

Monolingual and Multilingual Reduction of Gender Bias in Contextualized Representations

Anonymous COLING submission

Abstract

Pretrained language models (PLMs) learn stereotypes held by humans and reflected in text from their training corpora, including gender bias. When PLMs are used for downstream tasks such as picking candidates for a job, people’s lives can be negatively affected by these learned stereotypes. Prior work usually identifies a linear gender subspace and removes gender information by removing the subspace. Following this line of work, we use DensRay, an analytical method for obtaining interpretable ultradense subspaces. We show that DensRay performs on-par with prior approaches, but is more stable. In addition, DensRay can be used to obtain interpretable gender scores on token level for all representations. Finally, we demonstrate that we can remove bias multilingually, e.g., from Chinese, using only English training data.

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 been found that word embeddings exhibit biases such as gender bias, which are present in their training corpora (Wang et al., 2018). Contextual word embedding models, such as BERT (Devlin et al., 2019), have become increasingly common and achieved new state-of-the-art results in many NLP tasks. Researchers have also found gender bias in contextualized embeddings (Wang et al., 2019).

A common approach for removing gender information in static embeddings is to identify a linear gender subspace (e.g., a gender direction) and subsequently setting all values on the gender direction to 0. Successful approaches rely on simple principal component analysis (Wang et al., 2018). (Wang et al., 2018) require pairs of gendered words to compute a direction (e.g., “man”-“woman”) and (Wang et al., 2019) rely on computing a PCA of a set of gender words hoping that the main variation occurs across gender. We propose to use DensRay (Wang et al., 2019): the main advantage is that DensRay only requires two or multiple groups of gendered words. In contrast to (Wang et al., 2018), it does not require explicit pairs. Compared to (Wang et al., 2019), it has explicit supervision with gender labels. We show in §2.3 that DensRay is more stable.

In summary our contributions are: i) We adjust DensRay to work on contextualized embeddings. We apply DensRay to every BERT layer and evaluate two tasks: a set of templates we constructed and the Word Embedding Association Test (WEAT) (Wang et al., 2018). Our experiments find that debiasing with DensRay effectively mitigates gender bias and performs on par with prior approaches. ii) We show that DensRay is more robust and interpretable than prior approaches. iii) We investigate whether debiased models maintain the performance of BERT on language modeling and GLUE (Wang et al., 2018). iv) We apply our debiasing method to the multilingual-BERT (mBERT) model: we show that English training data can be used to effectively debias Chinese.

2 Methodology

2.1 Hard Debiasing

2.2 Debiasing Conceptor

2.3 DensRay

DensRay is an analytical method for identifying the embedding subspace of certain linguistic features. Similar to the methods mentioned in the previous section, we aim to identify the “gender subspace”

using a set of gendered words $V := \{v_1, v_2, \dots, v_n\}$ and their embeddings $E \in R^{n \times d}$, thus for word v_i we have the corresponding embedding vector e_{v_i} . We introduce a function l for the gender attribute: $l : V \rightarrow \{-1, 1\}$; e.g. $l(\text{father}) = 1$, $l(\text{sister}) = -1$. 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.

Let $L_+ := \{(v, w) \in V \times V | l(v) = l(w)\}$ and define L_- analogously. The DensRay objective in Eq. 1 is to maximize the distance of the word pairs from the same gender group (L_+) and minimize the distance of the word pairs from the different gender group (L_-).

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

where we define $d_{vw} := e_v - e_w$. We also have $q \in R^d$ and $q^T q = 1$ since Q is orthogonal. $\alpha_-, \alpha_+ \in [0, 1]$ are hyperparameters. Observing that $\|x\|_2^2 = x^T x$, objective Eq. 1 can be simplified to:

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

The objective in Eq. ?? is maximizing the Rayleigh quotient of A and q . Since A is symmetric, we can get an analytical solution q by the eigenvector with the max eigenvalue of A (?). Thus the matrix of k eigenvectors of A ordered by the corresponding eigenvalues yields the matrix Q .

