# *Sense Embeddings are also Biased* – Evaluating Social Biases in Static and Contextualised Sense Embeddings

**Anonymous ACL submission**

## Abstract

Sense embedding learning methods learn different embeddings for the different senses of an ambiguous word. One sense of an ambiguous word might be socially biased while its other senses remain unbiased. In comparison to the numerous prior work evaluating the social biases in pretrained word embeddings, the biases in sense embeddings have been relatively under studied. In this paper, we create a benchmark dataset for evaluating the social biases in sense embeddings and propose novel sense-specific bias evaluation measures. We conduct an extensive evaluation of multiple static and contextualised sense embeddings for various types of social biases using the proposed measures. Our experimental results show that even in cases where no biases are found at word-level, there still exist worrying levels of social biases at sense-level, which are often ignored by the word-level bias evaluation measures.

## 1 Introduction

Word embedding methods can be broadly classified into *static* (Pennington et al., 2014; Mikolov et al., 2013) vs. *contextualised* (Devlin et al., 2019a; Peters et al., 2018) embeddings depending on whether a word is represented by the same vector in all of its contexts. On the other hand, sense embedding learning methods use different vectors to represent the different senses of an ambiguous word (Reisinger and Mooney, 2010; Neelakantan et al., 2014; Loureiro and Jorge, 2019). Although numerous prior work has studied social biases in static and contextualised word embeddings, social biases in sense embeddings remain under explored (Kaneko and Bollegala, 2019, 2021a,a; Ravfogel et al., 2020; Dev et al., 2019; Schick et al., 2021; Wang et al., 2020).

Even if a word embedding is unbiased, some of its senses could still be associated with unfair social biases. For example, consider the ambiguous word *black*, which has two adjectival senses according to the WordNet (Miller, 1998): (1) black as a *colour* (*being of the achromatic colour of maximum darkness*, sense-key=**black%3:00:01**) and (2) black as a *race* (*of or belonging to a racial group especially of sub-Saharan African origin*, sense-key=**black%3:00:02**). However, only the second sense of *black* is often associated with racial biases. Owing to (a) the lack of evaluation benchmarks for sense embeddings, and (b) it is not being clear how to extend the bias evaluation methods proposed for static and contextualised embeddings to evaluate sense embeddings, existing social bias evaluation datasets and metrics do not consider multiple senses of words, thus not suitable for evaluating biases in sense embeddings.

To address this gap, we evaluate social biases in state-of-the-art (SoTA) static sense embeddings such as LMMS (Loureiro and Jorge, 2019) and ARES (Scarlini et al., 2020), as well as contextualised sense embeddings such as Sense-BERT (Levine et al., 2020). To the best of our knowledge, we are the first to conduct a systematic evaluation of social biases in sense embeddings. Specifically, we make two main contributions in this paper:

- First, to evaluate social biases in static sense embeddings, we extend previously proposed benchmarks for evaluating social biases in static (sense-insensitive) word embeddings by manually assigning sense ids to the words considering their social bias types expressed in those datasets (§ 3).

- Second, to evaluate social biases in sense-sensitive contextualised embeddings, we create Sense-Sensitive Social Bias (**SSSB**) dataset, a novel template-based dataset containing sentences annotated for multiple senses of an ambiguous word considering its

stereotypical social biases (§ 5).

Our experiments show that, similar to word embeddings, sense embeddings also encode worrying levels of social biases. Using SSSB, we show that the proposed bias evaluation measures for sense embeddings capture different types of social biases encoded in existing SoTA sense embeddings. More importantly, we see that even when social biases cannot be observed at word-level, such biases are still prominent at sense-level, raising concerns on existing evaluations that consider only word-level social biases.

## 2 Related Work

Our focus in this paper is the evaluation of social biases in English and *not* debiasing methods. We defer the analysis for languages other than English and developing debiasing methods for sense embeddings to future work. Hence, we limit the discussion here only to bias evaluation methods.

**Biases in Static Embeddings:** The Word Embedding Association Test (**WEAT**; Caliskan et al., 2017) evaluates the association between two sets of target concepts (e.g. *male* vs. *female*) and two sets of attributes (e.g. Pleasant (*love, cheer*, etc.) vs. Unpleasant (*ugly, evil*, etc.)). Here, the association is measured using the cosine similarity between word embeddings. Ethayarajh et al. (2019) showed that WEAT systematically overestimates the social biases and proposed relational inner-product association (**RIPA**), a subspace projection method, to overcome this problem. Word Association Test (**WAT**; Du et al., 2019) calculates a gender information vector for each word in an association graph (Deyne et al., 2019) by propagating information related to masculine and feminine words. Additionally, word analogies are used to evaluate gender bias in static embeddings (Bolukbasi et al., 2016; Manzini et al., 2019; Zhao et al., 2018). Loureiro and Jorge (2019) showed specific examples of gender bias in static sense embeddings. However, these datasets do not consider word senses, hence unfit for evaluating social biases in sense embeddings.

