# Exploring and Mitigating Gender Bias in GloVe Word Embeddings

**Ma Francesca Luisa C Vera**
Department of Computer Science
Stanford University
*fvera@stanford.edu*

## Abstract

When societal biases are discovered within items, it is natural to consider ways in which it is possible to remove those biases. In the past, language has been shown to carry certain biases including those that perpetuate gender and racial stereotypes. Consequently, much research has been done on how to better counteract these biases within language. In conjunction with this research is natural language processing, and the growing use of technology to solve linguistic tasks has led to many considerations about the biases algorithms, models, and tools may carry. One example of such a tool is word embeddings, which give words corresponding numerical values. Within these embeddings are gender biases against occupations that ought to be gender-neutral, but are often stereotyped towards male or female genders. We demonstrate that within the GloVe word embeddings, these occupations are stereotyped since similarities of these occupations to embeddings for "he" and "she" produce clearly different results. The correlation coefficient for "masculine" occupations is 0.91 against the 0.83 coefficient for "feminine" occupations. To debias these words, we propose finding a gender subspace in which gender-neutral words are placed at the 0 position – while still maintaining the embeddings of gender-specific words such as "man" and "woman". Because the new gender-neutral words are given embeddings such that they are equidistant between corresponding pairs of gender-specific words, it is expected that when looking at the stereotyped occupations, the embeddings of these occupations will be similar to males and females when tested. When we test for correlation coefficients against "masculine" and "feminine" occupations, we get 0.99 and 0.97 respectively, demonstrating how gender bias has been mitigated.

## 1    Introduction

In recent years, there has been increasing concern regarding the degrees of bias that might exist in technological tools used to perform tasks for societal problems. It is important to understand these biases because these tools may then perpetuate even further bias in the results of a task. One area that has drawn the attention of academics in recent years is gender bias in word embeddings. Word embeddings include a set of vectors that correspond to a specific word, and these vectors relate to each other in the same ways the words do. Thus, word embeddings are powerful because they can be used to solve natural language processing tasks that may require this transformation between words and numerical values. However, recent work has called into question whether or not these word embeddings promote societal biases including those related to race and gender. If this is the case, one critical responsibility of the users of these word embeddings is to ensure that these biases are mitigated before the embeddings are put to use.

There are many different sets of pretrained word embeddings available for public use, including the GloVe pretrained word embeddings. The GloVe word embeddings include sets that were trained on billions of tokens, some up to 840 billion tokens. It is available for download online, making it a popular source for word embeddings in the NLP space. When looking at such an influential tool, any semblance of bias can influence the results of a task

50  gravely.

51  This paper aims to highlight some of the gender biases that exist in the GloVe word
52  embeddings before adapting a past debiasing method to mitigate the biases in these
53  embeddings. It also experiments with the concept of gender-specific words using GloVe
54  word embeddings, exploring whether or not it is possible to classify a large set of
55  embeddings as gender-specific or not.

56
## 57  2    Related Work

58  In 2016, Bolukbasi et al. released a paper, *Man is to Computer Programmer as Woman is to*
59  *Homemaker? Debiasing Word Embeddings*, which pioneered this space of debiasing word
60  embeddings.[1] Their work involved looking at the lists of occupations that most clearly
61  exhibited gender stereotypes and looking at gender stereotyped he : she analogies. In this
62  context, bias was defined as showing what should be gender neutral occupations or analogies
63  are favored in the direction of one gender over another. Using w2vNEWS embeddings, they
64  first identify a gender subspace before "neutralizing" and "equalizing". (Further explained in
65  Section 4) For the purposes of this paper, this is the method I will be adapting and applying
66  to the GloVe word embeddings.

67  Another important piece of work by Chakraborty et al., *Reducing gender bias in word*
68  *embeddings*, uses GloVe vectors to .[2] Similarly to Bolukbasi et al., Chakraborty et al. look
69  at modifying the actual embeddings so that they do not exhibit gender biases. To do so, they
70  alter the settings in which the embeddings are trained, taking the co-occurrence matrix of an
71  occupation and scaling these so that the co-occurrence probabilities of these occupations
72  becomes 1. After using cosine similarity between embeddings and a gender direction, they
73  also aim to reduce bias by adding a regularization term to the objective function that
74  "penalizes" similarity to the gender direction, doing so for only the biased occupational
75  words.