Here we compare the difference among DensRay, debiasing conceceptor, and hard debiasing. Figure 1 shows artificially created two dimensional embeddings. The lines show the gender directions identified by hard debiasing, conceceptor and DensRay. Hard debiasing and conceceptor does not use the labels of the male-female word pairs. Instead, hard debiasing relies upon the assumption that the first principal component of the considered vectors is a meaningful gender direction. This can fail in some cases.

1 <<< hs: since this is a crucial point of the paper, we need a strong summary sentence here >>>

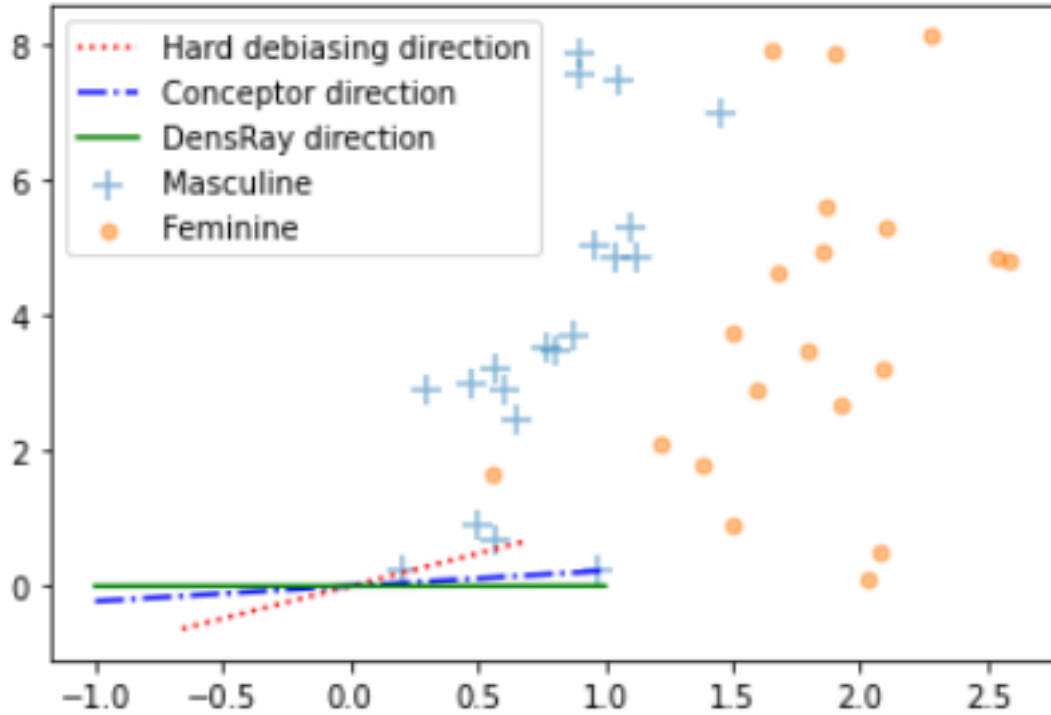


Figure 1: Gender direction on gendered words.

2.4 Adapting DensRay to Contextualized Language Models

We now describe how we adapt DensRay to contextualized language models. Given a set of gendered words V , we extract sentences containing a word in V from a corpus. We run a contextualized language model with M layers on each sentence $t_1, \dots, t_j, \dots, t_n$ (where $t_j \in V$) and compute the contextualized representations $e^m, 1 \leq m \leq M$ of t_j , one for each layer. We compute an orthogonal rotation matrix Q_m for the m th BERT layer using Eq. 2. Finally, for debiasing, we set the dimensions of the gender subspace to 0 with the goal of eliminating or at least reducing gender information that may cause bias; for measuring bias, we use the distance to the zero point of the gender subspace as the measurement. In this paper, we take the first dimension of the rotated space as the gender subspace.

3 Experiments

3.1 Setup and Data

In the experiments we downcase all text and use the BERT models “bert-base-uncased” and “bert-large-uncased”. We implemented all experiments using the transformers library (?).