**Biases in Contextualised Embeddings:** May et al. (2019) extended WEAT to sentence encoders by creating artificial sentences using templates and used cosine similarity between the sentence embeddings as the association metric. Kurita et al. (2019) proposed the log-odds of the target and prior probabilities of the sentences computed by masking respectively only the target vs. both target and attribute words. Nadeem et al. (**StereoSet**; 2020) created a human annotated contexts of social bias types, while Nangia et al. (2020) proposed Crowdsourced Stereotype Pairs benchmark (**CrowS-Pairs**). These benchmarks use sentence pairs of the form "*She is a nurse/doctor*". StereoSet calculates log-odds by masking the modified tokens (*nurse, doctor*) in a sentence pair, whereas CrowS-Pairs calculates log-odds by masking their unmodified tokens (*She, is, a*). Kaneko and Bollegala (2021b) proposed All Unmasked Likelihood (**AUL**) and AUL with Attention weights (**AULA**), which calculate log-likelihood by predicting all tokens in a test case, given the contextualised embedding of the unmasked input.

## 3 Evaluation Metrics for Social Biases in Static Sense Embeddings

We extend the WEAT and WAT datasets that have been frequently used in prior work for evaluating social biases in static word embeddings such that they can be used to evaluate sense embeddings. These datasets compare the association between a target word $w$ and some (e.g. pleasant or unpleasant) attribute $a$, using the cosine similarity, $\cos(\boldsymbol{w}, \boldsymbol{a})$, computed using the static word embeddings $\boldsymbol{w}$ and $\boldsymbol{a}$ of respectively $w$ and $a$. Given two same-size sets of *target* words $\mathcal{X}$ and $\mathcal{Y}$, with two sets of *attribute* words $\mathcal{A}$ and $\mathcal{B}$. The bias score, $s(\mathcal{X}, \mathcal{Y}, \mathcal{A}, \mathcal{B})$, for each target is calculated as follows:

$$s(\mathcal{X}, \mathcal{Y}, \mathcal{A}, \mathcal{B}) = \sum_{\boldsymbol{x} \in \mathcal{X}} w(\boldsymbol{x}, \mathcal{A}, \mathcal{B}) - \sum_{\boldsymbol{y} \in \mathcal{Y}} w(\boldsymbol{y}, \mathcal{A}, \mathcal{B}) \quad (1)$$

$$w(\boldsymbol{t}, \mathcal{A}, \mathcal{B}) = \operatorname*{mean}_{\boldsymbol{a} \in \mathcal{A}} \cos(\boldsymbol{t}, \boldsymbol{a}) - \operatorname*{mean}_{\boldsymbol{b} \in \mathcal{B}} \cos(\boldsymbol{t}, \boldsymbol{b}) \quad (2)$$

Here, $\cos(\boldsymbol{a}, \boldsymbol{b})$ is the cosine similarity between the embeddings $\boldsymbol{a}$ and $\boldsymbol{b}$. The one-sided $p$-value for the permutation test for $\mathcal{X}$ and $\mathcal{Y}$ is calculated as the probability of $s(\mathcal{X}_i, \mathcal{Y}_i, \mathcal{A}, \mathcal{B}) > s(\mathcal{X}, \mathcal{Y}, \mathcal{A}, \mathcal{B})$. The effect size is calculated as the normalised measure given by (3).

$$\frac{\operatorname*{mean}_{x \in \mathcal{X}} w(x, \mathcal{A}, \mathcal{B}) - \operatorname*{mean}_{y \in \mathcal{Y}} w(y, \mathcal{A}, \mathcal{B})}{\operatorname*{sd}_{t \in \mathcal{X} \cup \mathcal{Y}} w(t, \mathcal{A}, \mathcal{B})} \quad (3)$$

We repurpose these datasets for evaluating *sense* embeddings as follows. For each word in WEAT, we manually assign a sense id considering the bias type in which it is used for evaluation. For example, the word "violet" in the *Flowers* group is

assigned the sense id "violet%1:20:00::", which has the meaning – *any of numerous low-growing violas with small flowers*, according to the WordNet. We then measure the cosine similarity between two words using their corresponding sense embeddings.

WAT considers only gender bias and calculates the gender information vector for each word in a word association graph created with Small World of Words project (SWOWEN; Deyne et al., 2019) by propagating information related to masculine and feminine words $(w_m^i, w_f^i) \in \mathcal{L}$ using a random walk approach (Zhou et al., 2003). It is non-trivial to pre-specify the sense of a word in a large word association graph considering the paths followed by a random walk. The gender information is encoded as a vector $(b_m, b_f)$ in 2 dimensions, where $b_m$ and $b_f$ denote the masculine and feminine orientations of a word, respectively. The bias score of a word is defined as $\log(b_m/b_f)$. The gender bias of word embeddings are evaluated using the Pearson correlation coefficient between the bias score of each word and the score given by (4), computed as the average over the differences of cosine similarities between masculine and feminine words.

$$\frac{1}{|\mathcal{L}|} \sum_{i=1}^{|\mathcal{L}|} \left( \cos(w, w_m^i) - \cos(w, w_f^i) \right) \qquad (4)$$

To evaluate gender bias in sense embeddings, we compare each sense $s_i$ of the target word $w$ against each sense $a_j$ of a word selected from the association graph using their corresponding sense embeddings, $\boldsymbol{s}_i, \boldsymbol{a}_j$, and use the maximum similarity over all pairwise combinations (i.e. $\max_{i,j} \cos(\boldsymbol{s}_i, \boldsymbol{a}_j)$) as the word association measure. Measuring similarity between two words as the maximum similarity over all candidate senses of each word is based on the assumption that two words in a word-pair would mutually disambiguate each other in an association-based evaluation (Pilehvar and Camacho-Collados, 2019), and has been used as a heuristic for disambiguating word senses (Reisinger and Mooney, 2010).