76  Recently, there has been growing interest in the applications of adversarial learning. In 2018,
77  Zhang et al. released a paper, *Mitigating Unwanted Biases with Adversarial Learning*, which
78  discusses fairness and attempts to create "more fair" algorithms in different societal
79  settings.[3] One such setting was natural language processing, particularly gender bias in
80  word embeddings. Using embeddings trained from Wikipedia to generate input data, Zhang
81  et al. set up the "he : she" analogy test as a supervised learning task where the model would
82  pick the word corresponding to "she" after being given a word analogous to "he." To debias
83  the system, they added an "adversarial discriminator network" that would have trouble
84  guessing the gender direction of an output y. Thus, the embeddings remain unchanged but
85  the space they exist in differs such that they do not strongly perpetuate gender stereotypes
86  from the analogy task.

87
## 88  3    Demonstrating Bias in GloVe Word Embeddings

89  The first step to mitigating bias in anything is to show that there currently exists some form
90  of bias in the thing one is trying to debias. In this case, we aim to show that the GloVe word
91  embeddings demonstrate an inherent gender bias.

92  Borrowing a method from Bolukbasi et al (mentioned in the related work section), we will
93  show that when looking at occupations that are supposed to be gender-neutral, there are
94  inherent gender biases built into the word embeddings of these occupations. To show these
95  biases, we will map out the similarities these occupations have to the embeddings of "he"
96  and "she" – if an embedding appears more similar to "he" than "she", it would appear that
97  that occupation tends towards males – and vice versa. The occupations have been split into
98  "he" occupations and "she" occupations, with the understanding that these are the "most
99  biased" occupations in either a he or she direction. To calculate similarity, we use both the
100 inner product of the embeddings and the cosine similarity. Finally, to show that gender
101 biases exist across different types of embeddings, we use both the GloVe word embeddings
102 pretrained on Wiki and pretrained on Common Crawl.

103 Plotting all the occupations and their respective similarities to "he" and "she" gives us a

104 qualitative overview on how biased these words are. To quantify the results, we also look at
105 the correlation coefficients for the "he" occupations against the "she" occupations, with
106 coefficients that are nearly equal indicating a lack of bias.

107
### 108 3.1 Using occupations from Bolukbasi et al.

109 It makes sense that the results rely heavily on which predetermined "he" and "she"
110 occupations are used, so we decided to experiment using different groups. The first set of
111 occupations was taken from the Bolukbasi et al paper.[1] To determine these occupations,
112 they found the most extreme occupations as projected onto a gender direction, and labeled
113 these as "Occupational Stereotypes." The lists of occupations are as follows:

114 "He" Occupations: *["maestro", "skipper", "protege", "philosopher", "captain", "architect",*
115 *"financier", "warrior", "broadcaster", "magician", "pilot", "boss"]*

116 "She" Occupations: *["homemaker", "nurse", "receptionist", "librarian", "socialite",*
117 *"hairdresser", "nanny", "bookkeeper", "stylist", "housekeeper", "designer", "counselor"]*

