

# Removing Gender Information in BERT with Analytical Solution

Anonymous EMNLP submission

## Abstract

As one of the representatives of context word embedding, BERT has achieved the most advanced performance on many NLP tasks. Due to the strong feature extraction ability and the high demand for the amounts of training data, BERT can hardly avoid learning many of the human-generated stereotypes in the text data, including gender bias. In this paper, we (1) propose a template based method that is well suited to quantify gender bias in language models; (2) adapt DensRay (a vector space projection analysis method) to contextualized embeddings, and use this method to eliminate gender information. (3) investigate how English training data can be used to remove gender bias in Chinese using multilingual BERT.

## 1 Introduction

Word embeddings, which represent the semantic meaning of text data as vectors, are used as input in natural language processing tasks. It has disclosed that word embeddings exhibit unexpected social biases, such as gender bias, present in their training corpora (Bolukbasi et al., 2016; Caliskan et al., 2017; Garg et al., 2018). An example is that man is associated with computer programmer on the embedding space, and woman is associated with homemaker (Bolukbasi et al., 2016). Contextual word embedding models, such as BERT (Devlin et al., 2018), have become increasingly common and achieved new state-of-the-art results

in the many NLP tasks. Researches have also found gender bias in contextualized embeddings (Zhao et al., 2019; May et al., 2019).

In this work, we aim to mitigate gender bias on BERT embedding in a straight-forward and interpretable way. We introduce a debiasing method on BERT using DensRay (Dufter and Schütze, 2019), which yields interpretable dimensions by rotating the embedding spaces. We show that gender information is captured in every BERT layer. We apply DensRay to every BERT layer and evaluate two tasks: a set of templates we constructed and the Word Embedding Association Test (WEAT) (Caliskan et al., 2017). Our experiments find that the DensRay debiasing method effectively mitigates gender bias, while at the same time maintains the performance of BERT on language modeling and the GLUE tasks (Wang et al., 2018). As an extension, we also applied this debiasing method to the multilingual-BERT (mBERT) model: we use English gender label for computing the rotation matrix, and debias on our Chinese templates. Our contributions are summarized as the following:

- i) We adapt the analytical method DensRay, which was designed for static embeddings, to contextualized embeddings.
- ii) We demonstrate that DensRay is effective for removing gender information.
- iii) We show that the DensRay debiasing method can be applied to mBERT for zero-shot debiasing for other languages.

## 2 Background

### 2.1 Quantifying gender bias

A typical way to measure gender bias is to evaluate on downstream tasks. For coreference resolution, Zhao et al. (2018a) designed WinoBias and Rudinger et al. (2018) designed WinoGender schemas. Different from the WinoBias, Winogender schemas include gender-neutral pronouns which WinoBias does not, and one Winogender schema has one occupational mention and one “other participant” mention while WinoBias has two occupational. Webster et al. (2018) released GAP, a balanced corpus of Gendered Ambiguous Pronouns. Gender bias can be measured as the ratio of F1 score on masculine to F1 score on feminine, however the bias ratios are too close to 1 (Chada, 2019; Attree, 2019), so that the bias can’t be presented obviously. For sentiment analysis, Equity Evaluation Corpus (EEC) dataset (Kiritchenko and Mohammad, 2018) was designed to measure gender bias by the difference in emotional intensity predictions between gender-swapped sentences.

Another kind of method to measure gender bias is based on the association test, which originated from sociological research. Greenwald et al. (1998) proposed the Implicit Association Test (IAT) to quantified societal bias. In the IAT, response times were recorded when subjects were asked to match two concepts. For example, subjects were asked to match black and white names with “pleasant” and “unpleasant” words. Subjects tended to have shorter response times for concepts they thought associated. Based on the IAT, Caliskan et al. (2017) proposed the Word Embedding Association Test (WEAT), which used word similarities between targets and attributes instead of the response times to get rid of the reequipment of human subjects. Later, May et al. (2019) extended WEAT to the Sentence Embedding Associ-

ation Test (SEAT); Kurita et al. (2019) proposed a template-based log probability bias score to measure the association between targets and attributes in BERT.

