

# Taking uncertainty in word embedding bias estimation seriously: a bayesian approach

Alicja Dobrzeniecka and Rafal Urbaniak

## 1 Cosine-based measures of bias

### 1.1 Word embeddings and their bias

- akapit filozoficzny, z odniesieniem do value free

Modern Natural Language Processing (NLP) models are used to complete various tasks such as providing email filters, smart assistants, search results, language translations, text analytics and so on. All of them need as an input words represented by means of numbers which is accomplished with word embeddings, in which particular lexical units are represented as vectors of real numbers.

- troche o tym jak sa skonstruowane (optimizing for co-occurrence frequency)
- przejrzec jak to jest wprowadzane w kilku innych artykulach, zwl. fair is better than sensational
- troche wiecej o tej konstrukcji i intuicji semantycznych, ale to nie super dobry pomysl
- undesirable bias with respect to a certain groups words if stereotypically connected words (for stereotypes that we don't want to be used or relied on in downstream tasks) are systematically closer to each other. + EXAMPLES

It has been suggested [1–6] that in the learning process such models can learn implicit biases that reflect harmful stereotypical thinking. A large chunk of the literature on the topic focuses on the geometry of the learnt word embeddings; in particular, unusually low (cosine) proximity of words belonging, intuitively, to a stereotype, is taken as a sign that this stereotype has been built into a given embedding. This methodology will be in the focus of our paper.

### 1.2 General challenges

### 1.3 WEAT and MAC

One of the first measures in the discussion has been developed by [1]. First, the gender direction  $gd$  is obtained by taking the differences of the vectors corresponding to ten different gendered pairs (such as  $\vec{she} - \vec{he}$  or  $\vec{girl} - \vec{boy}$ ) and then identifying their principal component (which is the vector obtained by projecting the data points on their linear combination in a way that maximizes the variance of the projections).<sup>1</sup> The gender bias of a word  $w$  is understood as its projection on the gender direction:  $\vec{w} \cdot gd$ . Given the supposedly gender neutral words  $N$ <sup>2</sup> and the gender direction  $gd$  the direct gender bias is defined as the average cosine similarity of the words in  $N$  from  $gd$  ( $c$  is a parameter determining how strict we want to be):

$$\text{directBias}_c(N, gd) = \frac{\sum_{w \in N} |\cos(\vec{w}, gd)|^c}{|N|} \quad (1)$$

The use of projections has been criticized for instance by [4], who point out that while the distance to the gender direction might be an indicator of bias, it is only one possible manifestation of it, and reducing the cosine distance to such a projection might be insufficient. For instance, “math” and “delicate” might

<sup>1</sup>In the notebook associated with the paper, the authors simply use  $\vec{she} - \vec{he}$  as a gender direction, though.

<sup>2</sup>We follow the methodology in assuming that there is a class of words that ideally should be identified as neutral, such as *ballpark, solution, lecture, science, book* can be identified. We will have a bit to say about this assumption when we describe our dataset construction.

get back to this!

be in equal distance to a pair of opposed explicitly gendered words, while being closer to quite different stereotypical attribute words. Further, it is observed in [4] that most word pairs preserve similarity under debiasing meant to minimize projection-based bias.<sup>3</sup>

A measure of bias in word embeddings which does not employ gender directions, the Word Embedding Association Test (WEAT), has been proposed in [2]. The idea is that the bias between two sets of target words,  $X$  and  $Y$  (we call them protected words), should be quantified in terms of the cosine similarity between the protected words and attribute words coming from two sets of stereotype attribute words,  $A$  and  $B$  (we'll call them attributes). For instance,  $X$  might be a set of male names,  $Y$  a set of female names,  $A$  might contain stereotypically male-related, and  $B$  stereotypically female-related career words. The association difference for a term  $t$  is:

$$s(t, A, B) = \frac{\sum_{a \in A} \cos(t, a)}{|A|} - \frac{\sum_{b \in B} \cos(t, b)}{|B|} \quad (2)$$

then, the association difference between  $A$  and  $B$  is:

$$s(X, Y, A, B) = \sum_{x \in X} s(x, A, B) - \sum_{y \in Y} s(y, A, B) \quad (3)$$

The effect size is computed by normalizing the difference in means as follows:<sup>4,5</sup>

$$\text{weat}(A, B) = \frac{\mu(\{s(x, A, B)\}_{x \in X}) - \mu(\{s(y, A, B)\}_{y \in Y})}{\sigma(\{s(w, A, B)\}_{w \in X \cup Y})} \quad (4)$$