118 These occupations were determined by Bolukbasi et al. as the ones that carried the greatest degree
119 of gender bias within them. To verify, we turn to the GloVe word embeddings and our method of
120 determining gender bias. When using both cosine and inner product to calculate similarity, the
121 gender bias is very clear. Below (Table 1) are the results for the correlation coefficients after
122 running our experiments on the Wiki-trained GloVe word embeddings. Because the coefficients
123 for "he" occupations and "she" occupations are not nearly equal for either type of similarity, there
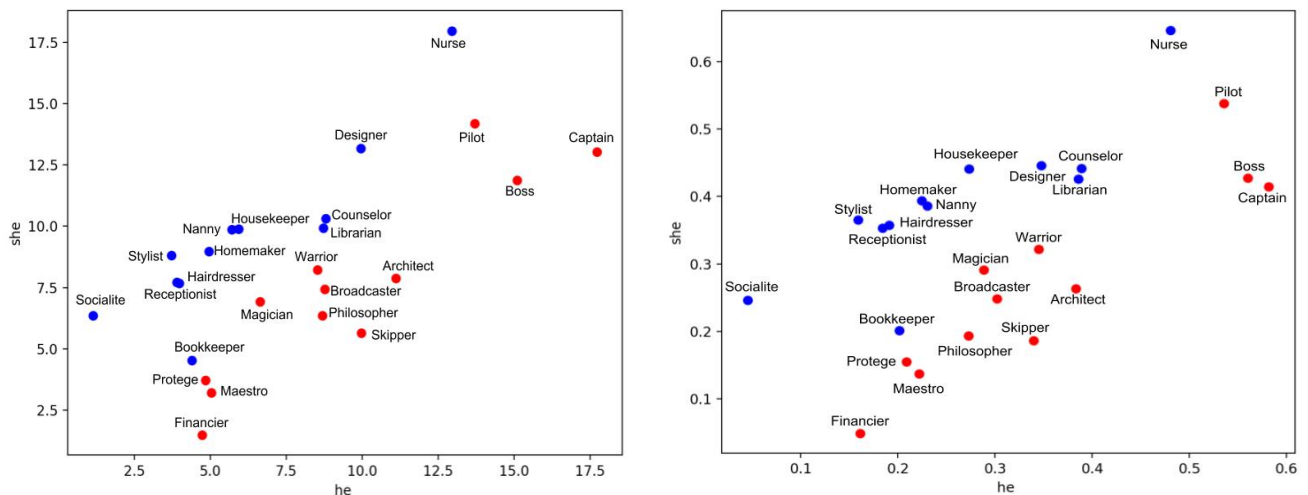124 is clearly some form of gender bias within the embeddings.

125 Table 1: Correlation coefficients using Wiki-trained Glove word embeddings

| Correlation Coefficients | "He" Occupations | "She" Occupations |
|---|---|---|
| Cosine Similarity | 0.9147 | 0.8322 |
| Inner Product Similarity | 0.9116 | 0.8820 |

126 An easier way to visualize gender bias is by plotting the occupation's similarities to "he" and
127 "she" against each other. This can be seen in figures 1-4, which show the results for both the
128 Wiki-trained GloVe word embeddings and the Common Crawl-trained GloVe word embeddings.
129 Each red point represents a "he" occupation and each blue point represents a "she" occupation –
130 with the charts showing which occupation corresponds to which point. For all charts, there is a
131 very obvious split between occupations (blue and red points), which corresponds to the calculated
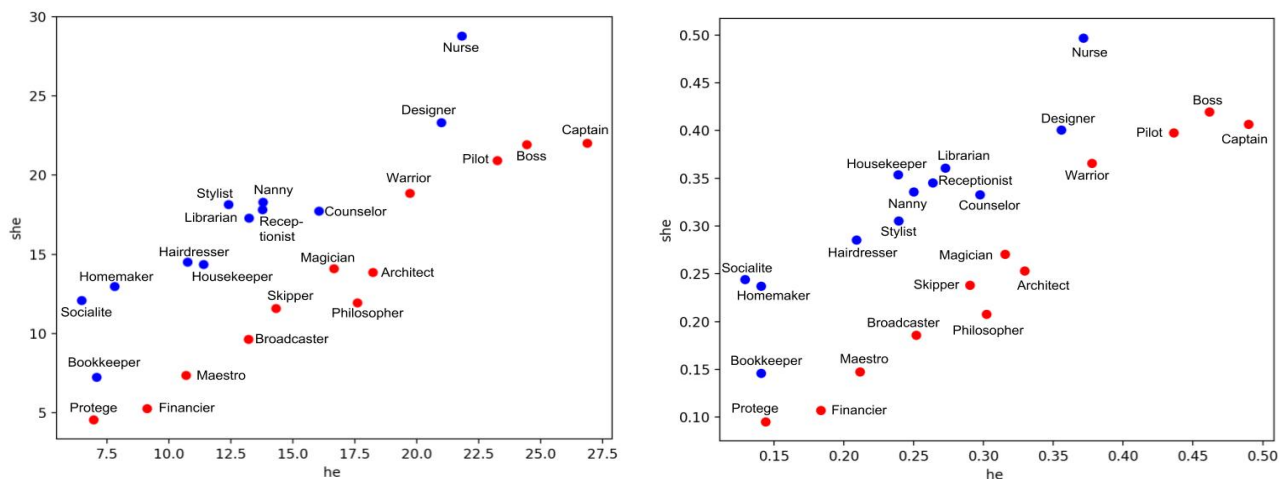132 coefficients.

