC categories are too broad to align with well-known stereotypes, results are consistent with Figure 3; for Public Figures, the high-bias data scores higher than the random sample for the “Sexual” (0.32 vs. 0.10) and “Body” (0.70 vs. 0.56). For Politicians, the high-bias comments score lower than the random sample in the “Drives” (8.76 vs. 9.71), which encompasses Affiliation, Achievement, Power, Reward, and Risk focus.

The difficulty in evaluating our model against existing lexicons as well as the differences between the two data sets motivates our goal in learning to detect bias automatically. Bias can differ in different contexts, making it difficult to crowdsource through annotations or define through lexicons.

**Examples** Table 3 shows training and test examples surfaced by our model. We identify them by selecting posts where O\_TXT is not strongly gendered (propensity score model described in §3 outputs a prediction score  $< 0.6$ ), but where COM\_TXT is strongly gendered ( $> 0.9$  prediction score). While posts from the Politicians data are diverse, posts from the Public Figures data focus on appearance. These comments reflect the broader trends shown

in the influential words and in Figure 3.

## 7 Related Work

Our work differs from prior work on bias detection in NLP in that we infer bias from data in an unsupervised way, whereas prior work relies on crowd-sourced annotations (Fast et al., 2016; Bolukbasi et al., 2016; Wang and Potts, 2019; Sap et al., 2020). This work typically focuses on specific types of bias, such as condescension (Wang and Potts, 2019) or microaggressions (Breitfeller et al., 2019) and involves carefully constructed annotations schemes that are difficult to generalize to other data sets or types of bias. In contrast, our unsupervised approach is not limited to any particular domain and does not rely on human annotations, which can be subjective.

Less-supervised approaches focus on corpus-level analyses, such as associations between gendered terms and occupational stereotypes (Wagner et al., 2015; Bolukbasi et al., 2016; Fu et al., 2016; Joseph et al., 2017; Nakandala et al., 2017; Friedman et al., 2019; Chaloner and Maldonado, 2019; Hoyle et al., 2019). Methodologies for identifying gender-related differences in text have varied, including word-embedding similarity (Bolukbasi et al., 2016), language model perplexity (Fu et al., 2016), and predictive words identified by logistic regression (Nakandala et al., 2017). These metrics are meaningful over a corpus-level, but are often difficult to interpret over short text spans. Additionally, none of these methods focus on controlling for confounds.

While matching is a well-established method for controlling for confounding variables in causality literature (Rosenbaum and Rubin, 1983, 1985; Stuart, 2010), considerably less work has drawn this methodology into NLP. Most work takes one of two approaches. In the first scenario, text maybe be



a confounding variable that needs to be controlled in order to measure the effect of a non-text variable (Roberts et al., 2020; Veitch et al., 2019). For example, Roberts et al. (2020) examine whether or not papers written by male authors are cited more than ones by female authors, while controlling for the content of the paper. Roberts et al. (2020) also offer a specific method for matching text, which relies on the output of a topic model. In this work, we use the output of an LSTM, which is generally more appropriate for short text, does not make the simplifying BOW assumption, and scales well to large data sets.

In the second scenario, it may be desirable to control for non-text confounds before analyzing text. Chandrasekharan et al. (2017) use matching to identify similar users on Reddit before comparing the content that they post. Our work requires both of these perspectives, as the variable we control for (O\_TXT) and the outcome we analyze (COM\_TXT) are both text. Egami et al. (2018) do consider a similar setting where text is both an outcome and a confound. While their goals differ greatly from ours, our framework is generally consistent with their recommendations. Keith et al. (2020) provide a more complete overview of using text to reduce the influence of confounding variables.