## 4 Sense-Sensitive Social Bias Dataset

Contextualised embeddings such as the ones generated by masked language models (MLMs) return different vectors for the same word in different contexts. However, the datasets discussed in § 3 do not provide contextual information for words and cannot be used to evaluate contextualised embeddings. Moreover, the context in which an ambiguous word occurs determines its word sense. Contextualised sense embedding methods such as Sense-BERT (fine-tuned using WordNet super senses), have shown to capture word sense information in their contextualised embeddings (Zhou and Bollegala, 2021).

| Category | noun vs. verb | ethnicity vs. colour | nationality vs. language |
|---|---|---|---|
| #pleasant words | 14 | 5 | 18 |
| #unpleasant words | 18 | 5 | 15 |
| #target words | 6 | 1 | 16 |
| #templates | 1 | 4 | 4 |
| #test cases | 324 | 733 | 2304 |

Table 1: Statistics of the the SSSB dataset.

CrowS-Pairs and StereoSet datasets were proposed for evaluating contextualised word embeddings. Specifically, an MLM is considered to be unfairly biased if it assigns higher pseudo log-likelihood scores for stereotypical sentences, $S^{\text{st}}$, than anti-stereotypical ones, $S^{\text{at}}$. However, both of those datasets do not consider multiple senses of words and cannot be used to evaluate social biases in contextualised sense embeddings.

To address this problem, we create the Sense-Sensitive Social Bias (SSSB) dataset, containing template-generated sentences covering multiple senses of ambiguous words for three types of social biases: *gender*, *race* and *nationality*. To the best of our knowledge, SSSB is the first-ever dataset created for the purpose of evaluating social biases in sense embeddings.[1] Table 1 shows the summary statistics of the SSSB dataset. Next, we describe the social biases covered in this dataset.

### 4.1 Nationality vs. Language Bias

These examples cover social biases related to a nationality (racial) or a language (non-racial). Each test case covers two distinct senses and the following example shows how they represent biases. *Japanese people are nice* is an anti-stereotype for *Japanese* as a nationality because it is associated with a pleasant attribute (i.e. *nice*) in this example sentence. On the other hand, *Japanese people are stupid* is a stereotype for *Japanese* as a nationality because it is associated with an unpleasant attribute (i.e. *stupid*). These can be considered as examples of racial biases.

Likewise, for the language sense of Japanese we create examples as follows. *Japanese language is*

---

[1] The dataset and evaluation scripts will be publicly released upon paper acceptance.

*difficult to understand* is a stereotype for *Japanese* as a language because it is associated with an unpleasant attribute (i.e. *difficult*). On the other hand, *Japanese language is easy to understand* is an anti-stereotype for *Japanese* as a language because it is associated with a pleasant attribute (i.e. *easy*).

In SSSB, we indicate the sense-type, WordNet sense-id and the type of social bias in each example as follows:

> *Japanese people are beautiful.*
> [nationality, japanese%1:18:00::, anti]

Here, sense-type is nationality, sense-id is *japanese%1:18:00::* and the bias is anti (we use the labels *anti* and *stereo* to denote respectively anti-stereotypical and stereotypical biases).

We use the likelihood scores returned by an MLM to nationality vs. language sentence pairs as described further in § 5 to evaluate social biases in MLMs. Essentially, if the likelihood score returned by an MLM for the example that uses an unpleasant attribute is higher than the one that uses a pleasant attribute for a member in the disadvantaged group, then we consider the MLM to be socially biased. Note that one could drop the modifiers such as *people* and *language* and simplify these examples such as *Japanese are stupid* and *Japanese is difficult* to generate additional test cases. However, the sense-sensitive embedding methods might find it difficult to automatically disambiguate the correct senses without the modifiers such as *language* or *people*. Therefore, we always include these modifiers when creating examples for nationality vs. language bias in the SSSB dataset.

## 4.2 Ethnicity vs. Colour Bias

The word *black* can be used to represent the ethnicity (black people) or the colour. We create examples that distinguish these two senses of black as in the following example. *Black people are friendly* represents an anti-stereotype towards *black* because it is associated with a pleasant attribute (i.e. *friendly*) of a disadvantaged group whereas, *Black people are arrogant* represents a stereotype because it is associated with an unpleasant attribute (i.e. *arrogant*).

On the other hand, for the colour black, *The black dress is elegant* represents an anti-stereotype because it is associated with a pleasant attribute (i.e. *elegant*), whereas *The black dress is ugly* represents a stereotype because it is associated with an

unpleasant attribute (i.e. *ugly*). As described in § 5, if the likelihood score returned by an MLM for the example that uses an unpleasant attribute is higher than a one that uses a pleasant attribute, then we consider the MLM to be socially biased.

## 4.3 Gender Bias in Noun vs. Verb Senses

To create sense-related bias examples for gender[2], we create examples based on occupations. In particular, we consider six occupations *engineer*, *nurse*, *judge*, *mentor*, *(tour) guide*, and *carpenter*, which can be used in a noun sense (e.g. *engineer is a person who uses scientific knowledge to solve practical problems*, *nurse is a person who looks after patients*, etc.) as well as in a verb sense expressing the action performed by a person holding the occupation (e.g. *design something as an engineer*, *nurse a baby*, etc.). Note that the ambiguity here is in the occupation (noun) vs. action (verb) senses and not in the gender, whereas the bias is associated with the gender of the person holding the occupation.

