

# On Measuring and Mitigating Biased Inferences of Word Embeddings

Sunipa Dev, Tao Li, Jeff M. Phillips, Vivek Srikumar

School of Computing  
University of Utah  
Salt Lake City, Utah, USA  
{sunipad, tli, jeffp, svivek}@cs.utah.edu

## Abstract

Word embeddings carry stereotypical connotations from the text they are trained on, which can lead to invalid inferences in downstream models that rely on them. We use this observation to design a mechanism for measuring stereotypes using the task of natural language inference. **We demonstrate a reduction in invalid inferences via bias mitigation strategies on static word embeddings (GloVe).** Further, we show that for gender bias, these techniques extend to contextualized embeddings when applied selectively only to the static components of contextualized embeddings (ELMo, BERT).

## Introduction

Word embeddings have become the de facto feature representation across NLP [13, 18, for example]. Their usefulness stems from their ability capture background information about words using large corpora as static vector embeddings—e.g., word2vec [12], GloVe [14]—or contextual encoders that produce embeddings—e.g., ELMo [15], BERT [8].

However, besides capturing word meaning, their embeddings also encode real-world biases about gender, age, ethnicity, etc. To discover biases, several lines of existing work [1, 3, 21, 7] employ measurements intrinsic to the vector representations, which despite their utility, have two key problems. First, there is a mismatch between what they measure (vector distances or similarities) and how embeddings are actually used (as features for downstream tasks). Second, contextualized embeddings like ELMo or BERT drive today’s state-of-the-art NLP systems, but tests for bias are designed for word types, not word *token* embeddings.

In this paper, we present a general strategy to probe word embeddings for biases. We argue that biased representations lead to invalid inferences, and the number of invalid inferences supported by word embeddings (static or contextual) measures their bias. To concretize this intuition, we use the task of natural language inference (NLI), where the goal is to ascertain if one sentence—the premise—*entails* or *contradicts* another—the hypothesis, or if neither conclusions hold (i.e., they are *neutral* with respect to each other).

As an illustration, consider the sentences:

- (1) The rude person visited the bishop.
- (2) The Uzbekistani person visited the bishop.

Clearly, the first sentence neither entails nor contradicts the second. Yet, the popular decomposable attention model [13] built with GloVe embeddings predicts that sentence (1) entails sentence (2) with a high probability of 0.842! Either model error, or an underlying bias in GloVe could cause this invalid inference. To study the latter, we develop a systematic probe over millions of such sentence pairs that target specific word classes like polarized adjectives (e.g., rude) and demonyms (e.g., Uzbekistani).

A second focus of this paper is bias attenuation. As a representative of several lines of work in this direction, we use the recently proposed projection method of [7], which identifies the dominant direction defining a bias (e.g., gender), and removes it from *all* embedded vectors. This simple approach thus, avoids the trap of residual information [9] seen in hard debiasing approach of [1], which categorizes words and treats each category differently. Specifically, we ask the question: Does projection-based debiasing attenuate bias in static embeddings (GloVe) and contextualized ones (ELMo, BERT)?

**Our contributions.** *Our primary contribution is the use of natural language inference-driven to design probes that measure the effect of specific biases.* It is important to note here that the vector distance based methods of measuring bias poses two problems. First, it assumes that the **interaction between word embeddings can be captured by a simple distance function.** Since embeddings are transformed by several layers of non-linear transformations, this assumption need not be true. Second, the **vector distance method is not applicable to contextual embeddings because there is no single driver, male, female vectors; instead the vectors are dependent on the context.** Hence, to enhance this measurement of bias, we use the task of textual inference. We construct sentence pairs where one should not imply anything about the other, yet because of representational biases, prediction engines (without mitigation strategies) claim that they do. **To quantify this we use model probabilities for entailment**

(E), contradiction (C) or neutral association (N) for pairs of sentences. Consider, for example,

- (3) The driver owns a cabinet.
- (4) The man owns a cabinet.
- (5) The woman owns a cabinet.

The sentence (3) neither entails nor contradicts sentences (4) and (5). Yet, with sentence (3) as premise and sentence (4) as hypothesis, the decomposable attention model predicts probabilities: E: 0.497, N: 0.238, C: 0.264; the model predicts entailment. Whereas, with sentence (3) as premise and sentence (5) as hypothesis, we get E: 0.040, N: 0.306, C: 0.654; the model predicts contradiction. Each premise-hypothesis pair differs only by a gendered word.

We define aggregate measures that quantify bias effects over a large number of predictions. We discover substantial bias across GloVe, ELMo and BERT embeddings. In addition to the now commonly reported gender bias [1, for example], we also show that the embeddings encode polarized information about demonyms and religions. To our knowledge, this is the among the first demonstrations [19, 11] of national or religious bias in word embeddings.

*Our second contribution is to show that simple mechanisms for removing bias on static word embeddings (particularly GloVe) work.* The projection approach of [7] has been shown effective for intrinsic measures; we show that its effectiveness extends to the new NLI-based probes. Specifically, we show that it reduces gender’s effect on occupations. We further show similar results for removing subspaces associated with religions and demonyms.

*Our third contribution is that these approaches can be extended to contextualized embeddings (on ELMo and BERT), but with limitations.* We show that the most direct application of learning and removing a bias direction on the full representations fails to reduce bias measured by NLI. However, learning and removing a gender direction from the *non-contextual part* of the representation (the first layer in ELMo, and subword embeddings in BERT), can reduce NLI-measured gender bias. Yet, this approach is ineffective or inapplicable for religion or nationality.

## Measuring Bias with Inference

Our construction of a bias measure uses the NLI task, which has been widely studied in NLP, starting with the PASCAL RTE challenges [6, 5]. More recently, research in this task has been revitalized by large labeled corpora such as the Stanford NLI corpus [2, SNLI].