#### 2.1.1 Word Embedding Association Test

Here we introduce WEAT in detail. Consider two sets of target words  $X_1$  and  $X_2$  where  $|X_1| = |X_2|$  (equal size) and two sets of attribute words  $A_1$  and  $A_2$  where  $|A_1| = |A_2|$ . As a statistical test, the null hypothesis of WEAT is: There is no difference in the cosine similarity between  $X_1, X_2$  and  $A_1, A_2$ . Taking the measurement of gender bias as an example, word sets about science and art can be used as the two target sets, masculine and feminine names can be used as the two attribute sets, such that the null hypothesis means science and art are equally similar to each masculine and feminine names, so there is no gender bias.

The test statistic is defined as,

$$s(X_1, X_2, A_1, A_2) = \sum_{x \in X_1} (x, A_1, A_2) - \sum_{x \in X_2} (x, A_1, A_2),$$

where

$$s(x, A_1, A_2) = \text{mean}_{a \in A_1} \cos(\vec{x}, \vec{a}) - \text{mean}_{a \in A_2} \cos(\vec{x}, \vec{a}).$$

Intuitively,  $s(x, A_1, A_2)$  measures the association of a word with the attributes, so the test statistic measures the differential association of the two target sets with the attributes.

Let  $\{(X_{1i}, X_{2i})\}_i$  denote all the partitions of  $X \cup Y$ . The one-sided  $p$ -value of the permutation test is defined as

$$Pr_i[s(X_{1i}, X_{2i}, A_1, A_2)] > s(X_1, X_2, A_1, A_2)$$

The effect size  $d$ -value is a normalized measure of how separated the two distributions of associations between the target and attribute are. It is defined as

$$\frac{s(X_1, X_2, A_1, A_2)}{\text{std}_{x \in X_1 \cap X_2} s(x, A_1, A_2)}.$$

## 2.2 Debiasing methods

Researchers proposed various methods to remove gender bias, in which the most common way is to define a gender direction (or, more generally, a subspace) by a set of gendered words, and debias the word embeddings in post-processing projecting. Bolukbasi et al. (2016) proposed a hard debiasing method where they used the gendered words to compute the difference embedding vector as the gender direction, and a machine learning based soft debiasing method which combined the inner-products objective of word embedding and an objective to project the word embedding into a subspace that orthogonal to the gendered words while it performed not so good as the hard debiasing method. Dev and Phillips (2019) explored partial projection and some simple tricks to improve the hard debiasing method. Zhao et al. (2019) applied the data augmentation and hard debiasing method of Bolukbasi et al. (2016) to mitigate gender bias on ELMo (Peters et al., 2018). Karve et al. (2019) introduced the debiasing conceptor, in which they shrined each principal component of the covariance matrix of the word embeddings to achieved a soft debiasing. Besides the above post-processing methods, (Zhao et al., 2018b) proposed GN-Glove which debias during training to learn word embedding with protected attributes. In this work, we have the same idea as hard debiasing, both of which are to find and eliminate gender subspace in post-processing, but we make the process more concise by the analytical solution of DensRay.

## 2.3 DensRay

**1 <<< pd: I think this section could be shortened >>>** DensRay is an analytical method proposed to identify the embedding subspace of linguistic features. Same as the methods mentioned in the previous section, we aim to identify the gender bias subspace

using a set of gendered words with vocabulary  $V := \{v_1, v_2, \dots, v_n\}$  and embedding matrix  $E \in R^{n \times d}$ , thus for word  $v_i$  we have the corresponding embedding vector  $e_{v_i}$ . We denotes the gendered words into a map  $l : V \rightarrow \{-1, 1\}$  (e.g.  $l(father) = 1, l(sister) = -1$ ). The objective of DensRay is to find an orthogonal matrix  $Q \in R^{d \times d}$  such that  $EQ$  is gender-interpretable, specifically, the first  $k$  dimensions can be interpreted as the "gender subspace".