133 Figure 1(left): Inner product similarities using Wiki-trained word embeddings

134 Figure 2 (right): Inner product similarities using Common Crawl-trained word embeddings

135    Figure 3 (left): Cosine similarities using Wiki-trained word embeddings

136    Figure 4 (right): Cosine similarities using Common Crawl-trained word embeddings



137

## 3.2    Using "most similar" occupations

139    We also looked at the occupations that had the closest similarities to "he" and "she" word
140    vectors, and then plotted those occupations to view whether or not they carried as clear of a
141    bias as the occupations from Bolukbasi et al. did. To find these occupations, we first
142    downloaded a list of 1155 general occupations. For each occupation on the list, we
143    calculated its similarity to "he" word embedding and "she" word embedding. After sorting
144    the list, we found the twelve closest occupations to each embedding – deeming those "he"
145    and "she" occupations. The results are as follows. In bold are words that also appeared in
146    Bolukbasi et al.'s list. Underlined are words that appear in both "he" and "she" lists.
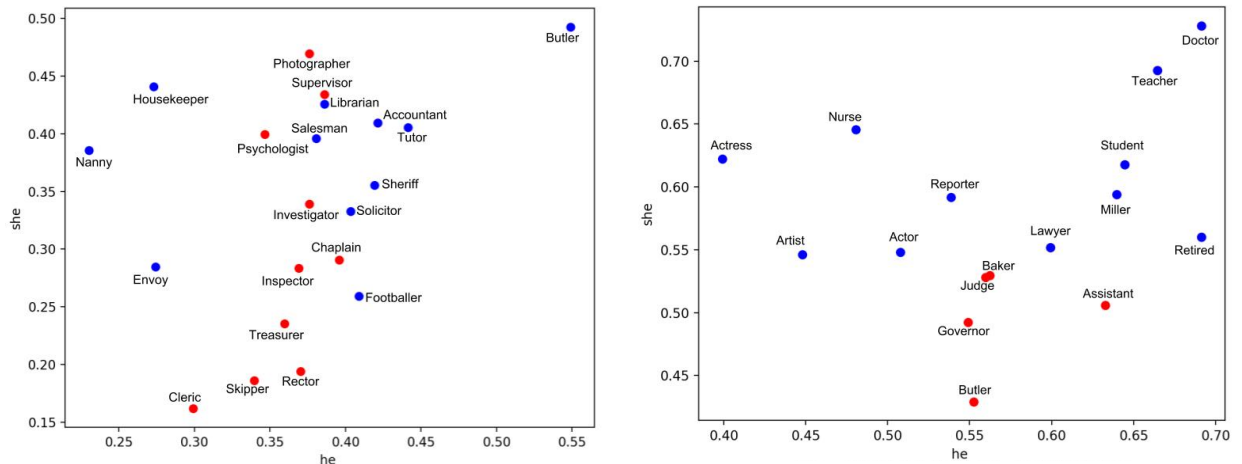
147    Table 2: List of occupations "most similar" to "he" and "she"

|  | "He" Occupations | "She" Occupations |
|---|---|---|
| Cosine Similarity | *["retired", "doctor", "teacher", "student", "miller", "assistant", "lawyer", "baker", "judge", "governor", "butler"]* | *["doctor", "**teacher**", "**nurse**", "actress", "student", "miller", "reporter", "retired", "lawyer", "actor", "artist"]* |
| Inner Product Similarity | *["cleric", "photographer", "**skipper**", "chaplain", "accountant", "inspector", "rector", "investigator", "psychologist", "treasurer", "supervisor"]* | *["**librarian**", "**housekeeper**", "**nanny**", "accountant", "sheriff", "envoy", "tutor", "salesman", "butler", "footballer", "solicitor"]* |

148    The results are very interesting because there is a lot of overlap between the two lists and
149    only few similarities with the Bolukbasi et al. list. The method also seemed to better at
150    identifying biased "she" occupations than "he" occupations – although some anomalies like
151    "actor" appear. There are several reasons for this outcome. Firstly, the results were
152    influenced by the initial list of 1155 occupations; a different list of initial occupations for
153    comparison would have yielded different results. Furthermore, Bolukbasi et al. added an
154    extra element of a "gender direction" in their calculations whereas we only use cosine or
155    inner product similarity.[1] That being said, it is encouraging that when plotted using the
156    same method to measure bias as above, there is still some observed gender bias among
157    occupations – although not as clear as before.