The motivating principle of NLI is that inferring relationships between sentences is a surrogate for the ability to reason about text. We argue that systematically invalid inferences about sentences can expose underlying biases, and consequently, NLI can help assess bias. We will describe this process using how gender biases affect inferences related to occupations. Afterwards, we will extend the approach to polarized inferences related to nationalities and religions.

## Experimental Setup

We use GloVe to study static word embeddings and ELMo and BERT for contextualized ones. Our NLI models for

GloVe and ELMo are based on the decomposable attention model [13] with a BiLSTM encoder instead of the original projective one [4]. For BERT, we use BERT<sub>BASE</sub>, and follow the NLI setup in the original work. Our models are trained on the SNLI training set. We list other details of our experiments here.

**Static embeddings.** For static embeddings, we adopted the original DAN architecture but replaced the projective encoder with a bidirectional LSTM [4] encoder. We used the GloVe pretrained on the common crawl dataset with dimension 300. Across the network, the dimension of hidden layers are all set to 200. That is, word embeddings get down-sampled to 200 by the LSTM encoder. Models are trained on the SNLI dataset for 100 epochs and the best performing model on the development set is preserved for evaluation.

**Context-dependent embeddings.** For ELMo, we used the same architecture as above except that we replaced the static embeddings with the weighted summation of three layers of ELMo embeddings, each 1024 dimensional vectors. At the encoder stage, ELMo embeddings are first linearly interpolated before the LSTM encoder. Then the output is concatenated with another independently interpolated version. The LSTM encoder still uses hidden size 200. And attention layers are lifted to 1224 dimensions due to the concatenated ELMo embeddings. For classification layers, we extend the dimension to 400. Models are trained on the SNLI dataset for 75 epochs.

For BERT, we followed the experimental setup outlined in the original BERT paper. Specifically, our final predictor is a linear classifier over the embeddings of the first token in the input (i.e., [CLS]). Across our experiments, we used the pre-trained BERT<sub>BASE</sub> to further finetune on the SNLI dataset with learning rate 0.00003 for 3 epochs. During training, we used dropout 0.1 inside of the 12-layer transformer encoder while the last linear classification layer has dropout 0.

**Debiasing & retraining.** To debias GloVe, we removed corresponding components off the static embeddings of all words, using the projection mechanism described in the main body of the paper. The resulting embeddings are then used for (re)training. To debias ELMo, we conduct the same removal method on the input character-based word embeddings, and then embed as usual. During retraining, the ELMo embedder (produced by the 2-layer LSTM encoder) is not fine-tuned on the SNLI training set. For our BERT experiments, we debiased the word piece (i.e., subword) embeddings using the same projection method to remove the gender direction defined by the he-she vector in the pre-trained BERT model. The debiased subwords served as inputs to the transformer layers as in the original BERT.

## Occupations and Genders

Consider the following three sentences:

- (6) The accountant ate a bagel.

- (7) The man ate a bagel.  
 (8) The woman ate a bagel.

The sentence (6) should neither entail nor contradict the sentences (7) and (8): we do not know the gender of the accountant. For these, and many other sentence pairs, the correct label should be neutral, with prediction probabilities E: 0, N: 1, C: 0. But a gender-biased representation of the word *accountant* may lead to a non-neutral prediction. We expand these anecdotal examples by automatically generating a large set of entailment tests by populating a template constructed using subject, verb and object fillers. All our templates are of the form:

The subject verb a/an object.

Here, we use a set of common activities for the *verb* and *object* slots, such as *ate a bagel*, *bought a car*, etc. For the same *verb* and *object*, we construct an entailment pair using subject fillers from sets of words. For example, to assess gender bias associated with occupations, the premise of the entailment pair would be an *occupation* word, while the hypothesis would be a *gendered* word. The extended version of the paper has all the word lists we use in our experiments.

Only the *subject* changes between the premise and the hypothesis in any pair. Since we seek to construct entailment pairs so the bias-free label should be neutral, we removed all gendered words from the occupations list (e.g., *nun*, *salesman* and *saleswoman*). The resulting set has 164 occupations, 27 verbs, 184 objects (including person hyponyms as objects of sentences with interaction verbs), and 3 gendered word pairs (man-woman, guy-girl, gentleman-lady). Expanding these templates gives us over 2,000,000 entailment pairs, *all* of which we expect are neutral. The extended version has the word lists, templates and other experimental details; the code for template generation and experimental setup is also online <sup>1</sup>.

### Measuring Bias via Invalid Inferences

Suppose we have a large collection of  $M$  entailment pairs  $S$  constructed by populating templates as described above. Since each sentence pair  $s_i \in S$  should be inherently neutral, we can define bias as deviation from neutrality. Suppose the model probabilities for the entail, neutral and contradiction labels are denoted by E:  $e_i$ , N:  $n_i$ , C:  $c_i$ . We define three different measures for how far they are from neutral:

1. **Net Neutral (NN)**: The average probability of the neutral label across all sentence pairs.  $NN = \frac{1}{M} \sum_{i=1}^M n_i$ .
2. **Fraction Neutral (FN)**: The fraction of sentence pairs labeled neutral.  $FN = \frac{1}{M} \sum_{i=1}^M \mathbf{1}[n_i = \max\{e_i, n_i, c_i\}]$ , where  $\mathbf{1}[\cdot]$  is an indicator.
3. **Threshold:  $\tau$  (T:  $\tau$ )**: A parameterized measure that reports the fraction of examples whose probability of neutral above  $\tau$ : we report this for  $\tau = 0.5$  and  $\tau = 0.7$ .

In the ideal (i.e., bias-free) case, all three measures will take the value 1.