## 8 Limitations and Future Work

While our work serves as an initial approach toward unsupervised detection of comment-level gender bias, we identify several limitations and areas for future work. We first focus on limitations within our proposed framework. First, while our results in §5 suggest that adversarial training does help reduce the influence of latent confounding variables, the analysis in §6 suggests that there is scope for improvement. Furthermore, while we focus on some confounds in the data, there may be additional ones that our model does not account for, such as the impact of videos, photos, or links shared with O\_TXT. Similarly, while our model uses O\_TXT for propensity matching in the training data, thus encouraging the model to encode indicators of bias, a model to classify comments as biased or unbiased should also incorporate O\_TXT when assessing test data. Additionally, we assume that all comments are directly addressed to OW, but some comments may be addressed to other commenters. Finally, our assumption that human judgements are not reliable for this task makes evaluation difficult, and this task

would benefit from the development of additional evaluation metrics.

There are additional avenues for future work beyond our proposed framework. Notably, we focus on the perspective of OW and examine what bias social media users may be exposed to, i.e. what comments men and women might expect to receive in response to their posts. We do not examine why comments addressed toward men and women may differ, whether because the same commenters write different comments to men and women, or because men and women attract comments from different types of people. This perspective would require controlling for traits of the commenter, such as gender, age, and occupation. Nevertheless, our work stands without this perspective: biased comments are harmful to the recipient, regardless of who wrote them.

## 9 Conclusions

Bias detection is useful for fostering civil communication on social media, as it allows recipients to screen out biased comments. Further, our intention is to detect implicit bias that people may not know they have - revealing these biases to social media users could proactively prevent them from posting unintentionally biased comments. Detecting and analyzing bias is a first step towards mitigating it, and we hope our work will encourage future work in this area.

## Acknowledgements

We would like to thank reviewers and area chairs, as well as Vidhisha Balachandran, Amanda Coston, Xiaochuang Han, Sachin Kumar, Artidoro Pagnoni, Chan Young Park, and Shuly Wintner for their helpful feedback on this work. This material is based upon work supported by the NSF Graduate Research Fellowship Program under Grant No. DGE1745016, the Google PhD Fellowship program, NSF grants IIS1812327 and SES1926043, an Okawa Grant, and the Public Interest Technology University Network Grant No. NVF-PITU-Carnegie Mellon University-Subgrant-009246-2019-10-01. We would also like to thank Amazon for providing GPU credits. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the NSF.