Consider  $L_+ := \{(v, w) \in V \times V | l(v) = l(w)\}$  and  $L_- := \{(v, w) \in V \times V | l(v) \neq l(w)\}$ , we defined  $d_{vw} := e_v - e_w$  as the difference vector of  $v$  and  $w$ . DensRay solves the following optimization problem,

$$\max_q \sum_{(v,w) \in L_-} \alpha_- \|q^T d_{vw}\|_2^2 - \sum_{(v,w) \in L_+} \alpha_+ \|q^T d_{vw}\|_2^2 \quad (1)$$

where  $q \in R^d$  and  $q^T q = 1$  since  $Q$  is orthogonal,  $\alpha_+, \alpha_- \in [0, 1]$  are hyperparameters. Intuitively the objective tries to maximize the distance of the word pairs from the same gender group (male words or female words) and minimize the distance of the word pairs from different gender group. Regard that  $\|x\|_2^2 = x^T x$ , objective 1 can be simplified to:

$$\begin{aligned} \max_q q^T \left( \sum_{(v,w) \in L_-} \alpha_- \|d_{vw} d_{vw}^T\|_2^2 - \sum_{(v,w) \in L_+} \alpha_+ \|d_{vw} d_{vw}^T\|_2^2 \right) q \\ =: \max_q q^T A q \end{aligned} \quad (2)$$

The objective 2 is maximizing the Rayleigh quotient of  $A$  and  $q$ . Since  $A$  is symmetric, instead of training the model by gradient decent, we can get an analytical solution  $q$  by the eigenvector with the max eigenvalue of  $A$  (Horn et al., 1990). Thus the matrix of  $k$  eigenvectors of  $A$  ordered by the corresponding eigenvalues yields the matrix  $Q$ .

### 3 Methodology

#### 3.1 DensRay debiasing on BERT

While DensRay has been successfully used for static embeddings. We now propose to apply it to contextualized representations. Given a gendered word list  $V := \{v_1, v_2, \dots, v_n\}$ , we feed sentences which contain the gendered words into BERT without applying any masking. Due to the different contexts, for each word we regard every occurrence in a corpus as  $v_{ij}$ , with the corresponding embedding  $e_{ij}$ . For each sentence that contains the word  $v_{ij}$ , the output of each BERT layer yields a contextual embedding  $e_{ij}^m, m \in \{1, 2, \dots, M\}$ , where  $M$  is the number of layers in the BERT model. We compute an orthogonal rotation matrix  $Q_m$  for the  $m$ th BERT layer as given in formula 2. We believe that by replacing parameters in the gender subspace by 0, the gender information that may cause bias can be mitigated. In this paper, we take the first dimension of the rotated space as the gender subspace.

#### 3.2 Evaluations

In this paper we will use WEAT to measure gender bias. Besides, to quickly test the results of debiasing, we also constructed a set of templates with 320 occupation names<sup>1</sup> provided by Bolukbasi et al. (2016) for masked language modeling: "[MASK] is a *occupation*." The gender bias in the templates is measured by the average difference between the probability of BERT predicts [MASK] as "he" and "she",

$$\text{diff} = \frac{\text{mean}}{i \in \text{templates}} (\text{prob}(\text{he}) - \text{prob}(\text{she}))$$

Through case study, we find that for most templates the probability of "he" is higher than "she", which qualitatively indicates that gender bias exists in

<sup>1</sup><https://github.com/tolga-b/debiaswe/blob/master/data/professions.json>

these templates. We also find that for most sentences the sum of the two probability is higher than 0.7, which means that the predictions will be stable. This templates set can be easily extended to other languages that do not have gendered profession names, like Chinese or Persian.

### 4 DensRay Debiasing Experiments on BERT Layers

#### 4.1 Experiments Setup

In the experiments we use the BERT models "bert-base-uncased" and "bert-large-uncased". We implemented all experiments using the transformers library (Wolf et al., 2019).

To compute the rotation matrices by DensRay, we need a gendered word list as label, and some corpus. For the word list, we get 23 masculine words and 23 feminine words from the "family" category<sup>2</sup> of the Google analogy test set (Mikolov et al., 2013), and label them as 1 and -1. As the input corpus, we collect text data from Wikipedia that contains 5,000 (10,000) occurrences of words in the gendered list. We carefully balance the occurrences such that the number of male and female samples are equal. We set  $\alpha_{\neq} = \alpha_{=} = 0.5$ , as we have balanced the training samples from the corpus.