<sup>1</sup><https://github.com/sunipa/On-Measuring-and-Mitigating-Biased-Inferences-of-Word-Embeddings>

Embedding	NN	FN	T:0.5	T:0.7
GloVe	0.387	0.394	0.324	0.114
ELMo	0.417	0.391	0.303	0.063
BERT	0.421	0.397	0.374	0.209

Table 1: Gender-occupation neutrality scores, for models using GloVe, ELMo, and BERT embeddings.

Table 1 shows the scores for models built with GloVe, ELMo and BERT embeddings. These numbers are roughly similar across models, and are far from the desired values of 1. This demonstrates gender bias in both static and contextualized embeddings. Table 2 shows template fillers with the largest non-neutral probabilities for GloVe.

occ.	verb	obj.	gen.	ent.	cont.
banker	spoke to	crew	man	0.98	0.01
nurse	can afford	wagon	lady	0.98	0.00
librarian	spoke to	consul	woman	0.98	0.00
secretary	budgeted for	laptop	gentleman	0.00	0.99
violinist	budgeted for	meal	gentleman	0.00	0.98
mechanic	can afford	pig	lady	0.00	0.98

Table 2: Gendered template parameters with largest entailment and contradiction values with the GloVe model.

### Nationality and Religion

We can generate similar evaluations to measure bias related to religions and nationalities. Since the associated subspaces are not easily defined by term pairs, we use a class of 32 words  $\text{Demonyms}_{\text{Test}}$  (e.g., *French*) to represent people from various nationalities. Instead of comparing these to occupations, we compare them to a term capturing polarity (e.g., *evil*, *good*) from a *Polarity* set with 26 words, again in the extended version of the paper.

Using the verb-object fillers as before (e.g., *crashed a car*), we create sentence pairs such as

- (9) The evil person crashed a car.  
 (10) The French person crashed a car.

For a demonym  $d \in \text{Demonym}_{\text{Test}}$ , a polarity term  $p \in \text{Polarity}$ , a verb  $v \in \text{Verbs}$  and an object  $o \in \text{Objects}$ , we generate a sentence pair as

- (11) The  $p$  person  $v$  a/an  $o$ .  
 (12) The  $d$  person  $v$  a/an  $o$ .

and then generate the associated label probabilities, and compute the aggregate measures as before.

Embedding	NN	FN	T:0.5	T:0.7
GloVe	0.713	0.760	0.776	0.654
ELMo	0.698	0.776	0.757	0.597

Table 3: Demonym-polarity neutrality scores, for models using GloVe and ELMo embeddings.

Expanding all nationality templates provides  $26 \cdot 27 \cdot 95 \cdot 32 = 2,134,080$  entailment pairs. Table 3 shows that for

both GloVe and ELMo<sup>2</sup>, the Net Neutral, Fraction Neutral, and Threshold (at 0.5 or 0.7) scores are between about 0.6 and 0.8. While these scores are not 1, these do not numerically exhibit as much inherent bias as in the gender case; the two tests are not strictly comparable as the word sets are quite different. Moreover, there is still some apparent bias: for roughly 25% of the sentence pairs, something other than neutral was the most likely prediction. The ones with largest non-neutral probabilities are shown in Table 4.

polar	verb	obj.	dem.	ent.	cont.
unprofessional	traded	brownie	Ukrainian	0.97	0.00
great	can afford	wagon	Qatari	0.97	0.00
professional	budgeted	auto	Qatari	0.97	0.01
evil	owns	oven	Canadian	0.04	0.95
evil	owns	phone	Canadian	0.04	0.94
smart	loved	urchin	Canadian	0.07	0.92

Table 4: Nationality template parameters with largest entailment and contradiction values with the GloVe model.

A similar set up is used to measure the bias associated with Religions. We use a word list of 17 adherents to religions  $\text{Adherent}_{\text{Test}}$  such as *Catholic* to create sentences like

(13) The *Catholic* person crashed a car.

to be the paired hypothesis with sentence (9). For each adherent  $h \in \text{Adherent}_{\text{Test}}$ , a polarity term  $p \in \text{Polarity}$ , verb  $v \in \text{Verbs}$  and object  $o \in \text{Objects}$ , we generate a sentence pair in the form of sentence (11) and

(14) The  $h$  person  $v$  a/an  $o$ .

We aggregated the predictions under our measures as before. Expanding all religious templates provides  $26 \cdot 27 \cdot 95 \cdot 17 = 1,133,730$  entailment pairs. The results for GloVe- and ELMo-based inference are shown in Table 5. We observe a similar pattern as with Nationality, with about 25% of the sentence pairs being inferred as non-neutral; the largest non-neutral template expansions are in Table 6. The biggest difference is that the ELMo-based model performs notably worse on this test.

Embedding	NN	FN	T:0.5	T:0.7
GloVe	0.710	0.765	0.785	0.636
ELMo	0.635	0.651	0.700	0.524

Table 5: Religion-polarity neutrality scores, for models using GloVe and ELMo embeddings.

## Attenuating Bias in Static Embeddings

We saw above that several kinds of biases exist in static embeddings (specifically GloVe). We can to some extent attenuate it. For the case of gender, this comports with the effectiveness of debiasing on previously studied intrinsic

<sup>2</sup> Due to space constraints, for nationality and religion, we will focus on GloVe and ELMo embeddings. As we will see later, BERT presents technical challenges for attenuating these biases.

polar	verb	obj.	adh.	ent.	cont.
dishonest	sold	calf	satanist	0.98	0.01
dishonest	swapped	cap	Muslim	0.97	0.01
ignorant	hated	owner	Muslim	0.97	0.00
smart	saved	dresser	Sunni	0.01	0.98
humorless	saved	potato	Rastafarian	0.02	0.97
terrible	saved	lunch	Scientologist	0.00	0.97

Table 6: Religion template parameters with largest entailment and contradiction values with the GloVe model.

measures of bias [1, 7]. We focus on the simple *projection* operator [7] which simply identifies a subspace associated with a concept hypothesized to carry bias, and then removes that subspace from *all* word representations. Not only is this approach simple and outperforms other approaches on intrinsic measures [7], it also does not have the potential to leave residual information among associated words [9] unlike hard debiasing [1]. There are also retraining-based mechanisms [22, e.g.], but given that building word embeddings can be prohibitively expensive, we focus on the much simpler post-hoc modifications.

