
[RE] Bad Seeds: Evaluating Lexical Methods for Bias Measurement

Anonymous Author(s)

Affiliation

Address

email

Reproducibility Summary

1

2 **Scope of Reproducibility**

3 In this work we verify the results of *Bad Seeds: Evaluating Lexical Methods for Bias Measurement* (Antoniak and
4 Mimno, 2021). We replicate the experiments conducted and verify the main claims made in the original paper: (1) Bias
5 measurements depend on seeds and models. (2) Shuffled seed pairs can result in a significant different bias subspace
6 compared to ordered seed pairs. (3) Set similarity is negatively correlated with the explained variance of the first PCA
7 component in the seed pairings subspace.

8 **Methodology**

9 We used skip-gram with negative sampling to train word2vec models with the same hyperparameters and data. We
10 implemented code for the experiments using the resulting word embeddings and the seed sets provided by the authors.

11 **Results**

12 Overall, only one claim was reproduced. We reproduced the claim that bias measurements is dependent on the choice
13 of seed set. We were not able to adequately reproduce the claims that shuffled pairs of seed sets generally result in
14 less clearly defined correlation and that for pair of seed sets set similarity is negatively correlated with the explained
15 variance of the first principal component.

16 **What was easy**

17 The paper is easy to follow. The data was publicly available. Authors replied frequently providing details about the
18 parameters and preprocessing steps. Also authors were open for the discussion regarding seeds on the Github repository
19 of the project.

20 **What was difficult**

21 In certain cases, the gathered seed sets json file contained errors. Specifically: 'daughters' was misspelled as 'daughters'
22 (has now been updated). 'ma', 'am' was used as two words instead of one word "ma'am".
23 The seed words for figure 4 as stated in the appendix are not the same as the one in the image itself. Table 2 from the
24 original paper was also difficult to reproduce as the preprocessing according to the authors' description gave close but
25 not equal results for the NYT dataset and significantly different results for two other datasets. Reproduced numbers are
26 presented in table 3.3. Because of the time constraints, 20 bootstrapped launches were not conducted.

27 **Communication with original authors**

28 Contact was made with the original authors on multiple occasions to ask for clarification questions regarding the
29 implementation of the experiments.

30 1 Introduction

31 Using local context in datasets to generate word embeddings for NLP problems carries the risk of incorporating the
32 original bias of the dataset into the generated model. Techniques to combat such biases requires the measurement of the
33 bias encoded in a model (Bolukbasi et al., 2016). Almost all bias measurement methods rely on lexicons of seed terms
34 to specify stereotypes, while the rationale for choosing specific seeds is often unclear.

35 In this work we verify the results of *Bad Seeds: Evaluating Lexical Methods for Bias Measurement* (Antoniak and
36 Mimno, 2021). We replicate the experiments conducted and attempt to verify the main claims made in the original paper.
37

38 2 Scope of reproducibility

39 The original paper collected many different seeds (gathered and generated ones) for bias measurement and provided
40 various ways these seeds were originally selected. Then the paper conducted different experiments showing how
41 different choices can result in different bias measurements. We focus on replicating these experiments, and addressing
42 the following claims:

- 43 • Claim 1: Bias measurements depend on seeds and models. The paper shows this by calculating similarity to
44 the *Unpleasantness* vector using different seeds and different embedding models. We verify this by recreating
45 this experiment and showing different seeds and models result in different bias measurements.
- 46 • Claim 2: Shuffled pairs of seed sets generally result in less clearly defined correlation in bias subspace than
47 unshuffled pairs of seed sets.
- 48 • Claim 3: Identifying bias is less effective when set pairs are similar. This claim is verified by generating seed
49 pairs and plotting set similarity against the explained variance of the first principal component in the bias
50 subspace.

51 Finally we explore if BERT embeddings results in similar results compared to word2vec models.

52 3 Methodology

53 3.1 Bias Measurement Algorithms

54 3.1.1 PCA

55 Principal component analysis can be used to measure the variability in the difference vectors between pair of word
56 vectors. The vector that represents this difference the best is used to generate the bias subspace. The PCA algorithm is
57 applied in the following manner: for each pair of seed sets, the mean vector of their two embedding vectors is calculated.
58 The half vectors, which are calculated by subtracting the embedding vectors from the mean, are added to a list. These
59 half vectors are used as the columns in the input matrix of the PCA algorithm. (Bolukbasi et al., 2016). We can further
60 calculate the explained variance ratio of resulting principal components to determine how much variance is explained
61 by the first few principal components.