#### 4.2 Results on Templates

Results about our experiments on the templates are summarized in table Table 1. Two example templates are given in table Table 2. The evaluation on our templates shows that DensRay can mitigate the gender bias on BERT: the average difference between predicting he/she drops by more than half (e.g., bert-base from 0.47 to 0.17)

<sup>2</sup><http://download.tensorflow.org/data/questions-words.txt>

model	prob(he)	prob(she)	diff	var
bert-base	0.66	0.19	0.47	0.16
bert-base-densray	0.51	0.34	0.17	0.01
bert-large	0.63	0.19	0.44	0.13
bert-large-densray	0.48	0.29	0.18	0.02

Table 1: BERT debiasing results on templates. *bert-base* and *bert-large* are the original model without debiasing. *prob(he)* is the mean probability that model predict *he* as the [MASK] in all templates. *var* is the variance of the differences between the probability of BERT predicts [MASK] as *he* and *she*.

sentence	model	prob(he)	prob(she)
[MASK] is a adjunct professor.	bert-base	0.72	0.19
	bert-base-densray	0.44	0.47
	bert-large	0.72	0.22
	bert-large-densray	0.40	0.53
[MASK] is a administrator.	bert-base	0.63	0.23
	bert-base-densray	0.50	0.38
	bert-large	0.65	0.23
	bert-large-densray	0.45	0.37

Table 2: Two example templates with prediction probabilities.

### 4.3 Results on WEAT

In WEAT we measure the effect size *d*-value and the one-side *p*-value of the permutation test. A higher absolute value of the *d*-value indicates larger gender bias between the target words with respect to the attribute words. So, for the *d*-value, the closer to zero, the less gender bias. Refer to the definition of the null hypothesis, if the *p*-value is less than 0.05 we will reject the null hypothesis so that there will be a significant gender bias. So, we would prefer a high *p*-value (at least 0.05) to indicate the lack of gender bias. Follow the same WEAT word lists setup as Karve et al. (2019), the results on WEAT is shown on table 3. For all the three categories, DensRay decreased absolute value of *d*-value and increased the *p*-value, although bert-large still showed strong bias in (*Career, Family*) vs (*Male, Female*) even after debiasing.

category	model	d	p
(Career, Family) vs (Male, Female)	bert-base	0.66	0.08
	bert-base-densray	0.64	0.11
	bert-large	1.57	0.00*
	bert-large-densray	1.00	0.02*
(Math, Arts) vs (Male, Female)	bert-base	0.60	0.11
	bert-base-densray	0.07	0.45
	bert-large	0.22	0.35
	bert-large-densray	-0.01	0.48
(Science, Arts) vs (Male, Female)	bert-base	0.78	0.08
	bert-base-densray	0.02	0.49
	bert-large	0.82	0.04*
	bert-large-densray	0.67	0.10

Table 3: BERT debiasing results on WEAT. \* shows significant gender bias.

### 4.4 Impact on Model Performance

It is crucial that debiasing methods do not harm downstream performance of BERT models. Thus we test the perplexity of language modeling on the Wikitext-2 dataset (Merity et al., 2016) which is a subset of Wikipedia with 2 million words, the results in table Table 4 show that DensRay caused a small increase in perplexity on Wikitext-2 for both BERT base and large model.

model	ppl
bert-base	3.77
bert-base-densray	3.81
bert-large	3.29
bert-large-densray	3.35

Table 4: Language modeling performance on BERT after debiasing with DensRay.

Following the same setup as Wolf et al. (2019)<sup>3</sup>, we also evaluate on the GLUE tasks (Wang et al., 2018), results are summarized in table 5.

### 4.5 Discussions

#### 4.5.1 Number of training samples

Through evaluation and inspection of the impact on the performance of downstream tasks, ex-

<sup>3</sup><https://huggingface.co/transformers/>

model	CoLA	SST-2	MRPC	STS-B	QQP	MNLI	QNLI	RTE	WNLI
bert-base	49.23	91.97	89.47	83.95	84.31	80.61	87.36	61.73	52.11
bert-base-densray	48.32	91.86	84.73	82.46	78.94	87.01	87.30	63.90	54.93
bert-large	47.93	94.90	89.30	87.60	72.10	86.70	92.70	70.10	65.10
bert-large-densray	48.91	94.02	88.84	85.630	70.54	86.24	90.61	67.78	64.48

Table 5: GLUE tasks performance on BERT with/without debiasing with DensRay.

periments show that DensRay is an effective debiasing method on BERT. Although DensRay is an analytical solution, the effect still depends on size of the training data. In the experiments, we regarded the occurrences of the same word in the corpus as independent words with the same gender label, and used balanced samples for masculine and feminine words. Now we analyze the impact of these processes.