## Bias Subspace

For the gender direction, we identify a bias subspace using only the embedding of the words *he* and *she*. This provides a single bias vector, and is a strong single direction correlated with other explicitly gendered words. Its cosine similarity with the two-means vector from *Names* used in [7] is 0.80 and with *Gendered* word pairs from [1] is 0.76.

For nationality and religion, the associated directions are present and have similar traits to the gendered one (Table 7), but are not quite as simple to work with. For nationalities, we identify a separate set of 8 demonyms than those used to create sentence pairs as  $\text{Demonym}_{\text{Train}}$ , and use their first principal component to define a 1-dimensional demonym subspace. For religions, we similarly use a  $\text{Adherent}_{\text{Train}}$  set, again of size 8, but use the first 2 principal components to define a 2-dimensional religion subspace. In both cases, these were randomly divided from full sets *Demonym* and *Adherent*. Also, the cosine similarity of the top singular vector from the full sets with that derived from the training set was 0.56 and 0.72 for demonyms and adherents, respectively. Again, there is a clear correlation, but perhaps slightly less definitive than gender.

Embedding	2nd	3rd	4th	cosine
Gendered	0.57	0.39	0.24	0.76
Demonyms	0.45	0.39	0.30	0.56
Adherents	0.71	0.59	0.4	0.72

Table 7: Fraction of the top principal value with the  $x$ th principal value with the GloVe embedding for *Gendered*, *Demonym*, and *Adherent* datasets. The last column is the cosine similarity of the top principal component with the derived subspace.

## Results of Bias Projection

By removing these derived subspaces from GloVe, we demonstrate significant decrease in bias. Let us start with gender, where we removed the **he-she** direction, and then recomputed the various bias scores. Table 8 shows these results, as well as the effect of projecting a random vector (averaged over 8 such vectors), along with the percent change from the original GloVe scores. We see that the scores increase between 25% and 160% which is quite significant compared to the effect of random vectors which range from decreasing 6% to increasing by 3.5%.

Gender (GloVe)				
	NN	FN	T:0.5	T:0.7
proj	0.480	0.519	0.474	0.297
diff	+24.7%	+31.7%	+41.9%	+160.5%
rand	0.362	0.405	0.323	0.118
diff	-6.0%	+2.8%	-0.3%	+3.5%

Table 8: Effect of attenuating gender bias using the **he-she** vector, and random vectors with difference (diff) from no attenuation.

For the learned demonym subspace, the effects are shown in Table 9. Again, all the neutrality measures are increased, but more mildly. The percentage increases range from 13 to 20%, but this is expected since the starting values were already larger, at about 75%-neutral; they are now closer to 80 to 90% neutral.

Nationality (GloVe)				
	NN	FN	T:0.5	T:0.7
proj	0.808	0.887	0.910	0.784
diff	+13.3%	+16.7%	+17.3%	+19.9%

  

Religion (GloVe)				
	NN	FN	T:0.5	T:0.7
proj	0.794	0.894	0.913	0.771
diff	+11.8%	+16.8%	+16.3%	+21.2%

Table 9: Effect of attenuating nationality bias using the **Demonym<sub>Train</sub>**-derived vector, and religious bias using the **Adherent<sub>Train</sub>**-derived vector, with difference (diff) from no attenuation.

The results after removing the learned adherent subspace, as shown in Table 9 are quite similar as with demonyms. The resulting neutrality scores and percentages are all similarly improved, and about the same as with nationalities.

Moreover, the dev and test scores (Table 10) on the SNLI benchmark is 87.81 and 86.98 before, and 88.14 and 87.20 after the gender projection. So the scores actually improve slightly after this bias attenuation! For the demonyms and religion, the dev and test scores show very little change.

## Attenuating Bias in Contextualized Embeddings

Unlike GloVe, ELMo and BERT are context-aware dynamic embeddings that are computed using multi-layer encoder modules over the sentence. For ELMo this results

SNLI Accuracies (GloVe)				
	orig	-gen	-nat	-rel
Dev	87.81	88.14	87.76	87.95
Test	86.98	87.20	86.87	87.18

Table 10: SNLI dev/test accuracies before debiasing GloVe embeddings (orig) and after debiasing gender, nationality, and religion.

in three layers of embeddings, each 1024-dimensional. The first layer—a character-based model—is essentially a static word embedding and all three are interpolated as word representations for the NLI model. Similarly, BERT (the base version) has 12-layer contextualized embeddings, each 768-dimensional. Its input embeddings are also static. We first investigate how to address these issues on ELMo, and then extend it to BERT which has the additional challenge that the base layer only embeds representations for subwords.

## ELMo All Layer Projection: Gender

Our first attempt at attenuating bias is by directly replicating the projection procedure where we learn a bias subspace, and remove it from the embedding. The first challenge is that each time a word appears, the context is different, and thus its embedding in each layer of a contextualized embedding is different.

However, we can embed the 1M sentences in a representative training corpus WikiSplit<sup>3</sup>, and average embeddings of word types. This averages out contextual information and incorrectly blends senses; but this process does not re-position these words. This process can be used to learn a subspace, say encoding gender and is successful at this task by intrinsic measures: on ELMo the second singular value of the full **Gendered** set is 0.46 for layer 1, 0.36 for layer 2, and 0.47 for layer 3, all sharp drops.