WEAT, however, has been developed to investigate biases corresponding to a pair of supposedly opposing stereotypes, and so the question arises as to how generalize the measure to contexts in which biases with respect to more than two stereotypical groups are to be measured. Such a generalization can be found in [6]. The authors introduce Mean Average Cosine similarity as a measure of bias (strictly speaking, in the paper they report cosine distances<sup>6</sup> rather than similarities). Let  $T = \{t_1, \dots, t_k\}$  be a class of protected word embeddings, and let each  $A_j \in A$  be a set of attributes stereotypically associated with a protected word). Then:

$$s(t_i, A_j) = \frac{1}{|A_j|} \sum_{a \in A_j} \cos(t_i, a) \quad (5)$$

$$\text{mac}(T, A) = \frac{1}{|T| |A|} \sum_{t_i \in T} \sum_{A_j \in A} s(t_i, A_j) \quad (6)$$

That is, for each protected word  $T$  and each attribute class, they first take the mean for this protected word and all attributes in a given attribute class, and then take the mean of thus obtained means for all the protected words.

Having introduced the measures, first, we will introduce a selection of general problems with this approach, and then we will move on to more specific but important problems related to the statistical significance of such measurement. We will focus on WEAT and MAC, as we want to put issues with the use of projections aside. WEAT will be useful in our criticism of the statistical methods involved, as it is simpler, so the explanation and visualizations will be more transparent (and *mutatis mutandis* this criticism will apply to MAC as well), and MAC will be useful in the development of our alternative method as its range of applicability is the widest.

<sup>3</sup>In [1] another method which involves analogies and their evaluations by human users on Mechanical Turk is also used. We do not discuss this method in this paper, see its criticism in [7].

<sup>4</sup>WEAT is a modification of the Implicit Association Test (IAT) [8] used in psychology, and it uses almost the same word sets, allowing for a *prima facie* sensible comparison with bias in humans. [2] argue that significant biases—thus measured—similar to the ones discovered by IAT can be discovered in word embeddings. [5] extended the methodology to a multilingual and cross-lingual setting, arguing that using Euclidean distance instead of similarity does not make much difference, while the bias effects vary greatly across embedding models (interestingly, with social media-text trained embeddings being less biased than those based on Wikipedia).

<sup>5</sup>A similar methodology is employed in [3]. The authors employ word embeddings trained on corpora from different decades to study the shifts in various biases. For instance, to compute the occupational embeddings bias for women the authors first compute the average vector of vector embeddings of words that represent women (e.g. *she*, and *female*), then calculate the Euclidean distance between this mean vector and words for occupations. Then they take the mean of these distances and subtract from it the analogously obtained mean for the average vector of vector embeddings of words that represent men. Formally they take the relative norm distance between  $X$  and  $Y$  to be:

<sup>6</sup>By the cosine distance in the literature the authors mean 1-cosine similarity; note however that this terminology is slightly misleading, as mathematically it is not a distance measure, because it does not satisfy the triangle inequality, as generally  $\text{dist}(A, C) \not\leq \text{dist}(A, B) + \text{dist}(B, C)$ ; we'll keep using this mainstream terminology though.

## 1.4 Methodological problems with cosine-based measures of bias

One issue to consider is the selection of attributes for bias measurement. The word lists used in the literature are often fairly small (5-50). While the papers do employ statistical tests to measure the uncertainty involved, we will later on argue that these methods are not proper for the goal at hand and show that a more appropriate use of statistical methods leads to estimates of uncertainty that are rather epistemologically pessimistic.

Let's think about MAC, using the case of religion-related stereotypes. In the original paper, words from all three religions were compared against all of the stereotypes. One reason this is problematic is that no distinction between cases in which the stereotype is associated with a given religion, as opposed to the situation in which it is associated with another one, is made. This is problematic, as not all of the stereotypical words have to be considered as harmful for all of the religions. One should investigate the religions separately as some of them may have stronger harmful associations than others.

The interpretation of the results is also a challenge. In [6] we can find summaries of average cosine distances per group (such as gender, race, or religion). For instance, for religion, here is the relevant fragment of table:

(MAC stand for mean average cosine similarity, although in reality the table contains mean cosine distances). What may attract attention is the fact that the value of cosine distance in "Biased" category is already quite high (i.e. close to 1) even before debiasing. High cosine distance indicates low cosine similarity between values. One could think that average cosine similarity equal to approximately 0.141 is not significant enough to consider it as bias. However, the authors still aim to mitigate such "bias" to make the distance even larger. Methodologically the question is, on what basis is this small similarity still considered as a proof of the presence of bias, and whether these small changes are meaningful.