158    Figure 5 (left): Cosine similarities of occupations found through inner product

159    Figure 6 (right): Cosine similarities of occupations found through cosine



160

161    Consequently, we decided to continue using the Bolukbasi et al. occupations for our
162    evaluation of debiasing.

163

## 4    Debiasing GloVe Word Embeddings

165    After demonstrating that the GloVe word embeddings carry some gender bias within them,
166    especially for "gendered" occupations, we then tried to debias these word embeddings so
167    that when looking at these occupations again, no clear bias would be reflected. As mentioned
168    in the related work section, there have already been a few academic studies regarding the
169    debiasing of word embeddings. For continuity, we adopt the method in Bolukbasi et al.'s
170    paper since our evaluation was based on findings from that paper.

171

### 4.1    Approach for mitigating bias

173    The method from Bolukbasi et al. was based on two important steps:

174        1.   Identifying the "gender subspace"
175        2.   Neutralizing or Equalizing words appropriately to form a new set of embeddings

176    The first step of identifying the gender subspace involves identifying a direction in which
177    the embedding carries some gender bias. To calculate this gender subspace, we first identify
178    a gender direction by looking at "definitional" pairs that help it orient towards the genders.
179    For those words that are not gender specific, the new embedding is the difference between
180    the former embedding and gender direction, multiplied by the embedding dot the gender
181    direction, divided by the gender direction dot itself. For words that are gender specific, their
182    embeddings are maintained.

183    Next, we decide whether or not a word should be equalized or neutralized. If a word is
184    gender neutral, we neutralize it such that they are at position 0 in the gender subspace.
185    Equalize then looks at pairs of corresponding gender-specific words so that we can enforce
186    any gender-neutral word to be equidistant from these corresponding pairs. From there, we
187    collect all the altered embeddings to create a new set of word embeddings with limited bias.

188

### 4.2    Results of debiasing GloVe word embeddings

190    To evaluate our method of debiasing, we will use the same methods as we did for showing
191    that the initial GloVe word embeddings carried some form of gender bias. This means that

192 when we use the new debiased embeddings in our methods of evaluation, we should not see
193 a clear divide between "he" and "she" occupations. This applies to plotting the
194 corresponding similarities and having the same correlation coefficients.

195 Table 3: Correlation coefficients using debiased word embeddings
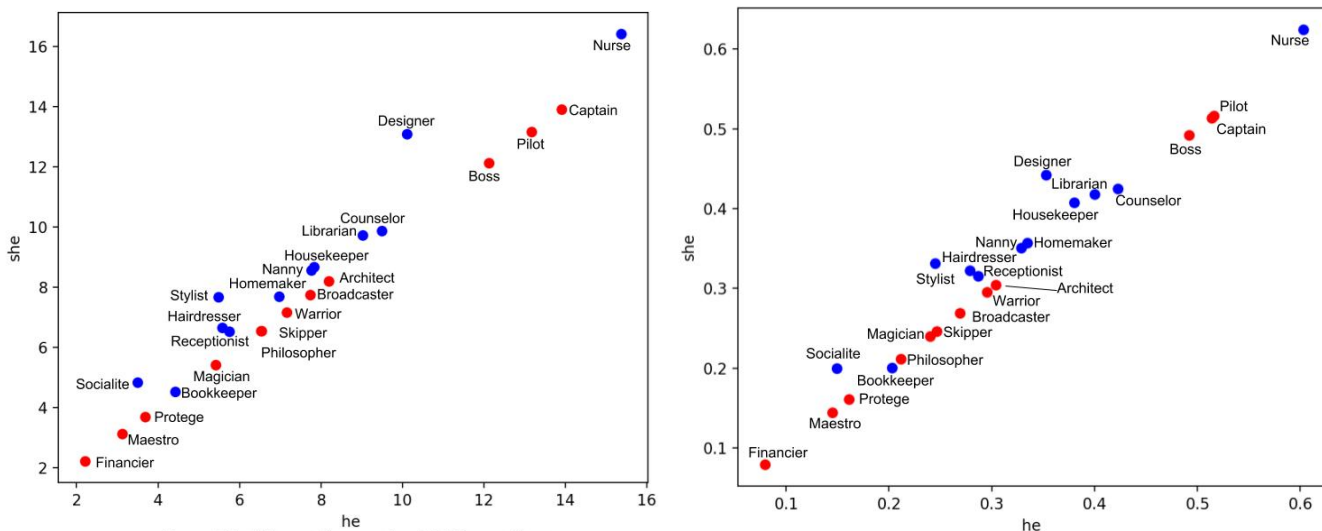
| Correlation Coefficients | "He" Occupations | "She" Occupations |
|---|---|---|
| Cosine Similarity | 0.9890 | 0.9688 |
| Inner Product Similarity | 0.9921 | 0.9723 |

