

Does Double-Hard Debias Keep its Promises? - A Replication Study

Kristina Kobrock¹, Meike Korsten¹, & Sonja Börgerding¹

¹ University of Osnabrück

Does Double-Hard Debias Keep its Promises? - A Replication Study

Introduction

Recent research has shown that word embeddings derived from natural language corpora inherit human biases. The first seminal study on this topic was by Bolukbasi, Chang, Zou, Saligrama, and Kalai (2016) who used the word analogy task developed by Mikolov, Yih, and Zweig (2013) to prove that “Man is to Computer Programmer as Woman is to Homemaker” - a clearly gender biased analogical relation derived from a word2vec embedding trained on Google News. Caliskan, Bryson, and Narayanan (2017) complement this finding with a large study on human-like semantic biases in Natural Language Processing (NLP) tools. They have shown that human biases as exhibited in psychological studies using, for example, the Implicit Association Test (IAT) (Greenwald, McGhee, & Schwartz, 1998) are learned by NLP algorithms designed to construct meaningful word representations. According to Zhao, Wang, Yatskar, Ordonez, and Chang (2018) these biases propagate to downstream tasks. As pre-trained word embeddings are frequently fed into NLP architectures used for more complex tasks encountered in everyday life like Machine Translation, the biased embeddings bear the risk of proliferating and strengthening existing biases and stereotypes in human cultures.

Debiasing Embeddings. But how can the problem of biased embeddings be solved? Several researchers have proposed post-processing techniques or algorithm modifications that promise to “debias” the word embeddings obtained by algorithms like word2vec or GloVe (e.g. Bolukbasi, Chang, Zou, Saligrama, & Kalai, 2016; Kaneko & Bollegala, 2019; Zhao, Zhou, Li, Wang, & Chang, 2018). Bolukbasi, Chang, Zou, Saligrama, and Kalai (2016) have developed an algorithm called Hard Debias that is based on the idea of removing biased gender directions from the embedding whilst preserving desired gender relations. For example, the encoded gender of the words “king” and “queen” is given by definition. So a relation between the concept *male* and “king” and between

female and “queen” is desired. On the other hand, words like “nurse” and “doctor” also tend to exhibit relations to a specific gender in semantic space, even though the words do not have a gender-specific meaning but can be used for all genders. The proposed Hard Debias algorithm first identifies a gender subspace of the embedding that best captures the bias. Subsequent neutralizing and equalizing steps ensure that the debiased embeddings satisfy neutrality constraints without impairing desired properties of the semantic space. Wang et al. (2020) build on Hard Debias in their newly proposed Double-Hard Debias algorithm. The main idea is to not only remove the gender direction of the biased embedding, but also the frequency direction as it has been shown that word frequency has a significant impact on the geometry of the embedding space (e.g. Gong et al., 2018; Mu, Bhat, & Viswanath, 2018, more on this later).

Motivation. In this project for the course “Implementing ANNs with Tensorflow,” we aimed to replicate the debiasing method presented by Wang et al. (2020) and to reproduce their experimental results. We chose this paper because it is the most recent one that proposes a post-processing technique for debiasing algorithms and because it builds on the seminal paper by Bolukbasi, Chang, Zou, Saligrama, and Kalai (2016). In the following, we would like to quickly motivate our choice of topic by relating it to the course and stating a motivation for replication and reproduction attempts in general. The course covered a huge breadth of topics ranging from the simplest neural network architectures to advanced Convolutional Neural Networks (CNN) (e.g. He, Zhang, Ren, & Sun, 2015; Huang, Liu, & Weinberger, 2016) and recently proposed Transformer models (Vaswani et al., 2017). One of the topics covered that stroke us the most interesting was Deep Learning for Natural Language Processing (NLP) covering word embeddings, language models, as well as Seq2seq models and preprocessing techniques. So our aim was to find a final project that would extend on the NLP knowledge gained in the course. The extra session on “Ethical aspects of ML” held by Pascal and Annie was a great starting point and we decided to settle on bias in word embeddings as a topic. Our search for papers proposing

debiasing algorithms led us to the work described above and we decided to reimplement the paper by Wang et al. (2020). This gave us the opportunity to apply a variety of skills learned in the course ranging from preprocessing datasets for NLP tasks, over skillfully working with word embeddings, to understanding, fine-tuning and training a NN model that we did not know before (GloVe (Jeffrey Pennington & Manning, 2014)). The reproduction of existing results is not only good practice in science, but it is also essential for gaining a deeper understanding on the methods used and it can help to validate existing results or to shed well-grounded doubt on them where needed. Recent studies of reproducibility in the field of Computer Science (e.g. Collberg, Proebsting, & Warren, 2015) and NLP (Belz, Agarwal, Shimorina, & Reiter, 2021; Cohen et al., 2018; e.g. Fokkens et al., 2013; Mieskes, 2017) explain why reproducibility endeavours are often failing and shed light on the exact nature of the ‘reproducibility crisis’ (term coined by e.g. Baker, 2015) in NLP research. Belz, Agarwal, Shimorina, and Reiter (2021), for example, report that the community’s interest in topics of reproducibility have risen even though reproduction attempts still tend to fail due to problems like missing data, missing code and incomplete documentation (see also Fokkens et al., 2013; Mieskes, Fort, Névél, Grouin, & Cohen, 2019). This project aims to contribute to the increasing body of reproducibility and replication attempts in NLP research.

