

Unmasking the Trade-off: Measuring Gender Bias Mitigation and Over-debiasing Effects in Pretrained Language Models

Anonymous ACL submission

Abstract

Pretrained language models (PLMs) have demonstrated success across many natural language processing tasks. However, they have been shown to encode gender bias present in the corpora they are trained on. Existing bias mitigation methods are usually devised to remove *all* associations related to gender. This can hurt the performance of PLMs, because of a possible loss of typical associations (e.g., not associating the word “mother” with female). To measure the extent of loss of typical gender associations (i.e. over-debiasing), we introduce the Typical Associations evaluation corpus for Gender (TA-Gender). We find that three popular debiasing methods result in substantial loss of typical gender associations. Our results highlight the importance of mitigating bias without removing typical gender associations, and our dataset constitutes the first benchmark to evaluate information loss.¹

1 Introduction

In recent years, pretrained language models (PLMs) (Devlin et al., 2019; Radford et al., 2019; Lewis et al., 2020) trained on large-scale corpora have become the de-facto backbone of modern NLP systems. These models are trained on minimally filtered real world text which reflects social biases of the real world (Sun et al., 2019; Bender et al., 2021), which in turn are encoded in PLMs and can propagate into downstream applications (Bolukbasi et al., 2016; Caliskan et al., 2017; Kiritchenko and Mohammad, 2018; May et al., 2019; Kurita et al., 2019). Considering the wide use of PLMs, the propagation of social bias in these models poses a danger of reinforcing stereotypes (Sun et al., 2019; Bender et al., 2021).

A number of methods have been introduced to remove social bias from PLMs (Zhao et al., 2017; Lu

et al., 2020; Zmigrod et al., 2019; Hall Maudslay et al., 2019; Liang et al., 2020; Huang et al., 2020), which are designed to bleach *all* associations with the debiasing target (e.g., gender) from PLMs in a non-discriminate fashion. A loss of typical associations (e.g., not associating “dad” with male) can result in over-debiasing and hurt performance of PLMs and their utility in downstream tasks such as question answering or common sense reasoning.

This paper presents a systematic analysis of gender *debiasing* and *over-debiasing* across a diverse range of PLMs using established debiasing methods. We conduct a suite of experiments to compare a number of bias evaluation metrics and introduce new metrics to measure the effects of over-debiasing in English. Specifically, we introduce the Typical Associations evaluation corpus for Gender (**TA-Gender**), a template-based challenge dataset that measures the degree of typical associations present in PLMs (Figure 1, top).

The boundary between *typical* and *inappropriate* is fluid, and we highlight weaknesses and put forward usage recommendations of our approach: (1) we treat gender as binary, as common in related work (Manzini et al., 2019; Bartl et al., 2020), but this by no means reflect the complexity of the concept (Cao and Daumé III, 2020); (2) our dataset encodes **typical, not factual, associations**: we acknowledge the fluidity of the concept of gender, and the fact that concepts like “pregnancy” and “fatherhood” are not fully restricted to one gender;² (3) our method cannot account for figurative language³ or cultural differences in language use across communities. Further discussion and recommendation on how to use our benchmark can be found in the Ethical Considerations (Section 9).

Our experiments test three widely used debiasing methods — counterfactual data substitution (Hall Maudslay et al., 2019), embedding regular-

¹Data and code are available at www.removed_for_anonymity.com

²E.g., trans men can and do become pregnant.

³E.g., “Dad is such a drama queen!”