196 We observe in table 3 that the correlation coefficients for "he" occupations and "she"
197 occupations when using the new debiased embeddings are nearly the same. This also holds
198 true for when we use either cosine similarity or inner product similarity.

199 Now that we have verified numerically that the initial gender bias has been reduced, we turn
200 to the qualitative methods of evaluation to see whether or not the differences between "he"
201 and "she" occupations have changed visually. Figures 7 and 8 show the points for all
202 gendered occupations when the similarities have been calculated with the new word
203 embeddings. As we can see, unlike the previous plots where "she" occupations tended
204 towards the "she" vector and "he" occupations tended towards the "he" vector, the ratio
205 between similarities to "he" and "she" are nearly 1:1 for all occupations – putting what were
206 once gendered occupations in the middle of the "he" and "she" genders.

207 Figure 7: Inner product similarities using newly debiased embeddings

208 Figure 8: Cosine product similarities using newly debiased embeddings



209

210 It is encouraging that both quantitative and qualitative methods of evaluation reflect some
211 change in the word embeddings. This means that we were successful in our task of
212 mitigating the gender bias in the GloVe wiki-trained word embeddings and can use the new
213 debiased embeddings in other NLP tasks.

214

215 **5      Extra Experiments (on Gender-Specific Words)**

216 In addition to the task of debiasing a set of word embeddings, we were inspired to explore
217 the concept of gender-specific words in a similar way to what was done in the Bolukbasi et
218 al. paper. Some words such as "he", "she", "mom", "dad" etc. are specific to gender and so
219 do not carry gender biases. There are two interesting tasks that can be done regarding
220 gender-specific words:

221    1. By using a pre-labeled set of gender-specific/non-gender-specific words and their
222       embeddings, can we classify a broader list of words as either gender-specific/non-
223       gender-specific?
224    2. Can we also use the existing set of gender-specific words to find additional gender-
225       specific words from a larger list of general words?

226 The list of gender-specific words were taken from Bolukbasi et al. [1], which took a subset
227 of 218 words from w2vNEWS and looked at their Wordnet definitions to determine if there
228 was some element of gender inherent to the definition. Some of the words included: *['he',
229 'his', 'her', 'she', 'him', 'man', 'women', 'men', 'woman', 'spokesman', 'wife', 'himself', 'son', 'mother',*
230 *'father', 'chairman', 'daughter', ..., 'fatherhood', 'councilwoman', 'princes', 'matriarch', 'colts', 'ma',*
231 *'fraternities', 'pa', 'fellas', 'councilmen', 'dowry', 'barbershop', 'fraternal', 'ballerina']*. As
232 Bolukbasi states, the words are highly "subjective" and encourages customization to the
233 application.[1] One can imagine that in the application of occupations, words such as
234 "councilwoman" will be more important to include than "colts". For the purposes of this paper, we
235 will use this set of words to perform the two tasks at hand.