Implementation

Shortly after starting the replication attempt, we noticed that the code provided on the Github repository linked in the paper by Wang et al. (2020) was poorly documented and seemed to deviate from the paper’s suggestions in some aspects. This is why we stucked closely to the ideas developed in the paper when reimplementing the proposed algorithm. In the following, we will explain in detail what we did and how specific implementation choices were motivated.

Pilot Study: Impact of Word Frequency. Double-Hard Debias is built on the claim that an embedding’s encoding of word frequency significantly influences the same embedding’s encoding of gender which can lead to a diminished efficacy of any debiasing algorithm (Wang et al., 2020). As interest in debiasing embeddings increased, so did awareness of the debiasing method’s weaknesses. Gonen and Goldberg (2019), for example, noticed that the leading methods barely scrape the surface of removing the full bias. Wang et al. (2020) now propose that the performance of debiasing methods, specifically of post-processing methods such as Hard Debias proposed by Bolukbasi, Chang, Zou, Saligrama, and Kalai (2016), can be significantly improved by mitigating frequency features before applying the debiasing step. To ground this assumption, Wang et al. (2020) perform a short pilot study. They manipulate word frequencies in the original data to investigate the impact of frequency on gender bias contained in the embedding. This is done by sampling certain sentences twice containing specific words of the set of definitional pairs introduced by Bolukbasi, Chang, Zou, Saligrama, and Kalai (2016), i.e. words with a definitional gender whose difference vectors are assumed to approximate the gender direction (Wang et al., 2020). Word embeddings were trained on a GloVe (Jeffrey Pennington & Manning, 2014) model implementation that was reworked by us (original from GradySimon). The model is trained using a minibatch size of 512, a learning rate of 0.05 and the Adam optimizer over 50 epochs. The embeddings created are 300-dimensional as we work with embeddings of that size in the remaining paper as well. These hyperparameter settings depicted a good trade-off between computational needs and performance on our machines. We could not replicate the authors’ approach as neatly as other parts of the paper due to hardware constraints which led us to working with a significantly smaller corpus. This left us with similar, but definitely less pronounced, results. Despite limited resources we were able to reproduce the effect that manually changing the word frequency statistics for a single word significantly influences the resulting gender direction of the embedding space. This supports the conclusion of Wang et al.

(2020) that enhanced emphasis needs to be put on word frequency in debiasing techniques.

Datasets and Preliminaries. The pre-trained, 300-dimensional GloVe embedding used by Wang et al. (2020) is provided via link on the authors’ Github repository. The original files included embeddings for 322.636 words. For easier access to embeddings throughout coding we created vocabularies and dictionaries mapping words to embedding indices (IDs). During pre-processing we followed common steps such as the removal of words containing digits or special characters and a restriction of the vocabularies to the 50.000 most common words which corresponds to the first 50.000 words in a GloVe embedding.

We made use of a number of word sets provided by Wang et al. (2020). Among the files used were `definitional_pairs`, `equalize_pairs`, `gender_specific_full`, `male_word_file` and `female_word_file`. The two files `definitional_pairs` and `equalize_pairs` are directly based on the work of Bolukbasi, Chang, Zou, Saligrama, and Kalai (2016), however for the remaining files it is unclear where the word sets were obtained. The gender pairs needed for the determination of the gender direction during the Hard Debias method (Bolukbasi, Chang, Zou, Saligrama, & Kalai, 2016), i.e. words which are associated with a gender by definition like “girl, boy” or “she, he,” are provided in `definitional_pairs`. The file `equalize_pairs` includes what Bolukbasi, Chang, Zou, Saligrama, and Kalai (2016) called the equality sets (word pairs such as “Women” and “Men” or “princess” and “prince”), which are needed for the equalization step. As an alternative to this, we created another set of pairs from the `female_word_file` and the `male_word_file` provided on the Github repository of Wang et al. (2020) since they also include corresponding words such as “priest” and “nun” or “spokesman” and “spokeswoman.” Lastly the Hard Debias algorithm expects to be passed a set of neutral words as well as sets of female and male words (Bolukbasi, Chang, Zou, Saligrama, & Kalai, 2016). We will later describe how the two gendered sets are obtained using the gender direction. To obtain the neutral set, we removed all words found in the five files

provided by Wang et al. (2020) from our vocabulary, since these files include words with a definitional gender.

Hard Debias. The authors make use of the Hard Debias algorithm proposed by Bolukbasi, Chang, Zou, Saligrama, and Kalai (2016). The paper at hand does not give much information on the exact implementation of Hard Debias used, which is why we stucked to the original paper from Bolukbasi, Chang, Zou, Saligrama, and Kalai (2016) in order to reimplement Hard Debias.