The underlying problem here is that in the paper there is no control group. One should also include control groups to have a way of comparing the results for a supposedly stereotyped group with the results for sets of neutral or human-related neutral words. In our approach later on, we distinguish between stereotypes associated with a given group, stereotypes associated with different groups, and introduce control groups: neutral words and stereotype-free human predicates.

## 1.5 Metrics that pre-average are a bad guide

In contrast, statistical intervals will help us decide whether a given cosine similarity is high enough to consider the words to be more similar than if we chose them at random. We will use highest posterior density intervals, in line with Bayesian methodology.

Crucially, these approaches use means of mean average cosine similarities to measure similarity between protected word and harmful stereotypes. If one takes a closer look at the individual values that are taken for the calculations, it turns out that there are quite a few outliers and surprisingly dissimilar words. This issue will become transparent when we inspect the visualizations of individual cosine distances, following the idea that one of the first steps to understand data is to look at it.

With such a method the uncertainty involved is not really considered which makes it even more difficult to give reasonable interpretations of the results. We propose the use of Bayesian method to obtain some understanding of the influence the uncertainty has on the interpretation of final results.

$s(X, Y, A, B)$  is the statistic used in the significance test, and the  $p$ -value is obtained by bootstrapping: it is the frequency of  $s(X_i, Y_i, A, B) > s(X, Y, A, B)$  for all equally sized partitions  $X_i, Y_i$  of  $X \cup Y$ . The effect size is computed by normalizing the difference in means as follows:

$$bias(A, B) = \frac{\mu(\{s(x, A, B)\}_{x \in X}) - \mu(\{s(y, A, B)\}_{y \in Y})}{\sigma(\{s(w, A, B)\}_{w \in X \cup Y})} \quad (7)$$

The t-tests they employ are run on average cosines used to calculate MAC.

Check Manzani

sensitivity to choice of words

crossref to a list of words an explanation

add p-values in the table

table the table

bliskość  
znaczeniowa a ko-  
kurencja

comparison should  
be made to actual  
frequencies, to sepa-  
rate bias caused by  
cosine similarity to  
upstream bias

unclear czy to się  
przekłada na down-  
stream effect

odległość nie musi  
uchwytywać relacji  
semantycznej żeby  
ocenić bias

założenie, bardziej  
że jeżeli odległości  
przekładają się na  
downstream, to  
warto patrzeć na  
odległości, nawet  
bez głębszej filo-  
zoficznej interpre-  
tacji ich