Since there are only 46 words in the gendered word list, if we average their embedding under different contexts, there will be only 46 training samples left for DensRay to calculate. DensRay is essentially a supervised learning method. In the case of insufficient labels, it is difficult for supervised learning to extract useful features. Treating different occurrences as different words greatly enriches training samples. As shown in figure, the debiasing results improve with an increased number of training samples.

The same as other projection-based debiasing methods (Bolukbasi et al., 2016; Zhao et al., 2019; Dev and Phillips, 2019; Karve et al., 2019), the premise of DensRay debiasing is that the bias direction should be correct. If the sample is unbalanced, the bias direction computed by DensRay will be biased towards either the male or the female, resulting in deleting the gender subspace during debiasing will reverse the gender bias (e.g. there are more masculine words in unbalanced text data, thus the embeddings will be biased towards female after biased). The figure also shows that balanced training sample improved the debiasing performers.

#### 4.5.2 Balancing Gender Bias

In this experiment, we used the method of removing the first dimension (replacing its value by 0) of the gender interpretable subspace to remove gender bias. Here we explore some other ways.

We explored two other ways to remove bias. The first is to replace the first dimension of the gender interpretable subspace with the mean value of the first dimension of the training samples. The second way is to standardize the first dimension. The results showed that both of these methods did not perform well. We further checked the mean and found that the mean of the different layers is not stable around 0, which is a problem worthy for further exploring. We also tried to delete more dimensions. However removing more dimensions does not improve the debiasing results significantly.

#### 4.5.3 Debiasing on different BERT layers

Here we only apply DensRay on one BERT layer at a time. We constructed a figure Figure ?? to illustrate the results of layers on our templates and the three WEAT categories.

### 5 DensRay Debiasing multilingual-BERT

#### 5.1 Setup

As an extension, we apply DensRay to mBERT for zero-shot debiasing on Chinese. Here we use the "bert-multilingual-uncased" model from (Wolf et al., 2019), we also use the same setup as the "bert-base-uncased" model in our previous experiments.

Figure 1: Here should be a graph.

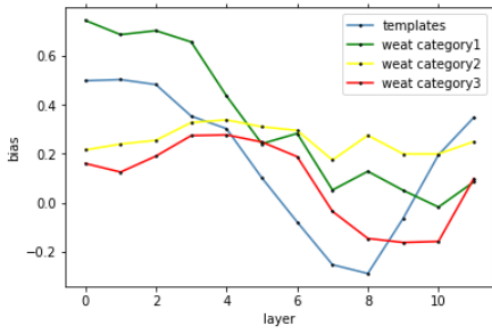


Figure 2: Debiasing on each single layer on BERT base.

As before, we compute the rotation matrices using the English gendered words from the “family” category of the Google analogy test set (Mikolov et al., 2013).

Since Chinese is a language that does not contain genus, we can construct the templates by directly translating from the English templates. So we got the following template: “[MASK]是一个*occupation*。”. For the occupation name, we referred to Tencent Translation<sup>4</sup> and made some manual adjustments to the translation. After removing the duplicates, we got 302 Chinese templates.

## 5.2 Results on templates

Results about our experiments on the templates are summarized in table Table 6. Two example templates are given in table Table 6. The evaluation on our templates shows that DensRay can mitigate the gender bias on BERT.

We also checked the perplexity for mBERT on Wikitext-2, see table Table 7. Results show that DensRay can be extended to mBERT as a zero-shot debiasing method for some other languages.

<sup>4</sup><https://fanyi.qq.com/>

model	prob(he)	prob(she)	diff	var
bert-multi-en	0.51	0.14	0.36	0.06
bert-multi-densray-en	0.33	0.12	0.21	0.03
bert-multi-cn	0.24	0.07	0.17	0.02
bert-multi-densray-cn	0.12	0.04	0.08	0.01

Table 6: Results of templates on mBERT after applied DensRay. Models with *-en* are tested on our English templates, and those with *-cn* are tested on our Chinese templates.

model	ppl
bert-multi	3.58
bert-multi-densray	3.72

Table 7: Language modeling performance on mBERT after applied DensRay.