The paper describes two steps: First, the gender direction (or, more generally, the subspace) has to be identified. This is achieved with the help of defining sets $D_1, D_2, \dots, D_n \subset W$ which consist of gender specific words. These are the words that can help to identify the gender direction by capturing the concept *female*, *male* in the embedding $\{\vec{w} \in \mathbb{R}^d\}_{w \in W}$. Whereas in some simple implementations of Hard Debias, only one definitional pair might be used, Bolukbasi, Chang, Zou, Saligrama, and Kalai (2016) suggest to compute the gender direction B across multiple pairs to more robustly estimate the bias. In our implementation we used the 10 word pairs suggested by Bolukbasi, Chang, Zou, Saligrama, and Kalai (2016) which were experimentally shown to agree with an intuitive concept of gender (Bolukbasi, Chang, Zou, Saligrama, & Kalai, 2016). For these 10 pairs the principal components (PCs) are calculated and the bias subspace B is made of the first $k \geq 1$ rows of the decomposition $\text{SVD}(C)$. According to Bolukbasi, Chang, Zou, Saligrama, and Kalai (2016), the first eigenvalue is significantly larger than the rest and so the top PC is hypothesized to capture the gender subspace. So $k = 1$ is chosen and our resulting gender subspace B is thus simply a direction. C is calculated in the following way:

$$C := \sum_{i=1}^n \sum_{w \in D_i} (\vec{w} - \mu_i)^T (\vec{w} - \mu_i) / |D_i|$$

where $\mu_i := \sum_{w \in D_i} \vec{w} / |D_i|$ are the means of the defining sets $D_1, D_2, \dots, D_n \subset W$.

As a second step Hard Debias neutralizes and equalizes the word embeddings. Neutralizing means to transform each word embedding \vec{w} such that every word $w \in N$ has

zero projection in the gender subspace. So for each word $w \in N$ in a set of neutral words $N \subseteq W$, we re-embed \vec{w} :

$$\vec{w} := (\vec{w} - \vec{w}_B) / \|\vec{w} - \vec{w}_B\|$$

. The equalize step ensures that desired analogical properties hold for both female and male words contained in the equality sets $\mathcal{E} = \{E_1, E_2, \dots, E_m\}$ where each $E_i \subseteq W$. For example, after debiasing we would like the embeddings of the pair $E = \{\textit{grandmother}, \textit{grandfather}\}$ contained in the equality sets to be equidistant from the embedding of “babysit.” This is enforced by equating each set of words $E \in \mathcal{E}$ outside of B to their simple average $\nu := \mu - \mu_B$ where $\mu := \sum_{w \in E} w / |E|$ before adjusting vectors such that they are of unit length. So for each word $w \in E$, \vec{w} is re-embedded to

$$\vec{w} := \nu + \sqrt{1 - \|\nu\|^2} \frac{\vec{w}_B - \mu_B}{\|\vec{w}_B - \mu_B\|}$$

. With the help of the original Bolukbasi, Chang, Zou, Saligrama, and Kalai (2016) paper, we were successfully able to reimplement Hard Debias.

Double-Hard Debias. Our implementation of Double-Hard Debias was guided by the pseudocode presented in the paper (Wang et al., 2020). The main idea of the algorithm is to remove two different components from the embedding. The first one supposedly encodes frequency information, the second one captures the gender direction, as proposed already by Bolukbasi, Chang, Zou, Saligrama, and Kalai (2016). In earlier work, Mu, Bhat, and Viswanath (2018) have shown that reducing the dimensionality of word embeddings can increase the linguistic regularities they capture. Specifically, they first subtract the common mean vector $\mu = \frac{1}{|V|} \sum_{\vec{w} \in V} \vec{w}$, with V being the set of word embeddings \vec{w} , from all embeddings. Then, a few top principal components of the embedding space are identified and removed. The resulting word embeddings actually prove to serve as stronger linguistic representations than the unprocessed ones. Mu, Bhat, and Viswanath (2018) noticed that some of the top PCA directions encode frequency to a significant degree. Wang et al. (2020) base their debiasing technique on those findings and apply the

described steps in their algorithm as a method that identifies the directions of the embedding space which encode word frequency.

The first part of Double-Hard Debias identifies the frequency direction that is eventually ought to be removed from all embeddings. This is done on a subset of 500 female and 500 male most biased words by examining different subspace projections and how these influence Hard Debias’ results. Following Mu, Bhat, and Viswanath (2018), the common mean vector is removed before PCA is applied. In the pseudocode presented, Wang et al. (2020) calculate as many principal components (PCs) as denoted by the dimensionality of the embedding. In the provided code, on the other hand, they only compute the top 20 principal components. As the components do not encode much variance of the data after around 20 PCs, we also decided to compute only the first 20. For each of these PCs, the corresponding direction \mathbf{u} is removed $w' = w - (\mathbf{u}^T w)\mathbf{u}$, Hard Debias is applied, and performance of the resulting embeddings on the Neighborhood Metric (Gonen & Goldberg, 2019) is measured. This metric will be discussed in a more detailed manner in the evaluation part. Then, the PC that leads to the best performance on the Neighborhood Metric, i.e. results in least biased embeddings, is identified. This direction, which is hypothesized to encode frequency and to significantly affect gender information, is needed for the second part of the algorithm. The second part focuses on debiasing the full set of word embeddings. The chosen frequency direction is removed from all embeddings before Hard Debias is applied. By first creating the frequency-mitigated projections, removing the gender direction shows to be more successful.