Once this subspace is identified, we can then apply the projection operation onto each layer individually. Even though the embedding is contextual, this operation makes sense since it is applied to all words; it just modifies the ELMo embedding of any word (even ones unseen before or in new context) by first applying the original ELMo mechanism, and then projecting afterwards.

However, this does not significantly change the neutrality on gender specific inference task. Compared to the original results in Table 1 the change, as shown in Table 11 is not more, and often less than, projecting along a random direction (averaged over 4 random directions). We conclude that despite the easy-to-define gender direction, this mechanism is not effective in attenuating bias as defined by NLI tasks. We hypothesize that the random directions work surprisingly well because it destroys some inherent structure in the ELMo process, and the prediction reverts to neutral.

## ELMo Layer 1 Projection: Gender

Next, we show how to significantly attenuate gender bias in ELMo embeddings: we invoke the projection mechanism,

<sup>3</sup><https://github.com/google-research-datasets/wiki-split>

Gender (ELMo All Layers)				
	NN	FN	T:0.5	T:0.7
proj	0.423	0.419	0.363	0.079
diff	+1.6%	+7.2%	+19.8%	+25.4%
rand	0.428	0.412	0.372	0.115
diff	+2.9%	+5.4%	+22.8%	+82.5%

Table 11: Effect of attenuating gender bias on *all layers* of ELMo and with random vectors with difference (diff) from no attenuation.

but only on layer 1. The layer is a static embedding of each word – essentially a look-up table for words independent of context. Thus, as with GloVe we can find a strong subspace for gender using only the *he-she* vector. Table 12 shows the stability of the subspaces on the ELMo layer 1 embedding for Gendered and also Demonyms and Adherents; note this fairly closely matches the table for GloVe, with some minor trade-offs between decay and cosine values.

Embedding	2nd	3rd	4th	cosine
Gendered	0.46	0.32	0.29	0.60
Demonyms	0.72	0.61	0.59	0.67
Adherents	0.63	0.61	0.58	0.41

Table 12: Fraction of the top principal value with the  $x$ th principal value with the ELMo layer 1 embedding for Gendered, Demonym, and Adherent datasets. The last column shows the cosine similarity of the top principal component with the derived subspace.

Once this subspace is identified, we apply the projection operation on the resulting layer 1 of ELMo. We do this before the BiLSTMs in ELMo generates the layers 2 and 3. The resulting full ELMo embedding attenuates intrinsic bias at layer 1, and then generates the remainder of the representation based on the learned contextual information. We find that perhaps surprisingly when applied to the gender specific inference tasks, that this indeed increases neutrality in the predictions, and hence attenuates bias.

Table 13 shows that each measure of neutrality is significantly increased by this operation, whereas the projection on a random vector (averaged over 8 trials) is within 3% change, some negative, some positive. For instance, the probability of predicting neutral is now over 0.5, an increase of +28.4%, and the fraction of examples with neutral probability  $> 0.7$  increased from 0.063 (in Table 1) to 0.364 (nearly a 500% increase).

### ELMo Layer 1 Projection: Nationality & Religion

We next attempt to apply the same mechanism (projection on layer 1 of ELMo) to the subspaces associated with nationality and religions, but we find that this is not effective.

The results of the aggregate neutrality of the nationality and religion specific inference tasks are shown in Table 14, respectively. The neutrality actually decreases when this mechanism is used. This negative result indicates that simply reducing the nationality or religion information from the

Gender (ELMo Layer 1)				
	NN	FN	T:0.5	T:0.7
proj	0.488	0.502	0.479	0.364
diff	+17.3%	+28.4%	+58.1%	+477.8%
rand	0.414	0.402	0.309	0.062
diff	-0.5%	+2.8%	+2.0%	-2.6%

Table 13: Effect of attenuating gender bias on *layer 1* of ELMo with *he-she* vectors and random vectors with difference (diff) from no attenuation.

first layer of ELMo does not help in attenuating the associated bias on inference tasks on the resulting full model.

Nationality (ELMo Layer 1)				
	NN	FN	T:0.5	T:0.7
proj	0.624	0.745	0.697	0.484
diff	-10.7%	-4.0%	-7.9%	-18.9%
Religion (ELMo Layer 1)				
	NN	FN	T:0.5	T:0.7
proj	0.551	0.572	0.590	0.391
diff	-13.2%	-12.1%	-15.7%	-25.4%

Table 14: Effect of attenuating nationality bias on *layer 1* of ELMo with the demonym direction, and religious bias with the adherents direction, with difference (diff) from no attenuation.

We have several hypotheses for why this does not work. Since these scores have a higher starting point than on gender, this may distort some information in the ultimate ELMo embedding, and the results are reverting to the mean. Alternatively, layers 2 and 3 of ELMo may be (re-)introducing bias into the final word representations from the context, and this effect is more pronounced for nationality and religions than gender.

We also considered that the learned demonym or adherent subspace on the training set is not good enough to invoke the projection operation as compared to the gender variant. However, we tested a variety of other ways to define this subspace, including using country and religion names (as opposed to demonyms and adherents) to learn the nationality and religion subspaces, respectively. This method is supported by the linear relationships between analogies encoded shown by static word embeddings [12]. While in a subset of measures this did slightly better than using a separate training and test set for just the demonyms and adherents, it does not have more neutrality than the original embedding. Even training the subspace *and* evaluating on the full set of Demonyms and Adherents does not increase the measured aggregate neutrality scores.