To compute the rotation matrices by DensRay, we need the labels of a gendered word list and a corpus. For the word list, we get 23 masculine words and 23 feminine words from the “family” category,¹ of the Google analogy test set (?) and label them as 1 and -1. As the input corpus, we collect text data from Wikipedia that contains 5,000 (resp. 10,000) occurrences of words in the gendered list for the BERT base (resp. large) model. 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.

We compare with the hard debiasing method (?) and the debiasing conceceptor (?) to eliminate gender bias as adapted to contextualized embeddings by (?).

3.2 Evaluations

3.2.1 OCCTMP

We use two evaluation datasets to measure gender bias: WEAT (Section 3.2.2) and OCCTMP.

OCCTMP is a new evaluation dataset based on occupation templates that we created specifically for the evaluation of contextualized language models. It has the added advantage that results are easier to interpret than those for WEAT.

To construct OCCTMP, we start with 320 occupation names² provided by (?). Each occupation name is converted into a template of the form “[MASK] is an *occupation*.” We measure gender bias in the templates as the average difference between the probability of BERT predicting [MASK] as “he” vs. “she”

$$\text{diff} = \frac{1}{|\mathcal{T}|} \sum_{T \in \mathcal{T}} (p(\text{he}|T) - p(\text{she}|T))$$

where \mathcal{T} is the set of 320 templates. We find that for most experiments and most templates the probability of “he” is higher than “she”, which qualitatively indicates that gender bias can be identified using these templates. We also find that in most cases the sum of the two probability is higher than 0.7; thus, this evaluation task is a good fit for BERT because it has learned that a pronoun is likely to occur in the masked position. Our templates can be easily extended to other languages as we later show for Chinese.

3.2.2 Association Tests

An alternative way to measure gender bias is based on association tests, which originated from sociological research. (?) proposed the Implicit Association Test (IAT) to quantified societal bias. In 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.

¹<http://download.tensorflow.org/data/questions-words.txt>

²<https://github.com/tolga-b/debiaswe/blob/master/data/professions.json>

Based on IAT, (?) 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 requirement of human subjects. Consider two sets of target words X_1, X_2 with equal size $|X_1| = |X_2|$, and two sets of attribute words A_1, A_2 . The null hypothesis in the statistical test of WEAT is: there is no difference in the similarity between X_1, X_2 and A_1, A_2 . In the prior literature it has been argued that if the null hypothesis can't be rejected, there is no significant gender bias.

2 <<< pd: we should criticize this reasoning. The null and alternative hypothesis should be swapped. Has other work criticized this setup? Maybe we can do the test in addition in an alternative way? >>>
3 <<< sl: as far as I see, all the association tests(IAT,WAET,SEAT,CEAT) use this hypothesis setup, I think it's just a 'tradition'. >>> The WEAT test statistic is defined as

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

4 <<< pd: I do not understand the notation fully. Is the p-value computed with respect to a single partition i? >>> **5 <<< sl: the probability is computed over the space of partitions >>>** where

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

in which $\cos(\vec{x}, \vec{a})$ denotes the cosine similarity between embedding vector \vec{x} and \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 $X_1 \cup X_2$. The one-sided p -value of the permutation test is defined as

$$p = \text{Pr}_i[s(X_{1i}, X_{2i}, A_1, A_2) > s(X_1, X_2, A_1, A_2)]$$

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

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

To extend WEAT to contextual embeddings, (?) extracted contextual embeddings from the template '[MASK] is word.' (?) proposed Sentence Embedding Association Test (SEAT), which designed more complex templates to extract word embeddings.

Dispensed with templates, (?) proposed Contextualized Embedding Association Test (CEAT), which extracted the embeddings of the stimulus' occurrences from the corpus, and computed the weighted mean of effect sizes and statistical significance by a random-effects model. The combination effect size is

$$d_c(X_1, X_2, A_1, A_2) = \frac{\sum_{i=1}^N v_i d_i}{\sum_{i=1}^N v_i}$$

where v_i is the weights in the random-effects model. To measure the statistical significance, they used two-tailed p -value $p_c = 2[1 - \phi(|\frac{d_c}{\text{std}(d_c)}|)]$.