Evaluation

This part deals with the reproduction of the experimental results by Wang et al. (2020) when evaluating their Double Hard debiased embedding compared to some baselines and on benchmark datasets using well-established tasks. We could reuse more of the code written by Wang et al. (2020) and provided on the authors’ Github repository for the

evaluations than for the implementation.

Baselines. Following Wang et al. (2020) we used GloVe embeddings that were obtained using several different debiasing methods, some applied during training and some during post-processing, as baselines for our evaluation along with the original, non-debiased embedding and the Double-Hard debiased embedding obtained by Wang et al. (2020). The embeddings were linked on the Github repository of Wang et al. (2020).

original GloVe: The non-debiased GloVe embedding used and provided by Wang et al. (2020), trained on the 2017 January dump of English Wikipedia.

Double-Hard-GloVe(Wang et al.): Double-Hard debiased GloVe embedding obtained by Wang et al. (2020) and reported to be the result of their implementation of Double-Hard Debias.

Double-Hard-GloVe(replication): Double-Hard debiased GloVe embedding obtained in this project. This is the result of applying our implementation of Double-Hard Debias to the original GloVe embedding.

GN-GloVe: Gender-Neutral GloVe embedding released by Zhao, Zhou, Li, Wang, and Chang (2018). This method restricts gender information in certain dimensions while neutralizing in the remaining dimensions.

GN-GloVe(a): A variant of Gender-Neutral GloVe obtained by Wang et al. (2020). Created by excluding the gender dimensions from the GN-GloVe embedding in an attempt to completely remove gender.

GP-GloVe: Gender preserving GloVe embedding released by Kaneko and Bollegala (2019). This method attempts to remove stereotypical gender bias and preserve non-discriminative gender information.

GP-GN-GloVe: Gender preserving, Gender-Neutral GloVe embedding provided by Kaneko and Bollegala (2019). This is the result of applying gender preserving debiasing to

an already debiased GN-GloVe embedding.

Hard-GloVe: Hard debiased GloVe embedding obtained by Wang et al. (2020) following the implementation of Bolukbasi, Chang, Zou, Saligrama, and Kalai (2016). This method aims to debias neutral words while preserving the gender specific words.

Strong-Hard-GloVe: Strong-Hard debiased GloVe embedding obtained by Wang et al. (2020). A variant of Hard debias which debiases all words instead of only neutral words.

Evaluation of Debiasing Performance. Wang et al. (2020) evaluated their embeddings on different tasks to test multiple of the characteristics that embeddings inherit. Even though we were unable to replicate all those tests, we focused on covering the main aspects.

Debiasing in Downstream Applications. As shortly mentioned above, a main concern regarding bias in word embeddings regards the application of pre-trained word embeddings in downstream tasks. The following method was developed to measure in how far bias contained in a pre-trained embedding influences downstream task performance.

Coreference Resolution. Wang et al. (2020) evaluated the performance of their Double-Hard debiased embedding in comparison to the baselines in coreference systems. Bias in coreference models has been shown by Zhao, Wang, Yatskar, Ordonez, and Chang (2018) as models which more successfully identify the semantic reference in “The physician hired the secretary because **he** was overwhelmed with clients” due to the biased association of “the physician” with the male gender, allowing it to be more readily related to the concept *he*. In contrast, non-consistent sentences such as “The physician hired the secretary because **she** was overwhelmed with clients” showed poorer performance. The idea is that the societal bias of “the physician” is implemented within its embedding, thereby influencing the model’s performance. Unbiased embeddings, on the other hand, should lead to equal model performance with both consistent and non-consistent samples. Wang et al. (2020) themselves neither provide any code for evaluations on the coreference

task, nor relate to any external source of the results they present in their paper. Actually, the Coreference Resolution is not mentioned at all in the accompanying material. Trying to implement the original Zhao, Wang, Yatskar, Ordonez, and Chang (2018) coreference task also showed to be not feasible due to its huge computational demands, which ultimately lead to us leaving it out.

Debiasing at Embedding Level. Besides evaluating the performance on different downstream NLP tasks, the authors also specifically investigated how much bias can be determined in the embeddings.