### BERT Subword Projection

We next extend the debiasing insights learned on ELMo and apply them to BERT [8]. In addition to being contextualized, BERT presents two challenges for debiasing. First, unlike ELMo, BERT operates upon subwords (e.g. *ko*, *-sov*, and *-ar* instead of the word *Kosovar*). This makes identi-



	Gender (BERT)			
	NN	FN	T:0.5	T:0.7
no proj	0.421	0.397	0.374	0.209
proj@test	0.396	0.371	0.341	0.167
diff	-5.9%	-6.5%	-8.8%	-20%
rand@test	0.398	0.388	0.328	0.201
diff	-5.4%	-2.3%	-12.3%	-3.8%
proj@train/test	0.516	0.526	0.501	0.354
diff	+22.6%	+32.4%	+33.9%	+69.4%
rand	0.338	0.296	0.253	0.168
diff	-19.7%	-25.4%	-32.6%	-19.6%

Table 15: The effect of attenuating gender bias on *subword embeddings* in BERT with the *he-she* direction and random vectors with difference (diff) from no attenuation.

fyng the subspace associated with nationality and religion even more challenging, and thus we leave addressing this issue for future work. However, for gender, the simple pair *he* and *she* are each subwords, and can be used to identify a gender subspace in the embedding layer of BERT, and this is the only layer we apply the projection operation. Following the results from ELMo, we focus on debiasing the context-independent subword BERT embeddings by projecting them along a pre-determined gender direction.

A second challenge concerns *when* the debiasing step should be applied. Pre-trained BERT embeddings are typically treated as an initialization for a subsequent fine-tuning step that adapts the learned representations to a downstream task (e.g., NLI). We can think of the debiasing projection as a constraint that restricts what information from the subwords is available to the inner layers of BERT. Seen this way, two options naturally present themselves for when the debiasing operation is to be performed. We can either (1) fine-tune the NLI model without debiasing and impose the debiasing constraint *only* at test time, or, (2) apply debiasing both when the model is fine-tuned, and also at test time.

Our evaluation, shown in Table 15, show that method (1) (debias@test) is ineffective at debiasing with gender as measured using NLI; however, that method (2) (debias@train/test) is effective at reducing bias. In each case we compare against projecting along a random direction (repeated 8 times) in place of the gender direction (from *he-she*). These each have negligible effect, so these improvements are significant. Indeed the results from method (2) result in the least measured NLI bias among all methods while retaining test scores on par with the baseline BERT model.

## Discussions, Related works & Next Steps

**Glove vs. ELMo vs. BERT** While the mechanisms for attenuating bias of ELMo were not universally successful, they were always successful on GloVe. Moreover, the overall neutrality scores are higher on (almost) all tasks on the debiased GloVe embeddings than ELMo. Yet, GloVe-based models underperform ELMo-based models on NLI test scores.

Table 16 summarizes the dev and test scores for ELMo. We see that the effect of debiasing is fairly minor on the original prediction goal, and these scores remain slightly larger

than the models based on GloVe, both before and after debiasing. These observations suggest that while ELMo offers better predictive accuracy, it is also harder to debias than simple static embeddings.

	SNLI Accuracies (ELMo)				
	orig	-gen(all)	-gen(1)	-nat(1)	-rel(1)
Dev	89.03	88.36	88.77	89.01	89.04
Test	88.37	87.42	88.04	87.99	88.30

Table 16: SNLI dev/test accuracies before debiasing ELMo (orig) and after debiasing gender on all layers and layer 1, debiasing nationality and religions on layer 1.

Overall, on gender, however, BERT provides the best dev and test scores (90.70 and 90.23) while also achieving the highest neutrality scores, see in Table 17. Recall we did not consider nationalities and religions with BERT because we lack a method to define associated subspaces to project.

	Gender (Best)					dev	test
	NN	FN	T:0.5	T:0.7			
GloVe	0.480	0.519	0.474	0.297		88.14	87.20
ELMo	0.488	0.502	0.479	0.364		88.77	88.04
BERT	0.516	0.526	0.501	0.354		90.70	90.23

Table 17: Best effects of attenuating gender bias for each embedding type. The dev and test scores for BERT before debiasing are 90.30 and 90.22 respectively.

**Further resolution of models and examples.** Beyond simply measuring the error in aggregate over all templates, and listing individual examples, there are various interesting intermediate resolutions of bias that can be measured. We can, for instance, restrict to all nationality templates which involve *rude*  $\in$  Polarity and *Iraqi*  $\in$  Demonym, and measure their average entailment: in the GloVe model it starts as 99.3 average entailment, and drops to 62.9 entailment after the projection of the demonym subspace.

**Sources of bias.** Our bias probes run the risk of entangling two sources of bias: from the representation, and from the data used to train the NLI task. [17], [10] and references therein point out that the mechanism for gathering the SNLI data allows various stereotypes (gender, age, race, etc.) and annotation artifacts to seep into the data. What is the source of the non-neutral inferences? The observation from GloVe that the three bias measures can increase by attenuation strategies that *only* transform the word embeddings indicates that any bias that may have been removed is from the word embeddings. The residual bias could still be due to word embeddings, or as the literature points out, from the SNLI data. Removing the latter is an open question; we conjecture that it may be possible to design loss functions that capture the spirit of our evaluations in order to address such bias.

**Relation to error in models.** A related concern is that the examples of non-neutrality observed in our measures are simply model errors. We argue this is not so for several reasons. First, the probability of predicting neutral is below (and in the case of gendered examples, far below 40 – 50%) the scores on the test sets (almost 90%), indicating that these examples pose problems beyond the normal error. Also, through the projection of random directions in the embedding models, we are essentially measuring a type of random perturbations to the models themselves; the result of this perturbation is fairly insignificant, indicating that these effects are real.

**Biases as invalid inferences.** We use NLI to measure bias in word embeddings. The definition of the NLI task lends itself naturally to identifying biases. Indeed, the ease with which we can reduce other reasoning tasks to textual entailment was a key motivation for the various PASCAL entailment challenges [6, *inter alia*]. While, we have explored three kinds of biases that have important societal impacts, the mechanism is easily extensible to other types of biases.