236
237 **5.1    Classifying gender-specific words**

238 The first task we experimented on was the classification of words as either gender-specific
239 or non-gender-specific. To do this, we trained a linear SVC using $C = 1.0$ on a subset of
240 GloVe word embeddings where any gender-specific word in the subset was labeled "1", and
241 the rest were considered gender-neutral so labeled "0". According to Bolukbasi et al.'s
242 results on the same task, the binary accuracy is expected to be "well over 99% due to the
243 imbalanced nature of the classes." In the results below, we find that we are able to achieve
244 the same score on different numbers of iterations. Another metric used by Bolukbasi et al.
245 was the F-score from 10-fold cross-validation. They achieved 0.627. Our results vary much
246 more in this metric. Although we are able to beat the score when there are only 5000 words
247 used in the training subset, it is expected that this is nowhere close to the 50 000 used by
248 Bolukbasi et al. Thus, we significantly underperform in this regard when the number of
249 words used to train the model increases.

250 Table 4: Results from classification of gender-specific words task

| Number of Words in Subset | Accuracy | 10-fold Cross-Validation Score |
|---|---|---|
| 5000 | 0.9996 | 0.7998 |
| 7500 | 0.9994 | 0.6998 |
| 10 000 | 0.9994 | 0.6597 |
| 20 000 | 0.9994 | 0.4998 |
| 100 000 | 0.9994 | 0.4998 |

251
252 **5.2    Identifying gender-specific words**

253 The second task involves finding even more gender-specific words from an initial set of
254 gender-specific "seed" words. To do so, we first train a linear SVC on a subset of GloVe
255 word embeddings and labels where any gender-specific word that appears is labeled "1". In
256 the previous subsection, we demonstrate that setting up a classifier of this sort achieves high
257 accuracy and reasonable cross-validation scores at certain numbers of iterations. After
258 training on this model, we found the model's coefficient and intercept. Subsequently, for
259 each word embedding in the larger list of general words, the dot product of that embedding
260 was taken with the model's coefficient and then compared to the intercept, resulting in a list
261 of words taken from the list supposedly with some gender-specificity.

262 This method yielded surprising results. There is some evidence of success as our model was

able to extract the following gender-specific words: *["macho", "dude", "gentleman", "dads", "guys", "gunman", "man", "mommy", "guy", "woman", "spokesman"]* among a few others from the model. However, some gender-neutral words such as "kid" and "somebody" were added to the list. An interesting observation is that the model also pulled some occupations as gender-specific, including "politician" and "cop" – hinting at the occupational gender biases that were discussed earlier in the paper.

## 6     Conclusion

It is important to be conscious of the inherent biases that might exist in technology because using these tools will perpetuate the inherent biases they hold in whatever tasks the tools are trying to perform. Word embeddings are no exception to this concept, and this work demonstrates how the GloVe word embeddings, a popular set of word embeddings, also carry biases with regards to gender stereotypes.

The ability to understand bias, determine where it exists, and then mitigate it is important in better understanding the biases that crop up in our world. In the context of NLP, our methods successfully create a new set of word embeddings that have limited gender bias. There is potential for further research using the foundation we built by testing the techniques in this paper against some standard NLP tasks such as the analogy task. The analogy task is useful in determining whether the embeddings encourage gender stereotypes within analogies.

Note that our work in debiasing changes the word embeddings directly. However, there has been work in debiasing that alters the gender bias space rather than altering the embeddings themselves. Consequently, it would be worthwhile to explore if this can be done using GloVe word embeddings. Although this method does not leave users with a new set of debiased embeddings to use, it demonstrates that gender biases have been recognized on many fronts and there exist successful attempts to mitigate these biases in the technological world. Overall, there is a lot of promising work out there that aims to mitigate societal biases within technology today.

### References

[1] Bolukbasi, T.; Chang, K.-W.; Zou, J. Y.; Saligrama, V.; and Kalai, A. T. Man is to computer programmer as woman is to homemaker? Debiasing word embeddings. In *Advances in Neural Information Processing Systems*, 4349– 4357. 2016.

[2] Chakraborty, T.; Badie, G.; Rudder, B.. Reducing gender bias in word embeddings. 2016.

[3] Zhang, B. H.; Lemoine, B.; Mitchell, M. Mitigating Unwanted Biases with Adversarial Learning. In *arXiv preprint arXiv:1801.07593*. 2018.

[4] Pennington, J.; Socher R.; Manning C. D.. GloVe: Global Vectors for Word Representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532-1543. 2014.