To illustrate this point further, consider the following examples. *She is a talented engineer* is considered as an anti-stereotypical example for the noun sense of *engineer* because females (here the disadvantaged group) are not usually associated with pleasant attributes (i.e. *talented*) with respect to this occupation (i.e. *engineer*). *He is a talented engineer* is considered as a stereotypical example for the noun sense of engineer because males (i.e. advantaged group) are usually associated with pleasant attributes with regard to this occupation. As described in § 5, if an MLM assigns a higher likelihood to the stereotypical example (second sentence) than the anti-stereotypical example (first sentence), then that MLM is considered to be gender biased.

On the other hand, *She is a clumsy engineer* is considered to be a stereotypical example for the noun sense of engineer because females (i.e. disadvantaged group) are historically associated with such unpleasant attributes (i.e. *clumsy*) with respect to such male-dominated occupations. Likewise, *He is a clumsy engineer* is considered as an anti-stereotypical example for the noun sense of engineer because males (i.e. advantaged group) are not usually associated with such unpleasant attributes (i.e. *clumsy*). Here again, if an MLM assigns a higher likelihood to the stereotypical example (first sentence) than the anti-stereotypical

---

[2] We consider only male and female genders in this work

4

example (second sentence), then it is considered to be gender biased. Note that the evaluation direction with respect to male vs. female pronouns used in these examples is opposite to that in the previous paragraph because we are using an unpleasant attribute in the second set of examples.

Verb senses are also used in the sentences that contain gender pronouns in SSSB. For example, for the verb sense of *engineer*, we create examples as follows: *She used novel material to engineer the bridge*. Here, the word engineer is used in the verb sense in a sentence where the subject is a female. The male version of this example is as follows: *He used novel material to engineer the bridge*. In this example, a perfectly unbiased MLM should not systematically prefer one sentence over the other between the two sentences both expressing the verb sense of the word *engineer*.

## 5  Evaluation Metrics for Social Biases in Contextualised Sense Embeddings

For a contextualised (word/sense) embedding under evaluation, we compare its pseudo-likelihood scores for stereotypical and anti-stereotypical sentences for each sense of a word in SSSB, using AUL (Kaneko and Bollegala, 2021b).[3] AUL is known to be robust against the frequency biases of words and provides more reliable estimates compared to the other metrics for evaluating social biases in MLMs. Following the standard evaluation protocol, we provide AUL the complete sentence $S = w_1, \ldots, w_{|S|}$, which contains a length $|S|$ sequence of tokens $w_i$, to an MLM with pretrained parameters $\theta$. We first compute $\text{PLL}(S)$, the Pseudo Log-Likelihood (PLL) for predicting all tokens in $S$ excluding begin and end of sentence tokens, given by (5).

$$\text{PLL}(S) \coloneqq \frac{1}{|S|} \sum_{i=1}^{|S|} \log P(w_i | S; \theta) \quad (5)$$

Here, $P(w_i | S; \theta)$ is the probability assigned by the MLM to token $w_i$ conditioned on $S$. The fraction of sentence-pairs in SSSB, where higher PLL scores are assigned to the stereotypical sentence than the anti-stereotypical one is considered as the AUL *bias score* of the MLM associated with the contextualised embedding, and is given by (6).

---

[3]The attention-weighted variant (AULA) is not used because contextualised sense embeddings have different structures of attention from contextualised embeddings, and it is not obvious which attention to use in the evaluations.

$$\text{AUL} = \left( \frac{100}{N} \sum_{(S^{\text{st}}, S^{\text{at}})} \mathbb{I}(\text{PLL}(S^{\text{st}}) > \text{PLL}(S^{\text{at}})) \right) - 50 \quad (6)$$

Here, $N$ is the total number of sentence-pairs in SSSB and $\mathbb{I}$ is the indicator function, which returns 1 if its argument is True and 0 otherwise. AUL score given by (6) falls within the range $[-50, 50]$ and an unbiased embedding would return bias scores close to 0, whereas bias scores less than or greater than 0 indicate bias directions towards respectively the anti-stereotypical or stereotypical examples.

## 6  Experiments

### 6.1  Bias in Static Embeddings

To evaluate biases in static sense embeddings, we select two current state-of-the-art embeddings: LMMS[4] (Loureiro and Jorge, 2019) and ARES[5] (Scarlini et al., 2020). In addition to WEAT and WAT datasets described in § 3, we also use SSSB to evaluate static sense embeddings using the manually assigned sense ids for the target and attribute words, ignoring their co-occurring contexts. LMMS and ARES sense embeddings associate each sense of a lexeme with a sense key and a vector, which we use to compute cosine similarities as described in § 3. To compare the biases in a static sense embedding against a corresponding sense-insensitive static word embedding version, we compute a static word embedding $\boldsymbol{w}$, for an ambiguous word $w$ by taking the average (**avg**) over the sense embeddings $\boldsymbol{s_i}$ for all of $w$'s word senses as given in (7), where $M(w)$ is the total number of senses of $w$.

$$\boldsymbol{w} = \frac{\sum_i^{M(w)} \boldsymbol{s_i}}{M(w)}. \quad (7)$$

This would simulate the situation where the resultant embeddings are word-specific but not sense-specific, while still being comparable to the original sense embeddings in the same vector space.