## References

- John Bargh. 1999. [The cognitive monster: The case against the controllability of automatic stereotype effects](#). *Dual-process theories in social psychology*, pages 361–382.
- Irene V. Blair. 2002. [The malleability of automatic stereotypes and prejudice](#). *Personality and Social Psychology Review*, 6(3):242–261.
- Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. [Man is to computer programmer as woman is to homemaker? Debiasing word embeddings](#). In *Proc. of NeurIPS*, pages 4349–4357.
- Luke Breitfeller, Emily Ahn, David Jurgens, and Yulia Tsvetkov. 2019. [Finding microaggressions in the wild: A case for locating elusive phenomena in social media posts](#). In *Proc. of EMNLP*, pages 1664–1674.
- Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. 2017. [Semantics derived automatically from language corpora contain human-like biases](#). *Science*, 356(6334):183–186.
- Kaytlin Chaloner and Alfredo Maldonado. 2019. [Measuring gender bias in word embeddings across domains and discovering new gender bias word categories](#). In *Proc. of ACL Workshop on Gender Bias in Natural Language Processing*, pages 25–32.
- Eshwar Chandrasekharan, Umashanthi Pavalanathan, Anirudh Srinivasan, Adam Glynn, Jacob Eisenstein, and Eric Gilbert. 2017. [You can’t stay here: The efficacy of Reddit’s 2015 ban examined through hate speech](#). In *Proc. of CSCW*, pages 1–22.
- MeiXing Dong, David Jurgens, Carmen Banea, and Rada Mihalcea. 2019. [Perceptions of social roles across cultures](#). In *Proc. of SocInfo*, pages 157–172.
- Naoki Egami, Christian J Fong, Justin Grimmer, Margaret E Roberts, and Brandon M Stewart. 2018. [How to make causal inferences using texts](#). *Working Paper*.
- Ethan Fast, Tina Vachovsky, and Michael S Bernstein. 2016. [Shirtless and dangerous: Quantifying linguistic signals of gender bias in an online fiction writing community](#). In *Proc. of ICWSM*, pages 112–120.
- Scott Friedman, Sonja Schmer-Galunder, Anthony Chen, and Jeffrey Rye. 2019. [Relating word embedding gender biases to gender gaps: A cross-cultural analysis](#). In *Proc. of ACL Workshop on Gender Bias in Natural Language Processing*, pages 18–24.
- Liye Fu, Cristian Danescu-Niculescu-Mizil, and Lillian Lee. 2016. [Tie-breaker: using language models to quantify gender bias in sports journalism](#). In *Proc. of IJCAI workshop on NLP meets Journalism*.
- Nikhil Garg, Londa Schiebinger, Dan Jurafsky, and James Zou. 2018. [Word embeddings quantify 100 years of gender and ethnic stereotypes](#). *Proc. of the National Academy of Sciences*, 115(16):E3635–E3644.
- Claudia Goldin. 1990. *Understanding the gender gap: an economic history of American women*. NBER series on long-term factors in economic development. Oxford University Press.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. [Generative adversarial nets](#). In *Proc. of NeurIPS*, pages 2672–2680.
- Anthony G Greenwald, Debbie E McGhee, and Jordan LK Schwartz. 1998. [Measuring individual differences in implicit cognition: the implicit association test](#). *Journal of personality and social psychology*, 74(6):1464.
- Xing Sam Gu and Paul R. Rosenbaum. 1993. [Comparison of multivariate matching methods: Structures, distances, and algorithms](#). *Journal of Computational and Graphical Statistics*, 2(4):405–420.
- Alexander Miserlis Hoyle, Lawrence Wolf-Sonkin, Hanna Wallach, Isabelle Augenstein, and Ryan Cotterell. 2019. [Unsupervised discovery of gendered language through latent-variable modeling](#). In *Proc. of ACL*, pages 1706–1716.
- Thomas Huckin. 2002. Critical discourse analysis and the discourse of condescension. *Discourse studies in composition*, pages 155–176.
- Kenneth Joseph, Wei Wei, and Kathleen M Carley. 2017. [Girls rule, boys drool: Extracting semantic and affective stereotypes from twitter](#). In *Proc. of CSCW*.
- David Jurgens, Libby Hemphill, and Eshwar Chandrasekharan. 2019. [A just and comprehensive strategy for using NLP to address online abuse](#). In *Proc. of ACL*, pages 3658–3666.
- Daniel Kahneman, Stewart Paul Slovic, Paul Slovic, and Amos Tversky. 1982. *Judgment under uncertainty: Heuristics and biases*. Cambridge University Press.
- 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](#). In *Proc. of ACL*.
- Nancy Krieger. 1990. [Racial and gender discrimination: risk factors for high blood pressure?](#) *Social science & medicine*, 30(12):1273–1281.
- Sachin Kumar, Shuly Wintner, Noah A. Smith, and Yulia Tsvetkov. 2019. [Topics to avoid: Demoting latent confounds in text classification](#). In *Proc. of EMNLP*, pages 4153–4163.