## 6 Conclusion

We introduced DensRay debiasing on BERT. Rather than training the model as common machine learning approaches, DensRay provides an analytical solution. With this interpretable method, we can debias in BERT straight-forward. Our experiments show that this method can effectively mitigate gender bias in BERT on our constructed templates and WEAT. By checking the perplexity on Wikitext-2 and the performers on GLUE tasks, we also found this method causes little loss to the model performance. We also extend this method to mBERT as zero-shot debiasing for Chinese. As to further research, we plan to explore the irregularity of the central point of the gender dimension found in the experiments. In addition, this method can also be extended to other linguistic features, which will also be one of the future works.

## References

Sandeep Attree. 2019. Gendered ambiguous pronouns shared task: Boosting model confidence by evi-

sentence	model	prob(他)	prob(她)
[MASK]是一个客座教授。	bert-multi-en	0.68	0.16
	bert-multi-densray-en	0.51	0.18
	bert-multi-cn	0.52	0.11
	bert-multi-densray-cn	0.30	0.08
[MASK]是一个管理员。	bert-multi-en	0.53	0.17
	bert-multi-densray-en	0.35	0.13
	bert-multi-cn	0.68	0.16
	bert-multi-densray-cn	0.51	0.18

Table 8: Sanity check on the Chinese templates, where 他 means *he* and 她 means *she*. The two sentences are translated from Table 2.

- dence pooling. *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*.
- Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Advances in neural information processing systems*, pages 4349–4357.
- Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186.
- Rakesh Chada. 2019. Gendered pronoun resolution using bert and an extractive question answering formulation. *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*.
- Sunipa Dev and Jeff Phillips. 2019. Attenuating bias in word vectors.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Philipp Dufter and Hinrich Schütze. 2019. Analytical methods for interpretable ultradense word embeddings. 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 1185–1191.
- N Garg, L Schiebinger, D Jurafsky, and J Zou. 2018. Word embeddings quantify 100 years of gender and ethnic stereotypes. *Proceedings of the National Academy of Sciences of the United States of America*, 115(16):E3635–E3644.
- 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.
- R.A. Horn, R.A. Horn, and C.R. Johnson. 1990. *Matrix Analysis*, chapter 4.2. Cambridge University Press.
- Saket Karve, Lyle Ungar, and João Sedoc. 2019. Conceptor debiasing of word representations evaluated on weat. In *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*, pages 40–48.
- Svetlana Kiritchenko and Saif Mohammad. 2018. Examining gender and race bias in two hundred sentiment analysis systems. *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*.
- Keita Kurita, Nidhi Vyas, Ayush Pareek, Alan W Black, and Yulia Tsvetkov. 2019. Measuring bias in contextualized word representations. *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*.
- Chandler May, Alex Wang, Shikha Bordia, Samuel Bowman, and Rachel Rudinger. 2019. On measuring social biases in sentence encoders. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 622–628.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2016. Pointer sentinel mixture models.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. In *Proceedings of International Conference on Learning Representations (ICLR)*.
- Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proc. of NAACL*.
- Rachel Rudinger, Jason Naradowsky, Brian Leonard, and Benjamin Van Durme. 2018. Gender bias in coreference resolution. *Proceedings of the 2018 Conference of the North American Chapter of the*



*Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers).*

Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*.

Kellie Webster, Marta Recasens, Vera Axelrod, and Jason Baldridge. 2018. Mind the gap: A balanced corpus of gendered ambiguous pronouns. *Transactions of the Association for Computational Linguistics*, 6:605–617.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and Jamie Brew. 2019. Huggingface’s transformers: State-of-the-art natural language processing.

Jieyu Zhao, Tianlu Wang, Mark Yatskar, Ryan Cotterell, Vicente Ordonez, and Kai-Wei Chang. 2019. Gender bias in contextualized word embeddings. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 629–634.

Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018a. Gender bias in coreference resolution: Evaluation and debiasing methods. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, volume 2.

Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. 2018b. Learning gender-neutral word embeddings. *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*.