From Table 2 we see that in WEAT[6] in all categories considered, sense embeddings always report a higher bias compared to their corresponding

---

[4]https://github.com/danlou/LMMS
[5]http://sensembert.org
[6]Three bias types (European vs. African American, Male vs. Female, and Old vs. Young) had to be excluded because these biases are represented using personal names that are not covered by LMMS and ARES sense embeddings.

5

| Dataset | LMMS word/sense | ARES word/sense |
|---|---|---|
| **WEAT** | | |
| Flowers vs Insects | 1.63/**2.00** | 1.58/**2.00** |
| Instruments vs Weapons | 1.42/**2.00** | 1.37/**1.99** |
| Math vs Art | 1.52/**1.83** | 0.98/**1.45** |
| Science vs Art | 1.38/**1.66** | 0.92/**1.44** |
| Physical vs. Mental condition | 0.42/**0.64** | -0.12/**-0.77** |
| **WAT** | **0.53**/0.41 | **0.46**/0.31 |
| **SSSB** | | |
| black (ethnicity) | **5.36**/4.64 | 5.40/**5.67** |
| black (colour) | **5.36**/1.64 | **5.40**/4.83 |
| nationality | **7.78**/7.01 | **6.94**/5.75 |
| language | 7.78/**8.23** | 6.94/**7.38** |
| noun | 0.34/**0.39** | 0.09/**0.16** |
| verb | **0.34**/0.26 | **0.09**/0.06 |

Table 2: Bias in LMMS and ARES Static Sense Embeddings. In each row, between sense-insensitive word embeddings and sense embeddings, the larger deviation from 0 is shown in bold.



Figure 1: Effect of the dimensionality of sense embeddings (LMMS) and word embeddings (LMMS-average).

sense-insensitive word embeddings. This shows that even if there are no biases at the word-level, we can still observe social biases at the sense-level in WEAT. However, in the WAT dataset, which covers only gender-related biases, we see word embeddings to have higher biases than sense embeddings. This indicates that in WAT gender bias is more likely to be observed in static word embeddings than in static sense embeddings.

In SSSB, the word embeddings always report the same bias scores for the different senses of the ambiguous word because static word embeddings are neither sense nor context sensitive. As aforementioned, the word "black" is a bias-neutral word with respect to the colour sense, while it often has a social bias in the racial sense. Consequently, for *black* we see a higher bias score for its ethnic sense than for its colour sense in both LMMS and ARES sense embeddings.

In the bias scores reported for *nationality* vs. *language* senses, we find that *nationality* obtains higher biases at word-level whereas that for *language* is higher at the sense-level in both LMMS and ARES. Unlike *black*, where the two senses (colour vs. ethnic) are distinct, the two senses *nationality* and *language* are much closer because in many cases (e.g. Japanese, Chinese, Spanish, French etc.) languages and nationalities are used interchangeably to refer to the same set of entities. Interestingly, the *language* sense is assigned a slightly higher bias score than the *nationality*
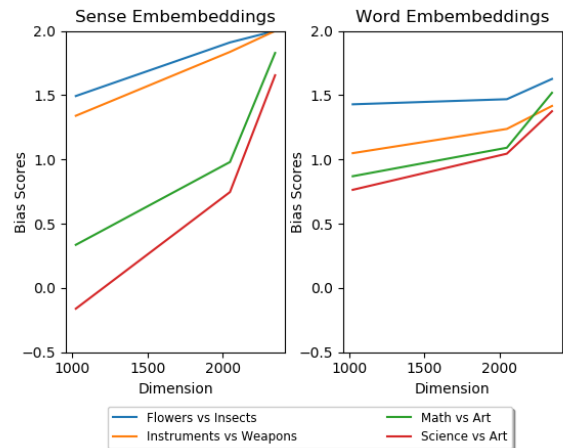
sense in both LMMS and ARES sense embeddings. Moreover, we see that the difference between the bias scores for the two senses in *colour* vs. *ethnicity* (for black) as well as *nationality* vs. *language* is more in LMMS compared to that in ARES sense embeddings.

Between noun vs. verb senses of occupations, we see a higher gender bias for the noun sense than the verb sense in both LMMS and ARES sense embeddings. This agrees with the intuition that gender biases exist with respect to occupations and not so much regarding what actions/tasks are carried out by the persons holding those occupations. Compared to the word embeddings, there is a higher bias for the sense embeddings in the noun sense for both LMMS and ARES. This trend is reversed for the verb sense where we see higher bias scores for the word embeddings than the corresponding sense embeddings in both LMMS and ARES. Considering that gender is associated with the noun than verb sense of occupations in English, this shows that there are hidden gender biases that are not visible at the word-level but become more apparent at the sense-level. This is an important factor to consider when evaluating gender biases in word embeddings, which has been largely ignored thus far in prior work.