6 <<< hs: i think there is a summary sentence missing here: how do we use this to evaluate / compare debiasing methods? >>>

3.2.3 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 Wikitext-2 (?), a subset of Wikipedia with 2 million words. We also test on GLUE tasks (?). For all the tests we follow the same setup as (?).³

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

4 Results

4.1 Results on OCCTMP

Table 1 gives results for OCCTMP. Two OCCTMP examples are given in Table 2. It shows that DensRay can mitigate the gender bias in BERT: the average difference between predicting he/she drops to around two third (e.g., for bert-base from 0.47 to 0.11).

7 <<< hs: there is a summarizing sentenc emissing jhere: performance of hard debiaisig and densray are comparable >>>

model	prob(he)	prob(she)	diff	var
bert-base	0.66	0.19	0.47	0.16
bert-base-hard	0.35	0.42	-0.07	0.03
bert-base-conceptor	0.18	0.11	0.08	0.01
bert-base-densray	0.48	0.37	0.11	0.02
bert-large	0.63	0.19	0.44	0.13
bert-large-hard	0.40	0.23	0.17	0.02
bert-large-conceptor	0.05	0.03	0.02	0.00
bert-large-densray	0.47	0.31	0.16	0.02

Table 1: BERT debiasing results on OCCTMP. *bert-base* and *bert-large* are the original models without debiasing. *prob(he)* is the average probability predicted for *he* as the [MASK] in OCCTMP. *var* is the variance of the differences between the probabilities of predicted for *he* and *she*.

9 <<< pd: why don't we compare with conceptor anymore? >>>

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: OCCTMP examples with prediction probabilities.

10 <<< hs: since “a adjunct” and “a administrator” are not correct english, can you please find examples that are correct english? >>>

4.2 Results on WEAT

In WEAT we measure the effect size *d*-value and the onesided *p*-value of the permutation test. A *d*-value closer to zero indicates less gender bias. We also prefer a high *p*-value (at least 0.05) to not reject the null hypothesis, i.e., we do not reject that there is no gender bias. We use (?)’s WEAT word list setup. Table 3 shows results on WEAT.

11 <<< hs: there is a summarizing sentenc emissing jhere: performance of hard debiaisig and densray are comparable (although results are somewhat random as we have discussed before) >>>

category	model	WEAT		SEAT	
		d	p	d	p
C6	bert-base	0.66	0.08	0.11	0.25
	bert-base-Mu	0.15	0.38	0.48	0.10
	bert-base-densray	0.62	0.12	-0.11	0.75
C7	bert-base	0.60	0.11	0.72	0.01
	bert-base-densray	-0.07	0.56	-0.09	0.72
	bert-base-densray	0.09	0.45	-0.11	0.26
C8	bert-base	0.78	0.08	1.00	0.01
	bert-base-Mu	-0.29	0.68	0.36	0.03
	bert-base-densray	0.03	0.47	0.75	0.01

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

12 <<< hs: “* shows significant gender bias”: i don’t see any stars >>>

4.3 Impact on Model Performance

Table 4 shows that DensRay debiasing gets comparable results with the original models on Wikitext-2 and GLUE tasks.

model	Wikitext-2	CoLA	SST-2	MRPC	STS-B	RTE	WNLI
bert-base	3.77	49.15	92.09	85.86	82.66	62.82	52.11
bert-base-mu	3.95	45.53	91.74	82.48	82.60	63.54	56.34
bert-base-densray	3.81	48.04	91.74	84.89	82.43	63.90	53.52
bert-large	3.29	47.93	94.90	89.30	87.60	70.10	65.10
bert-large-Mu	3.85	47.45	93.95	85.01	82.33	67.12	63.02
bert-large-densray	3.35	48.91	94.02	88.84	85.63	67.78	64.48

Table 4: Language modelling perplexity and GLUE tasks performance.