62 3.1.2 Word Embedding Association Test (WEAT)

To quantify bias, WEAT can be used to find which seed set is more associated to given attribute words (Caliskan et al., 2017). The WEAT score between sets \mathcal{X} and \mathcal{Y} , and sets of attributes, \mathcal{A} and \mathcal{B} , is defined as

$$s(\mathcal{X}, \mathcal{Y}, \mathcal{A}, \mathcal{B}) = \sum_{x \in \mathcal{X}} s(x, \mathcal{A}, \mathcal{B}) - \sum_{y \in \mathcal{Y}} s(y, \mathcal{A}, \mathcal{B})$$

where $s(w, \mathcal{A}, \mathcal{B})$ is equal to the difference in average cosine similarities between query w and each term in \mathcal{A} and \mathcal{B} . A WEAT subspace is created by calculating the vector between the average embeddings of two seed sets. This subspace is used to calculate coherence between pair of seeds. A seed pair is said to have a high coherence when the seeds are

highly separated. Coherence is calculated by ranking the mean ranks of every word in the seeds by the cosine similarity of the bias subspace:

$$\text{Coherence}(\mathcal{X}, \mathcal{Y}) = |\overline{R_{\mathcal{X}}} - \overline{R_{\mathcal{Y}}}|$$

63 where $\overline{R_{\mathcal{X}}}$ and $\overline{R_{\mathcal{Y}}}$ are the mean ranks of the seed sets in the bias subspace. These two values are normalized to be
64 between 0 and 1.

65 3.2 Datasets

66 For reproduction purposes, we used the links to the datasets mentioned in the paper. NYT was easily found on kag. The
67 original link to the WikiText-103 didn't work and a word-level dataset from wik was taken. Goodreads reviews worked
68 with the link provided.

69 3.3 Preprocessing

70 Most of the details of the preprocessing procedure were found in the paper. However, some steps became clear only
71 after reaching authors who were open to discussion.

72 The final preprocessing pipeline looked as follows:

- 73 • splitting of each dataset into logical documents:
 - 74 – articles separated by newline symbols and URLs in NYT;
 - 75 – articles separated by an article name inside "=" signs in WikiText;
 - 76 – separate json files in Goodreads;
- 77 • lowercasing;
- 78 • removing non-alphanumeric characters;
- 79 • splitting on whitespace;
- 80 • removing words that occur fewer than 10 times in the training dataset (via gensim word2vec model min_count
81 argument).

82 Apart from this, it was suggested to apply additional cleaning steps for WikiText:

- 83 • lists removal;
- 84 • HTML errors removal;
- 85 • math removal;
- 86 • code removal.

87 No details regarding this were found. Moreover, we were not able to recover \LaTeX expressions and HTML tags via
88 regular expressions. So the question about WikiText preprocessing remains open.

89 The datasets were preprocessed and summary statistics for them were calculated 3.3. It turned out that for WikiText
90 there is a decrease in the vocabulary size and the mean document length compared to the original paper. For Goodreads,
91 there is a difference in the vocabulary size, potentially caused by preprocessing nuances.

Team	Total Documents	Total Words	Vocabulary Size	Mean Document Length
NYT	8,888	7,210,433	136,404	811
WikiText	28,470	86,397,984	223,631	3,034
Goodreads (Romance)	197,000	24,116,941	291,623	122
Goodreads (History/Biography)	136,000	14,001,917	216,067	103

93 Table 2: Reproduced summary statistics for our datasets

94 3.4 Seeds

95 Algorithmic bias can be observed in machine learning models when terms associated with a particular gender, ethnicity,
96 political party or other such grouping are treated differently by the model than the average term. To quantify how much

97 bias is present in a model, seed terms are required that have a strong correlation with a particular bias dimension that
98 we wish to represent. A set of these terms, called a seed set, together form a reference for a particular bias group. A
99 seed set can be used to quantify how much bias is present in a model for a bias group, by looking at the difference in
100 outputs between the seed terms and all other terms. Clearly, it is important that the right seed terms are grouped to
101 measure bias for a particular bias group, otherwise, the bias quantification of a model might be flawed.

102 3.4.1 Gathered seeds

103 In the original paper, the authors use a collection of seed sets gathered from 18 'highly-cited papers'. We use the same
104 seed sets as used in the experiments for the original paper (Bolukbasi et al., 2016) (Kozłowski et al., 2019) (Garg et al.,
105 2018) (Caliskan et al., 2017) (Manzini et al., 2019) (Zhao et al., 2018) (Hoyle et al., 2019). Like the original work, we
106 use unigram seeds and omit words that were not present in the training set.