The Word Embeddings Association Test (WEAT). Following along the evaluation of Wang et al. (2020) we tested the bias of all embeddings with a permutation test. This Word Embedding Association Test (Caliskan, Bryson, & Narayanan, 2017) takes four sets of words, two sets of so-called target words and two sets of attribute words. It calculates the relative similarity between the target words and the attribute words respectively and outputs how likely it is that this result could have been obtained from a non-biased distribution. The two values it returns are the effect size d and the p-value p . A value of $p < 0.05$ indicates a significant bias and the bias is considered more pronounced for larger effect sizes. We used a pre-implemented version of WEAT for our evaluation that can be obtained here. Instead of reusing the code of Wang et al. (2020) we opted for a different WEAT implementation that is oriented on the original paper by Caliskan, Bryson, and Narayanan (2017). We decided not to use the authors’ original code since it was partly incomplete and poorly understandable. Instead we decided for a more compact implementation where d and p are obtained simultaneously. Wang et al. (2020) conducted WEAT three times for each embedding, once with the bias *Career & Family*, once with *Math & Arts* and lastly with *Science & Arts*. We replicated all three of these with the word lists taken directly from Caliskan, Bryson, and Narayanan (2017). Following Wang et al. (2020), we made one alteration to the sets, namely exchanging “Bill” in the male names of the *Career & Family* set for “Tom.” This is to avoid ambiguity due to the lower-casing of

GloVe. The exact lists can be found in our implementation of the evaluation. The results we obtained can be inspected in Table 1.

Table 1

WEAT Test

Embeddings	C&F d	C&F p	M&A d	M&A p	S&A d	S&A p
original GloVe	1.806	0.0001	0.6886	0.0837	1.1299	0.0116
Double-Hard-GloVe (Wang et al.)	1.531	0.0011	-0.5625	0.8702	-0.6503	0.9033
Double-Hard-GloVe (replication)	1.530	0.0012	-0.6895	0.9160	-0.9419	0.9707
GN-Glove	1.821	0.0001	-0.2564	0.6960	1.0690	0.0161
GN-Glove(a)	1.756	0.0002	0.5030	0.1586	0.8797	0.0398
GP-GloVe	1.806	0.0002	1.2085	0.0077	1.1064	0.0135
GP-GN-GloVe	1.797	0.0002	-0.0126	0.5102	0.8460	0.0447
Hard-GloVe	1.547	0.0010	-0.9826	0.9756	-0.5384	0.8586
Strong-Hard-GloVe	1.547	0.0010	-0.9857	0.9761	-0.5471	0.8620

With the WEAT implementation we used we were able to almost perfectly reproduce the results of Wang et al. (2020) for the *Career & Family* word sets. Both d and p differ only slightly from the reported values and we can clearly see that for the Double-Hard debiased embeddings p is larger not only in comparison to the original embedding but also to the other debiasing methods. Despite the bias still being significant ($p < 0.05$), the large p and small d compared to the baselines suggest a more effective debiasing. We were also able to obtain comparable results for the Double-Hard debiased embedding provided by Wang et al. (2020) and for our self-debiased embedding. For both *Math & Arts* and *Science & Arts* we obtained values for d and p that differ noticeably from the values reported by the authors. However, as in Wang et al. (2020), we can see that the bias in *Math & Arts* was already insignificant in the original GloVe embedding.

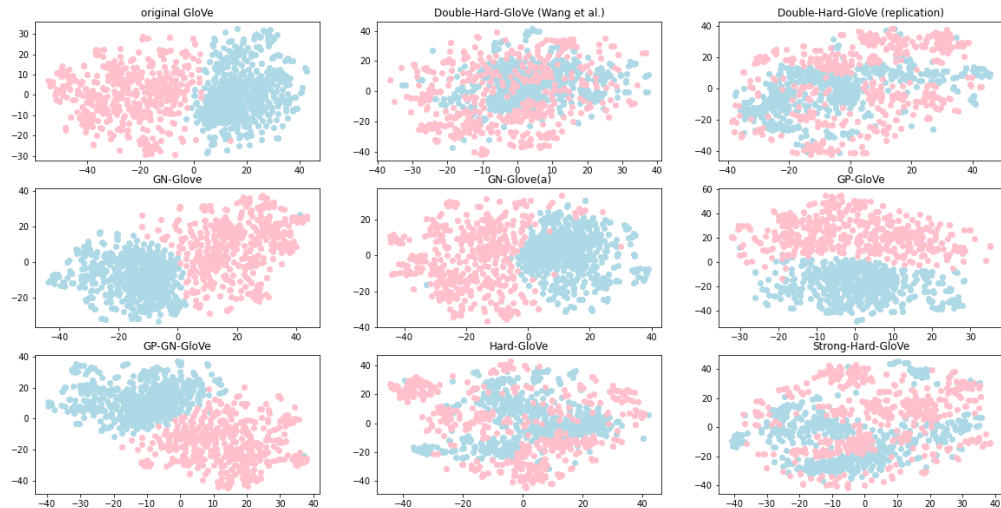
Neighborhood Metric. The Neighborhood Metric was originally introduced by Gonen and Goldberg (2019) and is based on the observation that even though debiased words no longer inherit their bias in form of a specific gender direction, it remains in the word embeddings’ spatial grouping. The supposedly “removed” bias is still manifested as one can observe those words socially-marked with the same gender situated closer to one another in the embedding space (Gonen & Goldberg, 2019). K-means clustering on a number of most biased words can often straightforwardly cluster “debiased” embeddings into the correct female and male categories, uncovering the remaining bias. Wang et al. (2020) evaluate K-means clustering performance on such most biased words by simply counting the samples correctly assigned to the gender bias allegedly removed in debiasing. The alignment score a is defined as $a = \frac{1}{2k} \sum_{i=1}^{2k} 1[\hat{g}_i == g_i]$ and set to $a = \max(a, 1 - a)$ with k resembling the number of samples for each gender, \hat{g}_i the estimated and g_i the correct label. According to this definition, an alignment score value of 0.5 indicates perfectly unbiased embeddings, as the clustering algorithm failed to replicate the original gender pattern (Wang et al., 2020). Even though Wang et al. (2020) make use of this metric within their debiasing algorithm, adjusting their embeddings to optimize performance on the Neighborhood Metric, they also apply it again in the scope of evaluating their debiasing performance in comparison to other embeddings. As one would expect, the Double-Hard debiased embeddings therefore do particularly well in comparison to the other baseline embeddings. Our results can be observed in Table 2, and are visualized with the help of tSNE below. Wang et al. (2020) applied the tSNE reduction technique to the top 500 female and male biased words to be able to plot the resulting clusters in two-dimensional space. Both values in Table 2 and the tSNE visualization show that Hard-Glove, Strong-Hard-Glove and both Double-Hard-GloVe embeddings obtained by us and by Wang et al. (2020) lead to nearly bias-free clustering performances, whereby our Double-Hard debiased embeddings slightly outperform the authors’ ones. The remaining baseline embeddings on the other hand can nearly perfectly be clustered,