TA-Gender	TA1	My ____ is one of the saleswomen. daughter (✓) son (✗)
	TA2	My ____ is a danseur. uncle (✓) aunt (✗)
StereoSet	SS1	Girls tend to be more ____ than boys. determined (✓) soft (✓) fish (✗)
	SS2	Every male acts ____. soft (✓) dumb (✓) target (✗)

Figure 1: Two examples each from our TA-Gender data set (top) vs. the StereoSet benchmark (bottom).

ization (Huang et al., 2020) and sentence debias (Liang et al., 2020) — to debias a number PLMs: BERT (Devlin et al., 2019), GPT2 (Radford et al., 2019) and BART (Lewis et al., 2020). Our results confirm prior results that all debiasing methods are effective in removing biases, but they also reveal that these methods lead to a substantial loss of typical gender associations. Lastly, we explore the trade-off between debiasing and over-debiasing. While our work focuses on gender as just one dimension of social biases, the presented framework can be extended to additional bias dimension such as racial or religious bias (exercising caution when defining *typicality*, as discussed above for the gender case). In sum, our contributions are

- A systematic evaluation of debiasing methods and bias metrics over a representative selection of PLMs of varying sizes.
- A benchmark data set, TA-Gender, designed to quantify the extent to which debiasing incurs a loss of typical gender associations.
- An analysis of the extent of over-debiasing across four debiasing methods and six PLMs.

2 Related Work

Bias Evaluation Caliskan et al. (2017) propose the Word Embedding Association Test (WEAT) to measure biases in word embeddings through the strength of association between target words (e.g., gender pronouns) and attribute words (e.g., gender neutral occupations). An unbiased model should exhibit no difference between the associations of attribute words with target words of different gender. May et al. (2019) extended this to biases in pretrained *contextualized* language models through the Sentence Encoder Association Test (SEAT), by encoding Caliskan et al. (2017)’s WEAT terms in

simple sentences, and measuring the associative strength of target sentences and attribute sentences as the cosine distances between their sentence embeddings. Focusing on masked language models, Kurita et al. (2019) propose logprob-score to evaluate bias in BERT. Instead of using cosine distances between embeddings, the association between target and attribute words is estimated by the probability of masked token predictions. We use both SEAT and logprob-score to evaluate bias in this work.

An alternative to template-based bias evaluation methods are crowdsourced datasets that capture societal notions of stereotypes across domains including gender, race or religion (Nadeem et al., 2020; Nangia et al., 2020). We consider the gender portion of StereoSet (Nadeem et al., 2020) in this work, which consists of crowd-sourced tuples of sentences capturing a stereotyped, an anti-stereotyped, and a semantically meaningless version of the same context (Figure 1, bottom). As such, StereoSet can be leveraged to evaluate bias (leveraging anti-stereotyped vs stereotyped alternatives) as well as the basic language modeling abilities of pretrained language models (leveraging the semantically meaningful vs meaningless alternatives). While StereoSet can provide an approximate measure of the extent of over-debiasing, the test is designed to probe stereotypical associations and general language modeling abilities. In contrast, our proposed dataset (TA-Gender), directly focuses on typical (i.e., (near) factual) information by using templates that directly probe for typical associations. Figure 1 illustrates the differences between TA-Gender and StereoSet. We compare the StereoSet LM-score against our proposed benchmark in our over-debiasing study (§ 7.2).

Bias Mitigation Several debiasing methods have been proposed in previous work. One line of work focuses on modifying training data to remove bias. Lu et al. (2020) propose counterfactual data augmentation to remove gender bias. For each sentence, it generates a counterfactual sentence which is identical except that all gender words are swapped to terms of the opposite gender. Then the union of original sentences and swapped sentences are used for training. Name-based counterfactual data substitution (CDS; Hall Maudslay et al. (2019)) is a variant which swaps the gender words including first names *in-place* for a random sub-set of sentences to counterbalance bias.

Another line of work introduces additional bias-

Template	Sent A (correct)	Sent B (incorrect)
Target is adj_attr .	<u>He</u> is male .	<u>She</u> is male .
Target is single_attr .	This man is a danseur .	This woman is a danseur .
As single_attr , <u>Target</u> deserves a wonderful life.	As a dad , <u>my husband</u> deserves a wonderful life.	As a dad , <u>my wife</u> deserves a wonderful life.
As one of the plural_attr , <u>Target</u> deserves a wonderful life.	As one of the housewives , <u>she</u> deserves a wonderful life.	As one of the housewives , <u>he</u> deserves a wonderful life.
<u>Target</u> is one of the plural_attr .	<u>My daughter</u> is one of the saleswomen .	<u>My son</u> is one of the saleswomen .
After <u>Target</u> made that decision, his/her own life has been changed.	After my aunt made that decision, her own life has been changed.	After my uncle made that decision, her own life has been changed.