- Brian K. Lee, Justin Lessler, and Elizabeth A. Stuart. 2010. [Improving propensity score weighting using machine learning](#). *Statistics in Medicine*, 29(3):337–346.
- Christine Logel, Emma C. Iserman, Paul G. Davies, Diane M. Quinn, and Steven J. Spencer. 2009. [The perils of double consciousness: The role of thought suppression in stereotype threat](#). *Journal of Experimental Social Psychology*, 45(2):299 – 312.
- Burt L. Monroe, Michael P. Colaresi, and Kevin M. Quinn. 2008. [Fightin’ words: Lexical feature selection and evaluation for identifying the content of political conflict](#). *Political Analysis*, 16(4):372–403.
- Supun Chathuranga Nakandala, Giovanni Luca Ciampaglia, Norman Makoto Su, and Yong-Yeol Ahn. 2017. [Gendered conversation in a social game-streaming platform](#). In *Proc. of ICWSM*, pages 163–171.
- James W. Pennebaker, Martha E. Francis, and Roger J. Booth. 2001. Linguistic inquiry and word count: Liwc 2001. *Mahway: Lawrence Erlbaum Associates*.
- Margaret E. Roberts, Brandon M. Stewart, and Richard A. Nielsen. 2020. Adjusting for confounding with text matching. *American Journal of Political Science (Forthcoming)*.
- Paul R. Rosenbaum. 1988. [Sensitivity analysis for matching with multiple controls](#). *Biometrika*, 75(3):577–581.
- Paul R. Rosenbaum and Donald B. Rubin. 1983. [The central role of the propensity score in observational studies for causal effects](#). *Biometrika*, 70(1):41–55.
- Paul R. Rosenbaum and Donald B. Rubin. 1985. [Constructing a control group using multivariate matched sampling methods that incorporate the propensity score](#). *The American Statistician*, 39(1):33–38.
- Maarten Sap, Saadia Gabriel, Lianhui Qin, Dan Jurafsky, Noah A. Smith, and Yejin Choi. 2020. Social bias frames: Reasoning about social and power implications of language. In *Proc. of ACL*.
- Natalie Schluter. 2018. [The glass ceiling in NLP](#). In *Proc. of EMNLP*, pages 2793–2798.
- Claude M. Steele and Joshua Aronson. 1995. [Stereotype threat and the intellectual test performance of african americans](#). *Journal of personality and social psychology*, 69(5):797.
- Elizabeth Stuart. 2010. [Matching methods for causal inference: A review and a look forward](#). *Stat Sci*, 25(1):1–21.
- Derald Wing Sue. 2010. *Microaggressions in everyday life: Race, gender, and sexual orientation*. John Wiley & Sons.
- Tony Sun, Andrew Gaut, Shirlyn Tang, Yuxin Huang, Mai ElSherief, Jieyu Zhao, Diba Mirza, Elizabeth Belding, Kai-Wei Chang, and William Yang Wang. 2019. [Mitigating gender bias in natural language processing: Literature review](#). In *Proc. of ACL*, pages 1630–1640.
- Victor Veitch, Dhanya Sridhar, and David M. Blei. 2019. [Using text embeddings for causal inference](#). *arXiv preprint arXiv:1905.12741*.
- Rob Voigt, David Jurgens, Vinodkumar Prabhakaran, Dan Jurafsky, and Yulia Tsvetkov. 2018. [RtGender: A corpus for studying differential responses to gender](#). In *Proc. of LREC*.
- Claudia Wagner, David Garcia, Mohsen Jadidi, and Markus Strohmaier. 2015. [It’s a man’s Wikipedia? assessing gender inequality in an online encyclopedia](#). In *Proc. of ICWSM*.
- Zijian Wang and Christopher Potts. 2019. [TalkDown: A corpus for condescension detection in context](#). In *Proc. of EMNLP*, pages 3711–3719.
- Funk MJ, Westreich D, Lessler J. 2010. [Propensity score estimation: neural networks, support vector machines, decision trees \(CART\), and meta-classifiers as alternatives to logistic regression](#). *J Clin Epidemiol*, 63(8):826–833.
- Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018. [Gender bias in coreference resolution: Evaluation and debiasing methods](#). In *Proc. of NAACL*, pages 15–20.

## A Data and Model Implementation Details

	Politicians	Pub. Figures
Raw train size	6.9M	4.2M
Test/Dev size	2.3M/2.5M	1.9M/0.55M
% M in train	71.3%	33.9%
Matched train size	256K	77K
Raw dem. dim.	240	63
Matched dem. dim.	239	60

Table 4: Data Statistics. “Matched train size” refers to the size of the training set after propensity matching, and “dem. dim.” refers to the size of the latent confound vector that is demoted during training.

All data is lowercased and tokenized, and we discard data points with fewer than 4 tokens. Table 4 reports details of our data set after preprocessing.