107 3.4.2 Generated seeds

108 Using generated seed sets is a way to get a large number of seed sets to use for experiments. To find suitable seed sets,
109 we want to find words that are close together in the embedding space. We also want to control for POS (obtained with
110 spaCy *en_core_web_sm* model is used) and frequency of the words in the dataset. The way we generate seed sets is by
111 first selecting a random noun in the vocabulary with a probability proportional to the frequency in the dataset. We check
112 if a word is a noun or not using the *spacy* library. Then we take the four nearest nouns of this word ranked by cosine
113 similarity

114 3.5 Experimental Setup

115 In this section, we will provide steps to replicate the results obtained by the original authors. We found that various
116 steps were not made explicit or used different seed sets as documented. We provide details on the steps we have taken
117 to replicate the experiments.

118 The codebase ¹ consist of two notebooks, one for training the models and one for the experiments and graph
119 generation.

120 3.5.1 Models

121 Different word embeddings have varying amounts of bias embedded inside the models. To verify the claims made
122 regarding the various effects of seeds on bias measurements, we limit the scope to the word embeddings used in the
123 original paper. For each dataset, we train a word2vec model using skip-gram with negative sampling (SGNS). Apart
124 from that, outside the scope of the original paper we tried using BERT-generated embeddings as well.

125 The word2vec models for Goodreads, NYT and WikiText are trained in *SGNS_train.ipynb* using *preprocessing.py*
126 script. The notebook consists of preprocessing of the original dataset, the training of the word2vec model and saving
127 the model. All SGNS model parameters were used with default values, namely *5 training epochs, 5 negative samples*
128 *per each positive, vector size of 100, window size of 5, word frequency not lower than 10*.

129 Script *bert_embeddings.py* was used to obtain pre-trained embeddings from a widely used *bert-base-uncased* model.
130 Devlin et al. (2018)

131 3.5.2 Experiments

132 The experiments of the original paper are reproduced in *Paper reproduction.ipynb*. To address claim 1, we replicate
133 figure 2 of the original paper to highlight how much different seeds in the same category affect the similarity to a certain
134 vector in the embedding space. This is done by calculating for both the word embedding using the romance Goodreads
135 reviews and the history + biography Goodreads reviews the cosine similarity of the words in various seed sets with the
136 word *Unpleasant* in the same word embedding. Note that we do not use the word *Unpleasantness* as was used in the
137 original paper, because that word is not contained in the embeddings generated with the Goodreads dataset.

¹The trained models and files for replication are available in an anonymized format at <https://anonymous.4open.science/r/xxxxxx-3788>

138 To address claim 2, we reproduce figure 3 of the original paper. Principal component analysis is used to calculate the
 139 explained variance of the first ten components of different pairs of seed sets. For each pair of seed set, both the original
 140 order is used and a pairing where the seed pairings are shuffled. The model used is trained on the NYT dataset, but
 141 we also replicate the same figure for the BERT embeddings. We also reproduce figure 4, which ranks word vectors of
 142 different shuffled, unshuffled and random seed pairs by cosine similarity to the first principal component.

143 To address claim 3, we reproduce table 4 from the original paper by generating 300 seed set pairs using the earlier
 144 described method of generating seed sets. We calculate the coherence of each pair using the WEAT bias subspace. We
 145 also calculate the coherence for a select amount of gathered seed sets.

146 The final figure reproduced for addressing claim 3 is created by generating 600 seed pairs. The set similarity for
 147 each seed pair is calculated using the cosine similarity between the set mean vectors. For each seed pair, the explained
 148 variance of the first PCA component is also calculated. The explained variance is plotted against the set similarity, and
 149 similarly to the original paper, a trendline is calculated and plotted over the data.

150 3.6 Computational requirements

151 Preprocessing, modelling, and visualisation were performed on a personal laptop with 32 GB RAM, 8 CPUs, and 512
 152 GB SSD. Uncompressed datasets took from 45 MB to 4 GB of memory and it is not a challenge to keep them without
 153 extra computation services. Datasets were loaded into memory and processed without any batching. Word2vec model
 154 was trained using all CPU cores. Plots for BERT embeddings were obtained on Google Colab via GPU.

155 3.7 Results

156 The figures and tables of the original paper are reproduced in this section. At the hand of the figures and tables, we will
 157 go over each claim and state whether we can verify the claim or not.