Table 2

Neighborhood Metric Clustering Accuracy

Embeddings	Top 100	Top 500	Top 1000
original GloVe	100.00	100.00	100.00
Double-Hard-GloVe (Wang et al.)	54.50	56.50	58.40
Double-Hard-GloVe (replication)	53.00	50.60	53.70
GN-Glove	100.00	99.90	99.90
GN-Glove(a)	100.00	99.70	97.90
GP-GloVe	100.00	100.00	100.00
GP-GN-GloVe	100.00	99.80	99.80
Hard-GloVe	51.50	51.50	50.15
Strong-Hard-GloVe	50.50	51.60	50.90

indicating remaining bias. Where our results deviate considerably to the ones from the paper is in the performance of (Strong-)Hard-Glove. In our evaluation these embeddings also show clustering accuracies alluding to basically no remaining bias, while Wang et al. (2020) present results that would indicate significantly more remaining bias.



In the graphic above you can see a plot that shows all the baseline datasets, as well as our Double-Hard debiased embedding and Wang et al. (2020)’s embedding. For the original GloVe embedding, GN-GloVe, GP-GloVe and GP-GN-GloVe the clustering accuracy is still very high, meaning that the top 500 female (lightblue) and male (pink) biased words are still biased after applying the respective debiasing techniques.

Analysis of Retaining Word Semantics. One of the most important properties of embeddings is that they represent meaningful word semantics, for example in the form of analogies or concepts. In this section it is tested in how far the embeddings still meet this criterion after debiasing.

Word Analogy. The word analogy task was introduced by Mikolov, Yih, and Zweig (2013). The task is to find the word D such that “A is to B as C is to D.” One example for an unbiased analogy is: “Man is to King as Woman is to Queen” whereas a biased analogy would be: “Man is to Computer Programmer as Woman is to Homemaker” (Bolukbasi, Chang, Zou, Saligrama, & Kalai, 2016). The debiased embeddings are evaluated on two word analogy test sets, the MSR (Mikolov, Yih, & Zweig, 2013) and the Google word analogy task (Mikolov, Chen, Corrado, & Dean, 2013), in order to find out

Table 3

Analogy Tasks

Embeddings	MSR	Sem	Syn	Total
original GloVe	54.40	80.14	55.69	63.90
Double-Hard-GloVe (Wang et al.)	42.40	71.36	47.99	56.67
Double-Hard-GloVe (replication)	62.12	80.20	66.55	71.13
GN-Glove	44.10	69.28	51.01	57.80
GN-Glove(a)	43.26	69.00	50.64	57.46
GP-GloVe	42.36	73.53	51.55	59.72
GP-GN-GloVe	44.65	69.44	53.18	59.22
Hard-GloVe	54.27	73.62	62.74	66.79
Strong-Hard-GloVe	53.94	71.23	62.43	65.70

whether they preserve desired unbiased analogies. The MSR word analogy dataset contains 8000 syntactic questions in the form presented above. The missing word D is computed by maximizing the cosine similarity between D and $C - A + B$. The evaluation metric is the percentage of correctly answered questions (see Wang et al., 2020). The Google word analogy dataset contains 19,544 (**Total**) questions, 8,869 of which are semantic (**Sem**) and 10,675 are syntactic (**Syn**) questions.