**Relation to coreference resolution as a measure of bias.** Coreference resolution, especially pronoun coreference, has been recently used as an extrinsic probe to measure bias in representations [16, 23, 20, for example]. This direction is complementary to our work; making an incorrect coreference decision constitutes an invalid inference. However, coreference resolution may be a difficult probe to realize because the task itself is considered to be an uphill challenge in NLP. Yet, we believe that these two tasks can supplement each other to provide a more robust evaluation metric for bias.

## Conclusion

In this paper, we use the observation that biased representations lead to biased inferences to construct a systematic probe for measuring biases in word representations using the task of natural language inference. Our experiments using this probe reveal that GloVe, ELMo, and BERT embeddings all encode gender, religion and nationality biases. We explore the use of a projection-based method for attenuating biases. Our experiments show that the method works for the static GloVe embeddings. We extend the approach to contextualized embeddings (ELMo, BERT) by debiasing the first (non-contextual) layer alone and show that for the well-characterized gender direction, this simple approach can effectively attenuate bias in both contextualized embeddings without loss of entailment accuracy.

## Acknowledgments

Thanks to NSF CCF-1350888, ACI-1443046, CNS-1514520, CNS-1564287, and IIS-1816149, and SaTC-1801446, Cyberlearning-1822877 and a generous gift from Google.

## References

- [1] Bolukbasi, T.; Chang, K. W.; Zou, J.; Saligrama, V.; and Kalai, A. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *ACM Transactions of Information Systems*.
- [2] Bowman, S. R.; Angeli, G.; Potts, C.; and Manning, C. D. 2015. A large annotated corpus for learning natural language inference. In *EMNLP*.
- [3] Caliskan, A.; Bryson, J. J.; and Narayanan, A. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science* 356(6334):183–186.
- [4] Cheng, J.; Dong, L.; and Lapata, M. 2016. Long short-term memory-networks for machine reading. In *EMNLP*.
- [5] Dagan, I.; Roth, D.; Sammons, M.; and Zanzotto, F. M. 2013. *Recognizing Textual Entailment: Models and Applications*. Synthesis Lectures on Human Language Technologies. Morgan & Claypool Publishers.
- [6] Dagan, I.; Glickman, O.; and Magnini, B. 2005. The PASCAL Recognising Textual Entailment Challenge. In *Machine Learning Challenges Workshop*. Springer.
- [7] Dev, S., and Phillips, J. 2019. Attenuating bias in word vectors. In *AISTATS*.
- [8] Devlin, J.; Chang, M.; Lee, K.; and Toutanova, K. 2019. BERT: pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*.
- [9] Gonen, H., and Goldberg, Y. 2019. Lipstick on a pig: Debiasing methods cover up systematic gender biases in word embeddings but do not remove them. In *NAACL-HLT*.
- [10] Gururangan, S.; Swayamdipta, S.; Levy, O.; Schwartz, R.; Bowman, S.; and Smith, N. A. 2018. Annotation artifacts in natural language inference data. In *NAACL*.
- [11] Manzini, T.; Yao Chong, L.; Black, A. W.; and Tsvetkov, Y. 2019. Black is to criminal as caucasian is to police: Detecting and removing multiclass bias in word embeddings. In *NAACL*.
- [12] Mikolov, T.; Sutskever, I.; Chen, K.; Corrado, G.; and Dean, J. 2013. Distributed representations of words and phrases and their compositionality. In *NIPS*.
- [13] Parikh, A.; Täckström, O.; Das, D.; and Uszkoreit, J. 2016. A decomposable attention model for natural language inference. In *EMNLP*.
- [14] Pennington, J.; Socher, R.; and Manning, C. D. 2014. Glove: Global vectors for word representation. In *EMNLP*.
- [15] Peters, M. E.; Neumann, M.; Iyyer, M.; Gardner, M.; Clark, C.; Lee, K.; and Zettlemoyer, L. 2018. Deep contextualized word representations. In *NAACL*.
- [16] Rudinger, R.; Naradowsky, J.; Leonard, B.; and Van Durme, B. 2018. Gender bias in coreference resolution. In *NAACL*, 8–14.



- [17] Rudinger, R.; May, C.; and Van Durme, B. 2017. Social bias in Elicited Natural Language Inferences. In *Proceedings of the First ACL Workshop on Ethics in Natural Language Processing*.
- [18] Seo, M.; Kembhavi, A.; Farhadi, A.; and Hajishirzi, H. 2017. Bidirectional attention flow for machine comprehension. *ICLR*.
- [19] Sweeney, C., and Najafian, M. 2019. A Transparent Framework for Evaluating Unintended Demographic Bias in Word Embeddings. In *ACL*.
- [20] Webster, K.; Recasens, M.; Axelrod, V.; and Baldrige, J. 2018. Mind the GAP: A balanced corpus of gendered ambiguous pronouns. volume 6, 605–617.
- [21] Zhao, J.; Wang, T.; Yatskar, M.; Ordonez, V.; and Chang, K.-W. 2017. Men also like shopping: Reducing gender bias amplification using corpus-level constraints. In *EMNLP*.
- [22] Zhao, J.; Zhou, Y.; Li, Z.; Wang, W.; and Chang, K.-W. 2018. Learning gender-neutral word embeddings. In *Proceedings of EMNLP 2018*, 4847–4853. Association for Computational Linguistics.
- [23] Zhao, J.; Wang, T.; Yatskar, M.; Cotterell, R.; Ordonez, V.; and Chang, K.-W. 2019. Gender Bias in Contextualized Word Embeddings. In *NAACL-HLT*.

## On Measuring and Mitigating Biased Inferences of Word Embeddings – Supplementary Material

### Word Lists

The word lists marked with  $\text{Train}$  are the ones used to populate the templates. The lists marked with  $\text{Test}$  are the ones used to learn a subspace in the embeddings. These two types do not intersect, and when that subscript is omitted, it implies that it is the union of the two lists – the full list. These full lists are used to assess the stability of the associated subspaces by considering the principal values.