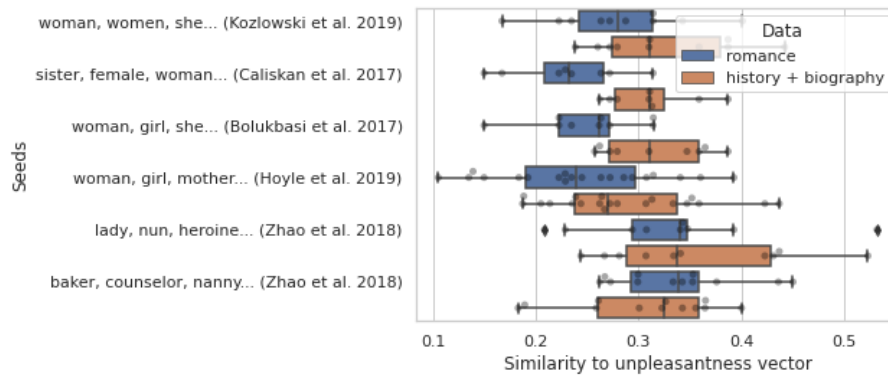


Figure 1: Cosine similarity is plotted for words in different seed sets containing *female* words and the *unpleasant* vector. The word embeddings used are created using both the romance and history+biography Goodreads dataset. The results show that different seed sets result in different bias measurements.

158 Claim 1: As shown in figure 1, different seed sets result in different bias measurement, even when all seed sets contain
 159 the same subject it tries to measure bias for. We can also see that the type of data used to generate the word embeddings
 160 can alter the bias measurement. With the replication of this image we can verify claim 1.

161 Claim 2: In table 1 the cosine similarities for ordered gender pairs identifies female and male gendered words slightly
 162 better than shuffled pairs. Figure 3 show that different pairs of seed sets have a different explained variance distribution
 163 for its principal components. While the original paper has a lower explained variance for the first principal components
 164 for the gender pairs seed set and a higher explained variance for the first principal components for the other two pairs,
 165 we have found a higher explained variance for the first principal component for all pairs of seed sets. Because the
 166 results in table 1 suggest the ordered gender pairs only slightly increase the ability to identify gendered words and the
 167 results in figure 3 does not suggest that ordered pairs can better define correlation in bias subspace. We therefore cannot
 168 verify claim 2.

herself	0.415	incentive	0.277	lily	0.36	herself	0.395
she	0.393	setback	0.248	theirs	0.34	she	0.393
female	0.354			fari	0.172	girl	0.367
her	0.348			meet	0.108	her	0.339
daughter	0.292			canoe	-0.043	daughter	0.297
boy	-0.165			bilingual	-0.126	his	-0.175
man	-0.176	likelihood	-0.106	brush	-0.202	himself	-0.191
himself	-0.176	tales	-0.162	dictates	-0.304	son	-0.207
son	-0.245	hood	-0.567	longest	-0.424	male	-0.255
he	-0.260	danced	-0.682	julianna	-0.556	he	-0.264

(a) Gender Pairs (b) Random Pairs (c) Random Pairs (d) Shuffled Gender Pairs

Table 1: Ranking word vectors by cosine similarity with respect to the PCA subspace. The top 5 and bottom 5 words ranked by cosine similarity are plotted, unless certain words are not found in the model trained on the NYT dataset, in which case fewer words are shown. (a) contains gender pairs, (b) and (c) contains random pairs and (d) contains shuffled gender pairs.



Figure 2: Using PCA, the explained variance for the first 10 principal components are plotted for 3 different pairs of seed sets. The pairs of seed sets are shuffled and the principal components for these new pairs are shown as well. The word embedding used are trained using the New York Times dataset.

169 Claim 3: For the third claim we first show how different pairs of seed sets have different coherence values. In table
 170 2, different pairs of seed sets, both generated and gathered, are shown. Seed sets with overlapping words or similar
 171 meanings are usually less coherent than more distinct pairs of seed sets. In figure 3, we can not see a significant
 172 correlation between the set similarity and the explained variance of the first principal component. We therefore cannot
 173 verify the claim made in the original paper.

174 In figure 4 we show the explained variance for the first 10 principle components using the BERT embedding. BERT
 175 shows a very high explained variance for the first principle component regardless of pairs of seed sets or whether it is
 176 shuffled or not.