To study the relationship between the dimensionality of the embedding space and the social biases it encodes, we compare 1024, 2048 and 2348 dimensional LMMS static sense embeddings and their corresponding word embeddings (computed using (7)) on the WEAT dataset in Figure 1. We see that all types of social biases increase with the di-

| | base | large |
|---|---|---|
| Dataset | BERT/SenseBERT | BERT/SenseBERT |
| **CrowS-Pairs** | **-1.66**/0.99 | **-3.58**/2.45 |
| **StereoSet** | -1.09/**8.31** | -1.47/**6.51** |
| **SSSB** | | |
| ethnicity | 10.19/**14.81** | **-17.59**/0.00 |
| colour | **-6.64**/-2.96 | -8.88/**9.84** |
| nationality | 5.79/**15.34** | 4.28/**8.10** |
| language | -0.17/**-2.95** | **6.25**/-3.82 |
| noun | 10.42/**14.06** | 3.13/3.13 |
| verb | **12.89**/-3.74 | 0.22/**-15.44** |

Table 3: Bias in BERT and SenseBERT contextualised word/sense embeddings. In each row, between the AUL bias scores for the word vs. sense embeddings, the larger deviation from 0 is shown in bold.

mensionality for both word and sense embeddings. This is in agreement with Silva et al. (2021) who also reported that increasing model capacity in contextualised word embeddings does not necessarily remove their unfair social biases. Moreover, in higher dimensionalities sense embeddings show a higher degree of social biases than the corresponding (sense-insensitive) word embeddings.

## 6.2 Bias in Contextualised Embeddings

To evaluate biases in contextualised sense embeddings, we use SenseBERT[7] (Levine et al., 2020), which is a fine-tuned version of BERT[8] (Devlin et al., 2019b) to predict supersenses in the WordNet. For both BERT and SenseBERT, we use base and large pretrained models of dimensionalities respectively 768 and 1024. Using AUL, we compare biases in BERT and SenseBERT using SSSB, CrowS-Pairs and StereoSet[9] datasets. Note that unlike SSSB, CrowS-Pairs and StereoSet *do not* annotate for word senses, hence cannot be used to evaluate sense-specific biases.

Table 3 compares biases in contextualised word/sense embeddings. For both base and large versions, we see that in CrowS-Pairs, BERT to be more biased than SenseBERT, whereas the opposite is true in StereoSet. Among the nine bias types included in CrowS-Pairs, gender bias related test instances are the second most frequent following *racial* bias. On the other hand, gender bias re-

---

lated examples are relatively less frequent in StereoSet (cf. gender is the third most frequent bias type in StereoSet after *racial* and *occupational* biases). This difference in the composition of bias types explains why the bias score of BERT is higher in CrowS-Pairs, while the same is higher for SenseBERT in StereoSet. In SSSB, in 8 out of the 12 cases SenseBERT demonstrates equal or higher absolute bias scores than BERT. This result shows that even in situations where no biases are observed at the word-level, there can still be significant degrees of biases at the sense-level. In some cases (e.g. *verb* sense in base models and *colour*, *language* and *verb* senses for the large models), we see that the direction of bias is opposite between BERT and SenseBERT. Moreover, comparing against the corresponding bias scores reported by the static word/sense embeddings in Table 2, we see higher bias scores reported by the contextualised word/sense embeddings in Table 3. Therefore, we recommend future work studying social biases to consider not only word embedding models but also sense embedding models.

## 7 Gender Biases in SSSB

In this section, we further study the gender-related biases in static and contextualised word and sense embeddings using the noun vs. verb sense instances (described in § 4.3) in the SSSB dataset. To evaluate the gender bias in contextualised word/sense embeddings we use AUL on test sentences in SSSB noun vs. verb category. To evaluate the gender bias in static embeddings, we follow Bolukbasi et al. (2016) and use the cosine similarity between (a) the static word/sense embedding of the occupation corresponding to its noun or verb sense and (b) the gender directional vector $g$, given by (8).

$$g = \frac{1}{|\mathcal{C}|} \sum_{(m,f) \in \mathcal{C}} (m - f) \qquad (8)$$

Here, $(m, f)$ are male-female word pairs used by Kaneko and Bollegala (2019) such as (*he*, *she*) and $m$ and $f$ respectively denote their word embeddings. Corresponding sense-insensitive word embeddings are computed for the 2048 dimensional LMMS embeddings using (7).

Figure 2 shows the gender biases in LMMS embeddings. Because static word embeddings are not sense-sensitive, they report the same bias scores for both noun and verb senses for each occupation. For all noun senses, we see positive (male) biases,

| stereo/anti-stereo sentences | BERT | | | SenseBERT | | |
|---|---|---|---|---|---|---|
| | stereo | anti | diff | stereo | anti | diff |
| he/she is a strong nurse | -0.45 | -0.67 | 0.22 | -15.71 | -16.64 | 0.93 |
| he/she is a professional nurse | -0.73 | -0.85 | 0.11 | -16.53 | 16.81 | 0.27 |
| As a mother/father of five, she/he carefully nurse all of her/his children | -0.16 | -0.15 | -0.01 | -18.07 | -18.24 | 0.18 |
| she/he made milk herself/himself to nurse the crying baby | -0.77 | -0.14 | -0.63 | -15.85 | -17.80 | 1.96 |

Table 4: Pseudo log-likelihood scores computed using Eq. (5) for stereo and anti-stereo sentences (shown together due to space limitations) using BERT-base and SenseBERT-base models. Here, diff = stereo - anti.