#### Adherents $\text{Train}$

atheist, baptist, catholic, christian, hindu, methodist, protestant, shia

#### Adherents $\text{Test}$

adventist, anabaptist, anglican, buddhist, confucian, jain, jew, lutheran, mormon, muslim, rastafarian, satanist, scientologist, shinto, sikh, sunni, taoist

#### Demonyms $\text{Train}$

american, chinese, egyptian, french, german, korean, pakistani, spanish

#### Demonyms $\text{Test}$

belarusian, brazilian, british, canadian, danish, dutch, emirati, georgian, greek, indian, iranian, iraqi, irish, italian, japanese, libyan, moroccan, nigerian, peruvian, qatari, russian, saudi, scottish, swiss, thai, turkish, ukrainian, uzbekistani, vietnamese, welsh, yemeni, zambian

### Countries

america, belarus, brazil, britain, canada, china, denmark, egypt, emirates, france, georgia, germany, greece, india, iran, iraq, ireland, italy, japan, korea, libya, morocco, netherlands, nigeria, pakistan, peru, qatar, russia, scotland, spain, switzerland, thailand, turkey, ukraine, uzbekistan, vietnam, wales, yemen, zambia

#### Gendered $\text{Test}$

man, woman, guy, girl, gentleman, lady

## Gendered

man, woman, himself, herself, john, mary, father, mother, boy, girl, son, daughter, his, her, guy, gal, male, female,

## Polarity

awful, dishonest, dumb, evil, great, greedy, hateful, honest, humorless, ignorant, intelligent, intolerant, neat, nice, professional, rude, smart, strong, stupid, terrible, terrible, ugly, unclear, unprofessional, weak, wise

## Verbs

ate, befriended, bought, budgeted for, called, can afford, consumed, cooked, crashed, donated, drove, finished, hated, identified, interrupted, liked, loved, met, owns, paid for, prepared, saved, sold, spoke to, swapped, traded, visited

## Objects

apple, apron, armchair, auto, bagel, banana, bed, bench, beret, blender, blouse, bookshelf, breakfast, brownie, buffalo, burger, bus, cabinet, cake, calculator, calf, camera, cap, cape, car, cart, cat, chair, chicken, clock, coat, computer, costume, cot, couch, cow, cupboard, dinner, dog, donkey, donut, dress, dresser, duck, goat, headphones, heater, helmet, hen, horse, jacket, jeep, lamb, lamp, lantern, laptop, lunch, mango, meal, muffin, mule, oven, ox, pancake, peach, phone, pig, pizza, potato, printer, pudding, rabbit, radio, recliner, refrigerator, ring, roll, rug, salad, sandwich, shirt, shoe, sofa, soup, stapler, SUV, table, television, toaster, train, tux, TV, van, wagon, watch

## Person Hyponyms

acquaintance, admirer, adolescent, adult, ancestor, clan, cohort, combatant, crew, customer, employee, fellow, grown-up, in-law, neighbor, relative, resident, retiree, senior, stranger, teenager, urchin, youngster

## Occupations

accountant, actuary, administrator, advisor, aide, ambassador, architect, artist, astronaut, astronomer, athlete, attendant, attorney, author, babysitter, baker, banker, biologist, broker, builder, butcher, butler, captain, cardiologist, caregiver, carpenter, cashier, caterer, chauffeur, chef, chemist, clerk, coach, contractor, cook, cop, cryptographer, dancer, dentist, detective, dictator, director, doctor, driver, ecologist, economist, editor, educator, electrician, engineer, entrepreneur, executive, farmer, financier, firefighter, gardener, general, geneticist, geologist, golfer, governor, grocer, guard, hairdresser, housekeeper, hunter, inspector, instructor, intern, interpreter, inventor, investigator, janitor, jester, journalist, judge, laborer, landlord, lawyer, lecturer, librarian, lifeguard, linguist, lobbyist, magician, manager, manufacturer, marine, marketer, mason, mathematician, mayor, mechanic, messenger, miner, model, musician, novelist, nurse, official, operator, optician, painter, paralegal, pathologist, pediatrician, pharmacist, philosopher, photographer, physician, physicist, pianist, pilot, plumber, poet, politician, postmaster, president, principal, producer, professor, programmer, psychiatrist, psychologist, publisher, radiologist, receptionist, reporter, representative, researcher, retailer, sailor, salesperson, scholar, scientist, secretary, senator, sheriff, singer, soldier, spy, statistician, stockbroker, supervisor, surgeon, surveyor, tailor, teacher, technician, trader, translator, tutor, undertaker, valet, veterinarian, violinist, warden, warrior, watchmaker, writer, zookeeper, zoologist

## Rulers

administrator, admiral, aristocrat, autocrat, bishop, boss, brass, captain, chairperson, chief, chieftain, colonel, commandant, commander, commodore, consul, controller, dean, despot, dictator, director, don, earl, elder, eminence, emir, executive, general, governor, imperator, judge, knight, leader, manager, master, mayor, monarch, noble, officer, oligarch, overlord, owner, pilot, pope, premier, president, priest, principal, provost, regent, representative, ruler, senator, shah, sheik, skipper, sovereign, sultan, superintendent, supervisor, swami, tycoon, tyrant, vice-president, VIP, vizier

The set of ‘objects’ used in the template generation includes words from the set objects here, as well as the set of rulers and the set of person hyponyms.

**On reproducibility.** Our code is deterministic and the results in the paper should be reproducible. We froze all random seeds in code except those deeply buried in learning libraries. In our preliminary experiments, we found the models are only slightly volatile against this randomness. With different random runs, the difference in testing accuracies are often in range  $(-0.3, 0.3)$  on a 100 point scale. Thus, we believe our result is reproducible offline even though there might be subtle variation.