The results can be inspected in Table 3 and are in line with the results reported by Wang et al. (2020). What is interesting is that our replicated Double-Hard debiased embedding scores best in all word analogy tasks and even outperforms the original GloVe embedding.

Concept Categorization. For Concept Categorization a set of words is provided where each word is assigned to one of k classes. After clustering by spacial closeness in the embedding space, it is evaluated how many of the words of the same class are assigned to

the same cluster. A high accuracy suggests that the words relating to the same concept have similar embeddings and this is a desirable trait of an embedding space. Wang et al. (2020) applied this evaluation method to all the different embedding baselines and their Double-Hard debiased embedding as well. However, we found the code provided by the authors to be incomplete. It was missing major parts, thereby hindering us from adequately replicating its application and due to time constraints we decided not to implement this method.

Discussion

The code provided by Wang et al. (2020) was incomplete at several points which hindered both our implementation of the Double-Hard Debias method and the Evaluation. For example the authors' notebook `GloVe_Debias` only showcases the application of Double-Hard Debias but does not provide the actual embedding that was used by Wang et al. (2020) in the evaluation as a result. This can be seen as several pre-processing steps are applied in `GloVe_Debias` such as the restriction of the vocabulary to 50.000 words or the removal of words containing digits. These steps go unmentioned in Wang et al. (2020)'s documentation and were not applied in the missing code that actually resulted in the Double-Hard debiased embedding as this is not restricted and of the same length as the original embedding. We decided to adopt some of the pre-processing steps oriented on their github repository for computational reasons. We also noted that from the code provided by Wang et al. (2020) it was not apparent whether they used the word pairs from `equalize_pairs` or those from the combination of `male_word_file` and `female_word_file` for the equalization step of the Hard debias method. We were able to report that this does not change the results, however it is another example of ill documentation. At many points throughout the evaluations, most noticeably the analysis of retaining word semantics, we were confronted with the `web` module used by Wang et al. (2020). However the authors did not provide this module, neither did they link or even

specify where the datasets and functions included could be found. This resulted in the need of self-implementing some functions as well as funding the datasets they used and pre-processing them such that they were usable with Wang et al. (2020)’s functions. Due to time constraints we were not able to cover everything included in this module which lead to us not implementing the Concept Categorization for our evaluation. Further, Wang et al. (2020) supplied no code of their evaluation using the coreference system. On Github the authors provided several links and suggestions how to replicate this, however no details were provided and the recreation of this proved time and resource intensive. Therefore we were not able to recreate this part of the evaluation either. We obtained results for three parts of the evaluation, namely WEAT, the Neighborhood Metric and the Word Analogy Task. For the Neighborhood Metric we obtained results that closely resemble those of Wang et al. (2020). Similarly we can report results comparable to those of the authors for the Word Analogy Task. However, after applying WEAT we obtained values that differ from those of Wang et al. (2020) for two of three word sets. The closely matching values we obtained for *Caree & Family* suggest a good replication of the authors’ results. The differences for the other word sets do not necessarily change the interpretation but are still concerning.

Preprocessing is one of the main aspects not typically reported in scientific NLP papers that can cause experimental variation identified by Fokkens et al. (2013). In line with this, the paper by Wang et al. (2020) did not specify how the relevant datasets were preprocessed. This could be one of the reasons why our results deviate from the results obtained by the original authors. While we were able to reproduce the main conclusions and findings of the paper, we were not able to reproduce concrete values. As Mieskes, Fort, Névél, Grouin, and Cohen (2019) point out: “Experiments in replication fail more often than not.” Whereas the replication attempt described in this report was not a complete failure, many difficulties were identified during the process of replication due to incomplete documentation of the original study and lack of provided code. The main reasons for this

appear to be common in the field: incomplete documentation, missing files and non-working links (Collberg, Proebsting, & Warren, 2015; see Mieskes, 2017).

Conclusion (kris)

References

- Baker, M. (2015). Reproducibility crisis: Blame it on the antibodies. *Nature*, *521*(7552), 274–276. <https://doi.org/https://doi.org/10.1038/521274a>
- Belz, A., Agarwal, S., Shimorina, A., & Reiter, E. (2021). *A systematic review of reproducibility research in natural language processing*. Retrieved from <http://arxiv.org/abs/2103.07929>
- Bolukbasi, T., Chang, K.-W., Zou, J. Y., Saligrama, V., & Kalai, A. (2016). Man is to computer programmer as woman is to homemaker? Debiasing word embeddings. *CoRR*, *abs/1607.06520*. Retrieved from <http://arxiv.org/abs/1607.06520>
- Caliskan, A., Bryson, J. J., & Narayanan, A. (2017). Semantics derived automatically from language corpora contain human-like biases. *Science*, *356*(6334), 183–186. <https://doi.org/10.1126/science.aal4230>
- Cohen, K. B., Xia, J., Zweigenbaum, P., Callahan, T., Hargraves, O., Goss, F., . . . Hunter, L. E. (2018). Three dimensions of reproducibility in natural language processing. *Proceedings of the eleventh international conference on language resources and evaluation (LREC 2018)*. Miyazaki, Japan: European Language Resources Association (ELRA). Retrieved from <https://www.aclweb.org/anthology/L18-1025>
- Collberg, C., Proebsting, T., & Warren, A. M. (2015). *Repeatability and benefaction in computer systems research*. studie. Retrieved from <http://reproducibility.cs.arizona.edu/v2/RepeatabilityTR.pdf>
- Fokkens, A., Erp, M. van, Postma, M., Pedersen, T., Vossen, P., & Freire, N. (2013). Offspring from reproduction problems: What replication failure teaches us. *Proceedings of the 51st annual meeting of the association for computational*