15 <<< pd: it seems on bert-large densray is always better than Mu? Can't we make the argument that DensRay affects performance less? >>>
16 <<< hs: great idea! >>>

4.4 Discussions

4.4.1 Debiasing on Attention Heads

We now apply DensRay to the attention heads in BERT to debias on OCCTMP, The heatmap Figure 2 shows that the debiasing effect of one single attention head is not obvious, with diff scores all in $[0.4, 0.5]$. Due to the lack of dimensions and the distribution of gender features in the attention heads, we chose to apply DensRay on layers as debiasing method. We conclude that there is no single attention head which is responsible for processing gender information.

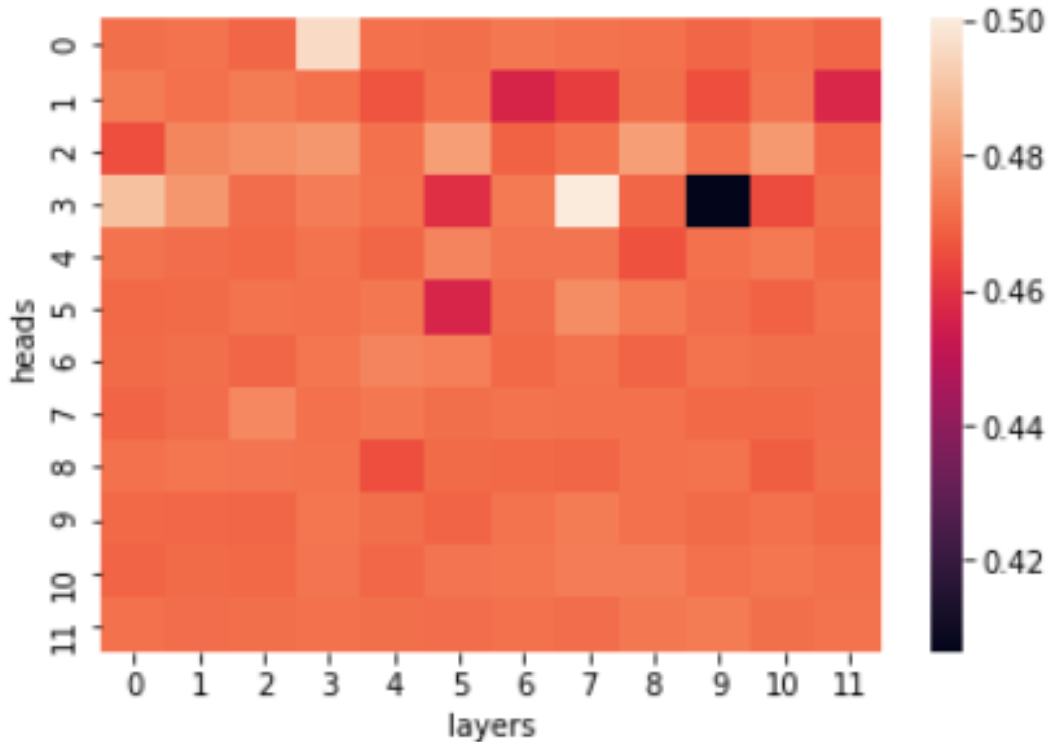


Figure 2: DensRay debiasing on each single attention head in BERT base, measured by diff on OCCTMP.

4.4.2 Number of Training Samples

In the experiments, we collected training samples for DensRay by considering occurrences of the same word in the corpus across different sentences. We collected equally many masculine and feminine words. Now we analyze the impact of these processes. 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 3, the debiasing results improve with an increased number of training samples.

Similar to other projection-based debiasing methods (Zhang et al., 2018; Wang et al., 2019; Wang et al., 2020), 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 and reversing the gender bias. For example, if there are more masculine words in unbalanced text data, then the embeddings will be biased towards female after debiasing. The figure also shows that a balanced training sample improves the debiasing performance.

17 <<< hs: i don't understand the last point: does the figure also show experimental results for balanced vs unbalanced training sets? >>>

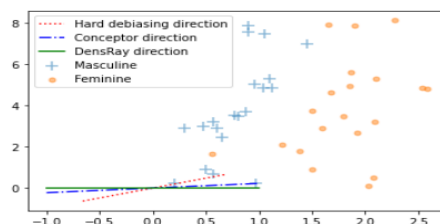


Figure 3: DensRay debiasing results on OCCTMP with different number of samples.