## 2 Pre-averaging and manufactured certainty

### 2.1 General problem with pre-averaging

### 2.2 Simulations for WEAT

### 2.3 Simulations for MAC

## 3 Bayesian method

### 3.1 Bernstein approach

### 3.2 Chasing metrics

### 3.3 Bayesian method introduction

### 3.4 Existing applications to NLP and perhaps to bias

### 3.5 Model

## 4 Results

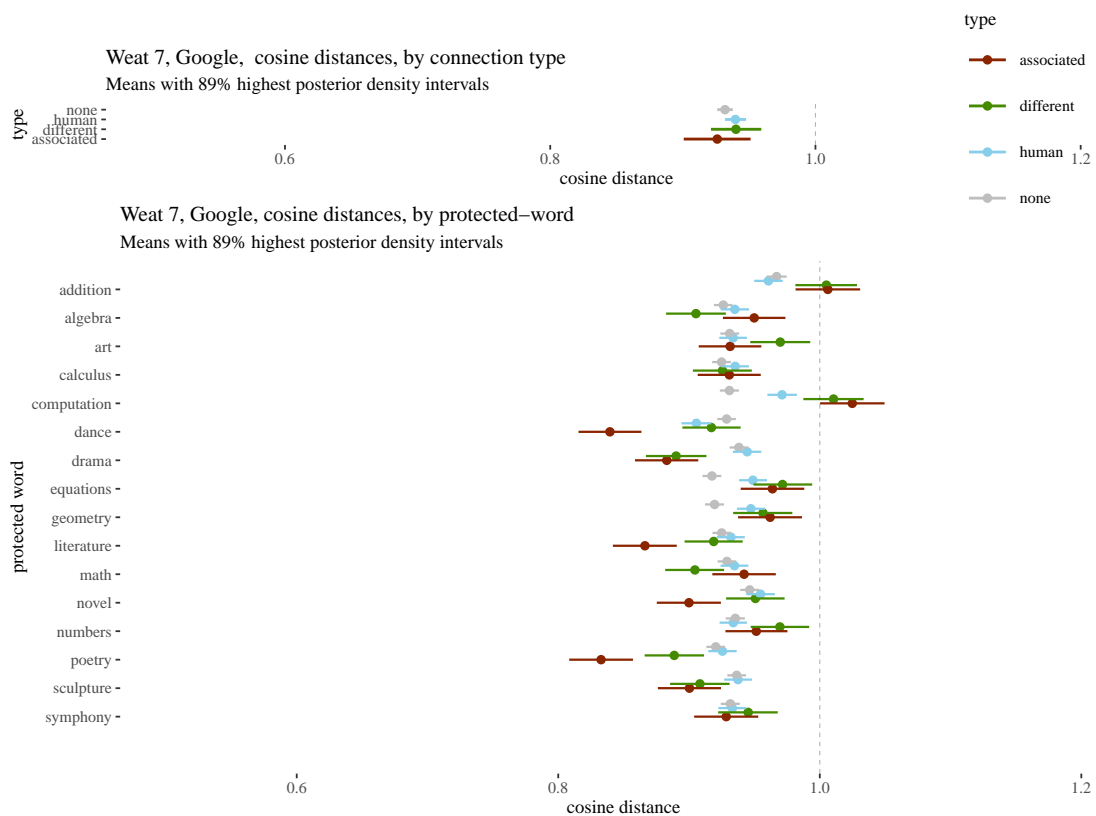


Figure 1: dsds

## 5 Discussion and summary

## 6 APPENDIX

### 6.1 Examples of WEAT and MAC calculations

#### Word lists

- [1] Tolga Bolukbasi, Kai-Wei Chang, James Y. Zou, Venkatesh Saligrama, and Adam Kalai. 2016. Man is to computer programmer as woman is to homemaker? Debiasing word embeddings. *CoRR* abs/1607.06520, (2016). Retrieved from <http://arxiv.org/abs/1607.06520>
- [2] Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science* 356, 6334 (April 2017), 183–186. DOI:<https://doi.org/10.1126/science.aal4230>
- [3] Nikhil Garg, Londa Schiebinger, Dan Jurafsky, and James Zou. 2018. Word embeddings quantify 100 years of gender and ethnic stereotypes. *Proceedings of the National Academy of Sciences* 115, 16 (April 2018), E3635–E3644. DOI:<https://doi.org/10.1073/pnas.1720347115>
- [4] Hila Gonen and Yoav Goldberg. 2019. Lipstick on a pig: Debiasing methods cover up systematic gender biases in word embeddings but do not remove them. In *Proceedings of the 2019 conference of the north American chapter of the association for computational linguistics: Human language technologies, volume 1 (long and short papers)*, Association for Computational Linguistics, Minneapolis, Minnesota, 609–614. DOI:<https://doi.org/10.18653/v1/N19-1061>
- [5] Anne Lauscher and Goran Glavas. 2019. Are we consistently biased? Multidimensional analysis of biases in distributional word vectors. *CoRR* abs/1904.11783, (2019). Retrieved from <http://arxiv.org/abs/1904.11783>
- [6] Thomas Manzini, Yao Chong Lim, Yulia Tsvetkov, and Alan W Black. 2019. Black is to criminal as caucasian is to police: Detecting and removing multiclass bias in word embeddings. Retrieved from <https://arxiv.org/abs/1904.04047>
- [7] Malvina Nissim, Rik van Noord, and Rob van der Goot. 2020. Fair is better than sensational: Man is to doctor as woman is to doctor. *Computational Linguistics* 46, 2 (June 2020), 487–497. DOI:[https://doi.org/10.1162/coli\\_a\\_00379](https://doi.org/10.1162/coli_a_00379)
- [8] Brian A. Nosek, Mahzarin R. Banaji, and Anthony G. Greenwald. 2002. Harvesting implicit group attitudes and beliefs from a demonstration web site. *Group Dynamics: Theory, Research, and Practice* 6, 1 (2002), 101–115. DOI:<https://doi.org/10.1037/1089-2699.6.1.101>