177 4 Discussion

178 Overall, the main quantitative result of the authors that seed sets contain biases is affirmed by our results, and the choice
 179 on which seed set to use to measure bias is an important decision. We were not able to verify the claim made that
 180 suggest ordered pairs of seed sets identifies bias more clearly than unordered pairs of seed sets. We were also not able
 181 to verify the claim that suggests set similarity to be strongly negatively correlated to the explained variance of the
 182 first principal component. The reason we were not able to verify the results might be because of various reasons. The
 183 trained models might be different because of slightly different datasets or preprocessing implementation. For figure 3
 184 specifically, the variation might be because of a different way of generating seed sets.

185 Finally, using embeddings from modern NLP models might show bias in a different way compared to word2vec.

Coherence	Generated Seed Set A	Generated Seed Set B
1.000	years months decades weeks decade	guards guard commandos gunmen soldiers
1.000	children parents families kids mothers	front rear porch lobby perimeter
1.000	statement spokeswoman letter telephone s	market sector volatility swaps markets
...		
0.320	blip hassle softness proverbial gloss	eyes lips sunglasses hips ears
0.280	case dismissal hague juror prosecution	court courts judge judges hague
Coherence	Gathered Seed Set A	Gathered Seed Set B
0.998	ASIAN: asian asian asian asia china asia	CAUCASIAN: caucasian caucasian white america
0.997	FEMALE: sister mother aunt grandmother	MALE 2: brother father uncle grandfather
0.886	CAREER: executive management professional	FAMILY: home parents children family cousins
0.611	FEMALE: countrywoman, sororal, witches	MALE: countryman fraternal wizards manservant
0.202	NAMES ASIAN: cho wong tang huang chu chu	NAMES CHINESE: chung liu wong huang ng
0.049	NAMES BLACK: harris robinson howard	NAMES WHITE: harris nelson robinson

Table 2: Using the WEAT subspace, the coherence is calculated for each pair of seed sets. The model used was trained using the NYT dataset. For the generated seed sets, 600 pairs were calculated of which the 3 with the highest and 2 with the lowest coherence were shown. We also show the coherence for a select amount of gathered seed sets.

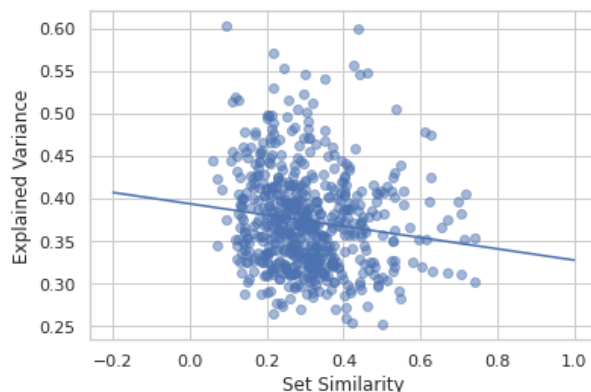


Figure 3: 600 random seed pairs are generated, after which the explained variance of the first principal component and set similarity of each pair are calculated. The WikiText dataset is used to train the word embeddings. A trendline is shown, however, there does not seem to be a significant correlation between the set similarity and explained variance.

186 4.1 What was easy

187 The paper was relatively easy to follow, and the math was quite straightforward.

188 As soon as we had determined the correct preprocessing steps after having contacted the original authors, implementing
189 the experiments from the original paper was not too challenging as the data was publicly available.

190 Another factor that made replicating the original paper easier, is that many points could be clarified by the authors,
191 who were kind enough to respond to our questions.

192 4.2 What was difficult

193 The code for the original implementation was not available. While implementing the experiments was not too difficult,
194 it was challenging to figure out the subtleties of the data preprocessing steps without having a code example, and in
195 some cases we had to resort to educated guesses. Some of these questions about implementation were cleared up after
196 contact with the authors.

197 While it was easy to find the datasets online, the references for some datasets were outdated. In the references section
198 of this work we detail the updated sources for the data that was used in this replication.



Figure 4: Using PCA, the explained variance for the first 10 principal components are plotted for 3 different pairs of seed sets. The pairs of seed sets are shuffled and the principal components for these new pairs are shown as well. The word embedding used are from a pretrained BERT word embedding with a vector size of 768.

219 According to the original paper, the WikiText-103 dataset was filtered for mathematical equations. No details
 200 were provided on how this was done, and due to the varying format of the equations found in the dataset, it was
 201 very challenging to filter them out. It is unclear whether the authors also proceeded using relatively 'dirty' data, or
 202 whether they used alternative methods to remove the mathematical equations. Table 2 from the original paper was also
 203 difficult to reproduce as the preprocessing according to authors gave close but not equal results for the NYT dataset and
 204 significantly different results for the WikiText-103 dataset. For the Goodreads dataset, the set of documents was not
 205 fixed but rather randomly sampled from a larger set using an unknown implementation, which made the comparison
 206 much more difficult.