4.4.3 Balancing Gender Bias

18 <<< pd: I think this section can be moved to the supplementary material >>>
19 <<< pd: are we still targeting a short paper? or a long paper? >>> 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 three other ways to remove bias: 1) replace the first dimension of the gender interpretable subspace with the mean value of the first dimension of the training samples. 2) standardize the first dimension. 3) replace the first dimension with a small random variable sampled from Gaussian distribution. All of them 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, while harming the model performance significantly.

4.4.4 Debiasing across different layers

So far we have applied DensRay to all BERT layers simultaneously. Figure 4 illustrates the effect of debiasing a single layer on our templates and the three WEAT categories. We see that the debiasing effect is stronger in layers 7–10 than in the other layers in BERT base.

4.4.5 Multilingual Debiasing

We now show that, in a multilingual contextualized language model like mBERT, we can use DensRay for zero-shot debiasing. Specifically, we train a DensRay model on English and use it to debias Chinese. We use bert-multilingual-uncased from (Liu et al., 2020). We use the same setup as for bert-base-uncased in our previous experiments.

20 <<< pd: If I recall correctly the uncased model is not good for Chinese and only the cased model should be used. This is because they used different preprocessing for both models. >>>

As before, we compute the rotation matrices using the English gendered words from the “family” category of the Google analogy test set (Liu et al., 2020).

Since Chinese is a language that does not mark gender, we can construct the OCCTMP templates by directly translating from the English templates. We use the following form: “[MASK]是一个occupation。”

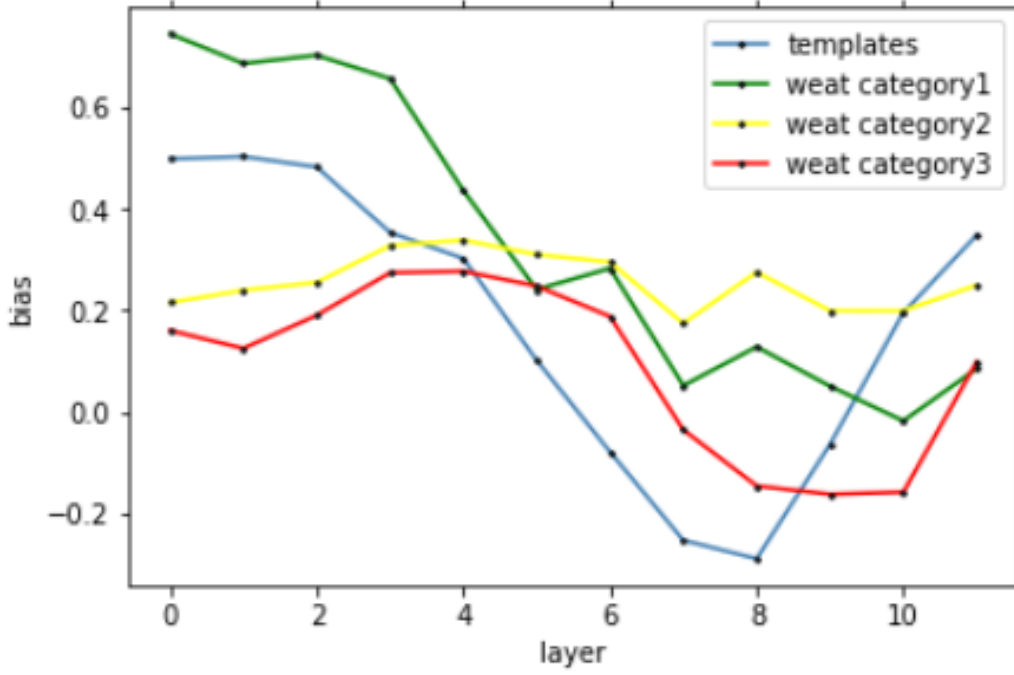


Figure 4: Debiasing on each single layer on BERT base. Bias is measured by diff on the templates and d -value on WEAT categories.