Figure 2: Gender biases found in the 2048-dimensional LMMS static sense embeddings and corresponding word embeddings computed using (7). Positive and negative cosine similarity scores with the gender directional vector (computed using (8)) represent biases towards respectively the *male* and *female* genders.
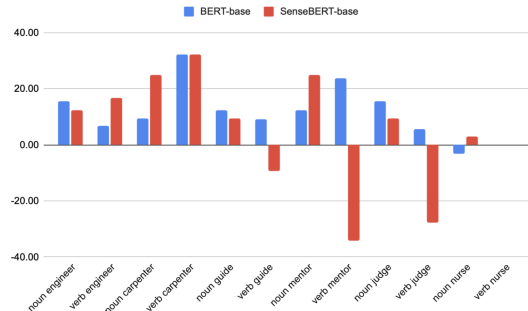


Figure 3: Gender biases found in 768-dimensional BERT-base and SenseBERT-base contextualised embeddings. Positive and negative AUL scores represent bias towards respectively the stereotypical and anti-stereotypical sentences.

except for *nurse*, which is strongly female-biased. Moreover, compared to the noun senses, the verb senses of LMMS are relatively less gender biased. This agrees with the intuition that occupations and not actions associated with those occupations are related to gender, hence can encode social biases. Overall, we see stronger biases in sense embeddings than in the word embeddings.

Figure 3 shows the gender biases in BERT/SenseBERT embeddings. Here again, we see that for all noun senses there are high stereotypical bises in both BERT and SenseBERT embeddings, except for *nurse* where BERT is slightly anti-stereotypically biased whereas SenseBERT shows a similar in magnitude but a stereotypical bias. Recall that *nurse* is stereotypically associated with the female gender, whereas other occupations are predominantly associated with males, which is reflected in the AUL scores here. Despite not fine-tuned on word senses, BERT shows different bias scores for noun/verb senses, showing its ability to capture sense-related information via contexts. The verb sense embeddings of SenseBERT of *guide*, *mentor* and *judge* are anti-stereotypical, while the corresponding BERT embeddings are stereotypical.

This shows that contextualised word and sense embeddings can differ in both magnitude as well as direction of the bias. Considering that SenseBERT is a fine-tuned version of BERT for a specific downstream NLP task (i.e. super-sense tagging), one must not blindly assume that an unbiased MLM to remain as such when fine-tuned on downstream tasks. *How social biases in word/sense embeddings change when used in downstream tasks* is an important research problem in its own right, which is beyond the scope of this paper.

A qualitative analysis is given in Table 4 where the top-two and bottom-two sentences selected from SSSB express respectively noun and verb senses of *nurse*. We see that SenseBERT has a higher preference (indicated by the high pseudo log-likelihood scores) for stereotypical examples than BERT over anti-stereotypical ones (indicated by the higher diff values).

## 8 Conclusion

We proposed novel datasets and metrics for evaluating social biases in sense embeddings. Our experiments show, for the first time that sense embeddings are also socially biased similar to word embeddings. In future work, we plan to develop debiasing methods for sense embeddings.

## 9 Ethical Considerations

In this paper we considered the relatively under explored aspect of social biases in pretrained sense embeddings. We created a new dataset for this purpose, which we name the Sense-Sensitive Social Bias (SSSB) dataset. The dataset we create is of a sensitive nature. We have included various sentences that express stereotypical biases associated with different senses of words in this dataset. We specifically considered three types of social biases in SSSB: (a) racial biases associated with a nationality as opposed to a language (e.g. *Chinese people are cunning*, *Chinese language is difficult*, etc.), (b) ethnic biases associated with the word *black* as opposed to its sense as a colour (e.g. *Black people are arrogant*, *Black dress is beautiful*, etc.) and (c) gender-related biases associated with occupations used as nouns as opposed to verbs (e.g. *She was a careless nurse*, *He was not able to nurse the crying baby*, etc.). As seen from the above-mentioned examples, by design, SSSB contains many offensive, stereotypical examples. It is intended to facilitate evaluation of social biases in sense embeddings and will be publicly released for this purpose only. We argue that SSSB should not be used to train sense embeddings. The motivation behind creating SSSB is to measure social biases so that we can make more progress towards debiasing them in the future. However, training on this data would defeat this purpose.

It is impossible to cover all types of social biases related to word senses in any single dataset. Given that our dataset is generated from a handful of manually written templates, it is far from complete. Moreover, the templates reflect the cultural and social norms of the annotators from a US-centric viewpoint. Therefore, SSSB should not be considered as an ultimate test for biases in sense embeddings. Simply because a sense embedding does not show any social biases on SSSB according to the evaluation metrics we use in this paper *does not* mean that it would be appropriate to deploy it in downstream NLP applications that require sense embeddings. In particular, task-specific fine-tuning of even bias-free embeddings can result in novel unfair biases from creeping in. Last but not least we state that the study conducted in this paper has been limited to the English language and represent social norms held by the annotators.

## References

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. In *NIPS*.

Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356:183–186.

Sunipa Dev, Tao Li, Jeff Phillips, and Vivek Srikumar. 2019. On Measuring and Mitigating Biased Inferences of Word Embeddings.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019a. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019b. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.