207 In certain cases, the gathered seed sets json file contained errors. Specifically: 'daughters' was misspelt as 'daughers',
 208 and 'ma', 'am' was used as two words instead of one word "ma'am". We have notified the authors about this and the
 209 seed sets have since been updated.

210 The seed words for Figure 4 as stated in the appendix of the original paper are not the same as the ones in the image
 211 itself, which caused some confusion in replicating the figure.

212 4.3 Communication with original authors

213 Communication with the original authors played a major role in the reproducibility of the original paper. After receiving
 214 details on the pre-processing steps and clarification of some sentences in the original paper, our results came a lot closer
 215 to those of the original authors. We have included these details in this paper to be used for future research.

216 4.4 Ablation studies

217 Ablation studies were not necessary for this paper, as most of the results were quite straightforward to derive, and no
 218 multi-faceted methods were used.

219 4.5 Recommendations and future work

220 Having an original implementation available publicly is invaluable for making research reproducible. If the codebase
 221 for research is available, many uncertainties can be immediately cleared up and a major source of uncertainty for
 222 reproduction can be eliminated. Furthermore, having the code available would enable other researchers to easily adapt
 223 the code to new problems or perform new experiments. The time saved by not having to recreate an implementation
 224 could be spent on novel research, enabling the scientific community evaluate and expand on papers more quickly.

225 Besides this general recommendation, it would be interesting to reproduce this particular research with a more recent
 226 word embedding, such as those generated by Transformer-based models.

227 **References**

- 228 Kaggle new york times dataset. <https://www.kaggle.com/nzalake52/new-york-times-articles>. Accessed: 2022-04-02.
- 229 Wikitext-103. [https://www.salesforce.com/products/einstein/ai-research/the-wikitext-dependency-language-modeling-](https://www.salesforce.com/products/einstein/ai-research/the-wikitext-dependency-language-modeling-dataset/)
230 [dataset/](https://www.salesforce.com/products/einstein/ai-research/the-wikitext-dependency-language-modeling-dataset/). Accessed: 2022-04-02.
- 231 Maria Antoniak and David Mimno. 2021. Bad seeds: Evaluating lexical methods for bias measurement. In *Proceedings*
232 *of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint*
233 *Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1889–1904, Online. Association for
234 Computational Linguistics.
- 235 Tolga Bolukbasi, Kai-Wei Chang, James Zou, Venkatesh Saligrama, and Adam Kalai. 2016. Man is to computer
236 programmer as woman is to homemaker? debiasing word embeddings.
- 237 Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language
238 corpora contain human-like biases. *Science*, 356(6334):183–186.
- 239 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: pre-training of deep bidirectional
240 transformers for language understanding. *CoRR*, abs/1810.04805.
- 241 Nikhil Garg, Londa Schiebinger, Dan Jurafsky, and James Zou. 2018. Word embeddings quantify 100 years of gender
242 and ethnic stereotypes. *Proceedings of the National Academy of Sciences*, 115(16):E3635–E3644.
- 243 Alexander Miserlis Hoyle, Lawrence Wolf-Sonkin, Hanna Wallach, Isabelle Augenstein, and Ryan Cotterell. 2019.
244 Unsupervised discovery of gendered language through latent-variable modeling. In *Proceedings of the 57th Annual*
245 *Meeting of the Association for Computational Linguistics*, pages 1706–1716, Florence, Italy. Association for
246 Computational Linguistics.
- 247 Austin C Kozlowski, Matt Taddy, and James A Evans. 2019. The geometry of culture: Analyzing the meanings of class
248 through word embeddings. *American Sociological Review*, 84(5):905–949.
- 249 Thomas Manzini, Lim Yao Chong, Alan W Black, and Yulia Tsvetkov. 2019. Black is to criminal as caucasian is to
250 police: Detecting and removing multiclass bias in word embeddings. In *Proceedings of the 2019 Conference of the*
251 *North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume*
252 *1 (Long and Short Papers)*, pages 615–621, Minneapolis, Minnesota. Association for Computational Linguistics.
- 253 Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. 2018. Learning gender-neutral word embeddings.
254 In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4847–4853,
255 Brussels, Belgium. Association for Computational Linguistics.