We translate the occupation name based on Tencent Translation⁴ and make some manual adjustments to the translation. After removing duplicates, 302 Chinese templates remain.

Table 5 gives results for the Chinese templates. Two examples are given in Table 5. We see that DensRay trained with English can mitigate gender bias in mBERT: the average difference drops from 0.17 to 0.08 on Chinese templates. Also, mBERT still gets comparable perplexities on Wikitext-2 after debiasing: see table Table 6.

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

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

Table 6: Language modeling performance on mBERT after applied DensRay. **22 <<< pd: on which language? Not sure whether this table is necessary in the main paper >>>**

5 Related Work

5.1 Quantifying Gender Bias

A typical way to measure gender bias is to evaluate on **downstream tasks**. For coreference resolution, (?) designed Winobias and (?) designed Winogender schemas. In contrast to WinoBias, Winogender

⁴<https://fanyi.qq.com/>

sentence	model		
	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
	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 7: Sanity check on the Chinese templates, where *he* means *he* and *she* means *she*. The two sentences are translated from Table 2.

schemas include gender-neutral pronouns. One Winogender schema has one occupational mention and one “other participant” mention while WinoBias has two occupational mentions. **23 <<< pd:**

is the difference between Winobias and Winogender relevant to this work? >>> (?)

released GAP, a balanced corpus of Gendered Ambiguous Pronouns, which measures gender bias as the ratio of F1 score on masculine to F1 score on feminine. However the ratio is very close to 1 (?: ?) making it hard to compare debiasing systems. For sentiment analysis, Equity Evaluation Corpus (EEC) (?) was designed to measure gender bias by the difference in emotional intensity predictions between gender-swapped sentences.

An alternative way to measure gender bias is based on **association tests**, which originated from sociological research. (?) proposed the Implicit Association Test (IAT) to quantify 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, (?) proposed the Word Embedding Association Test (WEAT), which uses word similarities between targets and attributes instead of the response times to get rid of the requirement of human subjects. (?) extended WEAT to the Sentence Embedding Association Test (SEAT); (?) proposed a template-based log probability bias score to measure the association between targets and attributes in BERT.

24 <<< hs: for many of the papers you discuss above it’s not clear what the relationship to the current work is. this should always be clear >>>

5.2 Debiasing Methods

Many methods to remove gender bias have been proposed. 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 a post-processing projection. (?) propose (i) *hard debiasing*: they use the gendered words to compute the difference embedding vector as the gender direction; and (ii) *soft debiasing*, a machine learning based method that combines the inner-products objective of word embedding and an objective to project the word embedding into an orthogonal gender subspace. Hard debiasing has been found to work better. **25**

<<< pd: should we mention hard-debiasing by mu et al here and explain the difference to bolukbasi? >>>

(?) explored partial projection and some simple tricks to improve the hard debiasing method. (?) applied the data augmentation and debiasing method of (?) to mitigate gender bias on ELMo (?). (?) introduce the debiasing conceptor: they shrink each principal component of the covariance matrix of the word embeddings to achieve a soft debiasing. Besides the above post-processing methods, (?) propose GN-Glove: it debiases during training to learn word embeddings with protected attributes. The method we use here, DensRay, is similar to hard debiasing in that we find and eliminate a gender subspace in post-processing. But DensRay can be solved efficiently in closed form and it is more stable than hard debiasing.

6 Conclusion

We introduced DensRay debiasing on BERT. 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.

26 <<< hs: three parts of the paper should be in sync:

(i) the abstract

(ii) the contributions at the end of the introduction

(iii) the conclusion

make sure to check that during final editing >>>

References

- 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, Apr.
- 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.
- 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.
- Wei Guo and Aylin Caliskan. 2020. Detecting emergent intersectional biases: Contextualized word embeddings contain a distribution of human-like biases.
- 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)*.
- Jiaqi Mu and Pramod Viswanath. 2018. All-but-the-top: Simple and effective post-processing for word representations. In *6th International Conference on Learning Representations, ICLR 2018*.
- 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, Dec.
- 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, 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*.
- 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.