For the primary prediction models, we use the same architectures as Kumar et al. (2019), including training multiple (2) adversaries. We perform minimal hyper-parameter tuning, primarily using the same parameters as Kumar et al. (2019), with the exception of the learning rate, which we changed slightly to decrease fluctuations in validation accuracy, and the number of training epochs for each phase of the model, which we increased or decreased as needed based on how long the validation accuracy improved for. These changes were determined by manual tuning over < 10 trials. For the propensity score model, we use a learning rate of 1e-3. For all other models we use a learning rate of 1e-4. For the models without confound demotion, we train for 5 epochs. For the models with confound demotion, we train the classifier for 3 epochs, the adversary for 10 epochs, and we repeat the alternating cycle for 3 epochs. For all models, we choose the best model as measured by W\_GEN classification accuracy over the validation set. Each model was trained using 1 GPU. The models without latent confound demotion and the propensity score estimation model have 4.2M parameters each. The adversary in the latent confound demotion models adds an additional 61.7K parameters to the Politicians model and 16.2K parameters to the Public Figures model.

## B Additional Evaluation Metrics

Table 5 provides the same results as Table 1, with the addition of precision and recall scores. Table 6 shows results for the same experiments as Table 5,

but provides metrics over the validation sets instead of the test sets. Table 7 extends Table 2 by additionally showing precision and recall scores.

	Prec.	Rec.	F1	Acc.
<b>Public Figures</b>				
base	67.3	84.5	74.9	63.8
+demotion	67.8	<b>86.7</b>	<b>76.1</b>	<b>65.1</b>
+match	65.9	65.0	65.4	56.0
+match+demotion	<b>69.0</b>	67.5	68.2	59.7
<b>Politicians</b>				
base	24.0	22.4	23.2	73.2
+demotion	<b>24.8</b>	13.4	17.4	<b>77.1</b>
+match	18.8	<b>58.8</b>	28.5	46.7
+match+demotion	19.5	54.4	<b>28.8</b>	51.4

Table 5: Evaluation over held-out test sets, where W\_GEN = F is considered the positive class, extending Table 1 by showing precision and recall.

	Prec.	Rec.	F1	Acc.
<b>Public Figures</b>				
base	61.4	76.8	68.3	57.6
+demotion	61.6	79.5	69.4	58.5
+match	61.0	61.2	61.1	53.8
+match+demotion	64.1	55.6	59.5	55.2
<b>Politicians</b>				
base	24.0	20.1	21.9	68.5
+demotion	26.3	13.0	17.4	72.9
+match	21.6	55.9	31.2	45.9
+match+demotion	22.8	53.3	31.9	50.2

Table 6: Evaluation over validation sets, where W\_GEN = F is considered the positive class, provided for reproducibility.

## C Microaggressions Data Set

The dataset of microaggressions is taken from Breittfeller et al. (2019), who collected the corpus from [www.microaggressions.com](http://www.microaggressions.com). On this website, posters describe a microaggression that they experienced. They can use quotes, transcripts, or narrative text to describe the experience, and these posts are tagged with type of bias expressed, such as “gender”, “ableism”, “race”, etc. We discard all posts that contain only narrative text, since it is not 2<sup>nd</sup> person perspective and thus very different than our training data, which leaves 1,604 posts for analysis.



	<b>Prec.</b>	<b>Rec.</b>	<b>F1</b>	<b>Acc.</b>
<b>Public Figures Training Data</b>				
base	50.9	77.0	61.3	57.3
+demotion	51.3	79.0	<b>62.2</b>	57.9
+match	49.7	32.0	38.9	55.9
+match+demotion	51.0	50.7	50.9	57.0
<b>Politicians Training Data</b>				
base	66.0	37.8	48.1	64.2
+demotion	59.6	48.9	53.7	61.5
+match	44.5	49.6	46.9	50.7
+match+demotion	45.7	75.3	<b>56.9</b>	49.9
Random	43.5	48.7	46.0	49.8
Class Random	41.4	42.9	42.1	48.3

Table 7: Evaluation over the microaggressions data set, extending Table 2 by showing precision and recall.