- linguistics (volume 1: Long papers)*, 1691–1701. Sofia, Bulgaria: Association for Computational Linguistics. Retrieved from <https://www.aclweb.org/anthology/P13-1166>
- Gonen, H., & Goldberg, Y. (2019). *Lipstick on a pig: Debiasing methods cover up systematic gender biases in word embeddings but do not remove them*. Retrieved from <http://arxiv.org/abs/1903.03862>
- Gong, C., He, D., Tan, X., Qin, T., Wang, L., & Liu, T.-Y. (2018). FRAGE: Frequency-agnostic word representation. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, & R. Garnett (Eds.), *Advances in neural information processing systems* (Vol. 31, p.). Curran Associates, Inc. Retrieved from <https://proceedings.neurips.cc/paper/2018/file/e555ebe0ce426f7f9b2bef0706315e0c-Paper.pdf>
- Greenwald, A. G., McGhee, D. E., & Schwartz, J. L. K. (1998). Measuring individual differences in implicit cognition: The implicit association test. *Journal of Personality and Social Psychology*, 74(6), 1464–1480. <https://doi.org/10.1037/0022-3514.74.6.1464>
- He, K., Zhang, X., Ren, S., & Sun, J. (2015). Deep residual learning for image recognition. *CoRR*, *abs/1512.03385*. Retrieved from <http://arxiv.org/abs/1512.03385>
- Huang, G., Liu, Z., & Weinberger, K. Q. (2016). Densely connected convolutional networks. *CoRR*, *abs/1608.06993*. Retrieved from <http://arxiv.org/abs/1608.06993>
- Jeffrey Pennington, Richar Socher, & Manning, C. (2014). Glove: Global vectors for word representation. *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, 1532–1543. Dohar, Qatar: Association for Computational Linguistics. Retrieved from

<https://nlp.stanford.edu/pubs/glove.pdf>

- Kaneko, M., & Bollegala, D. (2019). Gender-preserving debiasing for pre-trained word embeddings. *Proceedings of the 57th annual meeting of the association for computational linguistics*, 1641–1650. Florence, Italy: Association for Computational Linguistics. <https://doi.org/10.18653/v1/P19-1160>
- Mieskes, M. (2017). A quantitative study of data in the NLP community. *Proceedings of the first ACL workshop on ethics in natural language processing*, 23–29. Valencia, Spain: Association for Computational Linguistics. <https://doi.org/10.18653/v1/W17-1603>
- Mieskes, M., Fort, K., Névél, A., Grouin, C., & Cohen, K. B. (2019). NLP Community Perspectives on Replicability. *Recent Advances in Natural Language Processing*. Varna, Bulgaria. Retrieved from <https://hal.archives-ouvertes.fr/hal-02282794>
- Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). *Efficient estimation of word representations in vector space*. Retrieved from <http://arxiv.org/abs/1301.3781>
- Mikolov, T., Yih, W., & Zweig, G. (2013). Linguistic regularities in continuous space word representations. *Proceedings of the 2013 conference of the north American chapter of the association for computational linguistics: Human language technologies*, 746–751. Atlanta, Georgia: Association for Computational Linguistics. Retrieved from <https://www.aclweb.org/anthology/N13-1090>
- Mu, J., Bhat, S., & Viswanath, P. (2018). *All-but-the-top: Simple and effective postprocessing for word representations*. Retrieved from <http://arxiv.org/abs/1702.01417>
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., . . . Polosukhin, I. (2017). *Attention is all you need*. Retrieved from <http://arxiv.org/abs/1706.03762>

Wang, T., Lin, X. V., Rajani, N. F., McCann, B., Ordonez, V., & Xiong, C. (2020).

Double-hard debias: Tailoring word embeddings for gender bias mitigation.

Association for computational linguistics (ACL).

Zhao, J., Wang, T., Yatskar, M., Ordonez, V., & Chang, K.-W. (2018, January).

Gender bias in coreference resolution: Evaluation and debiasing methods. 15–20.

<https://doi.org/10.18653/v1/N18-2003>

Zhao, J., Zhou, Y., Li, Z., Wang, W., & Chang, K.-W. (2018). Learning

gender-neutral word embeddings. *Proceedings of the 2018 conference on*

empirical methods in natural language processing, 4847–4853. Brussels, Belgium:

Association for Computational Linguistics.

<https://doi.org/10.18653/v1/D18-1521>