Simon De Deyne, Danielle J. Navarro, Amy Perfors, Marc Brysbaert, and Gert Storms. 2019. The "small world of words" english word association norms for over 12,000 cue words. *Behavior Research Methods*, 51:987–1006.

Yupei Du, Yuanbin Wu, and Man Lan. 2019. Exploring human gender stereotypes with word association test. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6132–6142, Hong Kong, China. Association for Computational Linguistics.

Kawin Ethayarajh, David Duvenaud, and Graeme Hirst. 2019. Understanding undesirable word embedding associations. In *Proceedings of the 57th Conference of the Association for Computational Linguistics*, pages 1696–1705, Florence, Italy. Association for Computational Linguistics.

Masahiro Kaneko and Danushka Bollegala. 2019. Gender-preserving debiasing for pre-trained word embeddings. In *ACL*, pages 1641–1650.

Masahiro Kaneko and Danushka Bollegala. 2021a. Debiasing pre-trained contextualised embeddings. In *Proc. of 16th conference of the European Chapter of the Association for Computational Linguistics (EACL)*.

9

Masahiro Kaneko and Danushka Bollegala. 2021b. Unmasking the mask–evaluating social biases in masked language models. *arXiv preprint arXiv:2104.07496*.

Keita Kurita, Nidhi Vyas, Ayush Pareek, Alan W Black, and Yulia Tsvetkov. 2019. Measuring bias in contextualized word representations. In *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*, pages 166–172, Florence, Italy. Association for Computational Linguistics.

Yoav Levine, Barak Lenz, Or Dagan, Ori Ram, Dan Padnos, Or Sharir, Shai Shalev-Shwartz, Amnon Shashua, and Yoav Shoham. 2020. SenseBERT: Driving some sense into BERT. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4656–4667, Online. Association for Computational Linguistics.

Daniel Loureiro and Alipio Jorge. 2019. Language modelling makes sense: Propagating representations through wordnet for full-coverage word sense disambiguation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5682–5691.

Thomas Manzini, Lim Yao Chong, Alan W Black, and Yulia Tsvetkov. 2019. Black is to criminal as caucasian is to police: Detecting and removing multiclass bias in word embeddings. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 615–621, Minneapolis, Minnesota. Association for Computational Linguistics.

Chandler May, Alex Wang, Shikha Bordia, Samuel R. Bowman, and Rachel Rudinger. 2019. On measuring social biases in sentence encoders. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 622–628, Minneapolis, Minnesota. Association for Computational Linguistics.

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in NIPS*, pages 3111–3119.

George A Miller. 1998. *WordNet: An electronic lexical database*. MIT press.

Moin Nadeem, Anna Bethke, and Siva Reddy. 2020. StereoSet: Measuring stereotypical bias in pretrained language models.

Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R. Bowman. 2020. CrowS-pairs: A challenge dataset for measuring social biases in masked language models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1953–1967, Online. Association for Computational Linguistics.

Arvind Neelakantan, Jeevan Shankar, Alexandre Passos, and Andrew McCallum. 2014. Efficient nonparametric estimation of multiple embeddings per word in vector space. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1059–1069.

Jeffery Pennington, Richard Socher, and Christopher D. Manning. 2014. Glove: global vectors for word representation. In *EMNLP*, pages 1532–1543.

Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proc. of NAACL-HLT*.

Mohammad Taher Pilehvar and Jose Camacho-Collados. 2019. WiC: the word-in-context dataset for evaluating context-sensitive meaning representations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 1267–1273, Minneapolis, Minnesota. Association for Computational Linguistics.

Shauli Ravfogel, Yanai Elazar, Hila Gonen, Michael Twiton, and Yoav Goldberg. 2020. Null it out: Guarding protected attributes by iterative nullspace projection. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 7237–7256. Association for Computational Linguistics.

Joseph Reisinger and Raymond Mooney. 2010. Multiprototype vector-space models of word meaning. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 109–117.

Bianca Scarlini, Tommaso Pasini, and Roberto Navigli. 2020. With more contexts comes better performance: Contextualized sense embeddings for all-round word sense disambiguation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3528–3539.

Timo Schick, Sahana Udupa, and Hinrich Schütze. 2021. Self-diagnosis and self-debiasing: A proposal for reducing corpus-based bias in nlp. *Computing Research Repository*, arXiv:2103.00453.

Andrew Silva, Pradyumna Tambwekar, and Matthew Gombolay. 2021. Towards a comprehensive understanding and accurate evaluation of societal biases in pre-trained transformers. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2383–2389, Online. Association for Computational Linguistics.

10

Tianlu Wang, Xi Victoria Lin, Nazneen Fatema Rajani, Bryan McCann, Vicente Ordonez, and Caiming Xiong. 2020. Double-hard debias: Tailoring word embeddings for gender bias mitigation. In *Association for Computational Linguistics (ACL)*.

Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. 2018. Learning Gender-Neutral Word Embeddings. In *Proc. of EMNLP*, pages 4847–4853.

Dengyong Zhou, Olivier Bousquet, Thomas Navin Lal, Jason Weston, and Bernhard Schölkopf. 2003. Learning with local and global consistency. In *NIPS*.

Yi Zhou and Danushka Bollegala. 2021. Learning sense-specific static embeddings using contextualised word embeddings as a proxy. In *Proc. of the 35-th Pacific Asia Conference on Language, Information and Computation (PACLIC)*.