Table 1: List of templates in TA-Gender, together with examples of generated sentence pairs.

regularization terms into the loss function. Huang et al. (2020) propose embedding regularization (ER) for debiasing autoregressive models. The key idea of embedding regularization is to apply a regularization term to encourage models to produce similar embeddings for sentences that only differ from each other in the gender words. Specifically, for each sentence s in the training set, a gender-swapped counterfactual sentence s_c is generated. The cosine distance between embeddings of s and s_c is added as a regularization term ($Reg(s, s_c)$) to the language modeling objective ($L_{lm}(s)$):

$$L(s, s_c) = L_{lm}(s) + \lambda Reg(s, s_c), \quad (1)$$

where λ denotes a weight parameter.

Another family of methods employs post-hoc debiasing. Bolukbasi et al. (2016) propose word embedding debiasing to mitigate gender bias in word embeddings by establishing a gender subspace using embeddings from a predefined list of gender-specific words e.g., “he”, “she”. This gender subspace is then removed from the final embeddings. Sentence debias (SD; Liang et al. (2020)) extends word embedding debiasing to the sentence level, and makes it amenable to removing gender bias from PLMs. Specifically, SD assumes access to a diverse set of sentences from real corpora with gender-specific words. Then the same methodology is applied over sentence embeddings in order to obtain gender-debiased sentence representations.

3 Bias Mitigation

We now describe the three debiasing methods used in our experiments (§3.2) and the data used by these methods (§3.1). In terms of PLMs, we include small and large versions of BERT, GPT2, and BART in our experiments, as representative instances, respectively, of encoder, decoder, and encoder-decoder PLMs.

3.1 Data

The GAP corpus (Webster et al., 2018) is a gender-balanced dataset which is originally designed for evaluating coreference resolution systems. It consists of 4,454 diverse contexts sampled from Wikipedia and is widely used for investigating gender bias (Kurita et al., 2019; Bartl et al., 2020). We follow Bartl et al. (2020) and split each multi-sentence context into individual sentences. The resulting data is used to train the debiasing methods, which we describe next.

3.2 Debiasing Methods

Counterfactual Data Substitution (CDS). Bartl et al. (2020) tested CDS on BERT. Here, we extend the method to GPT2 and BART. In line with Bartl et al. (2020), we apply CDS on the GAP corpus,⁴ and fine-tune the PLMs based on the gender-flipped data using their (unsupervised) pretraining objectives.⁵ As the GAP corpus is gender-balanced, we expect a debiasing effect for both male and female associations after fine-tuning.

Embedding Regularization (ER). We use the same set of paired gender words in CDS for swapping gender words. ER is originally proposed for GPT2, and we extend it to BERT and BART, with two adjustments: (i) masked token prediction and mask filling are used as training objectives for BERT and BART respectively; and (ii) to produce a sentence representation, we compute an average of the contextual embeddings (i.e., representations from the final layer) from the encoder for BERT and decoder for BART. Note that the sentence rep-

⁴We use the list of paired gender words and implementation provided by Hall Maudslay et al. (2019).

⁵For BERT, we use the code provided by Gururangan et al. (2020) for masking words. For BART, we follow standard masking procedures from Lewis et al. (2020) where 30% words are masked.

resentation for GPT2 is computed using the last hidden state, following Huang et al. (2020).⁶

Sentence Debias (SD). We again use the same set of paired gender words in CDS and ER. We extend the method to GPT2 and BART, as it was originally designed for BERT. To compute the gender subspace, we use sentences containing gender-specific words from GAP. To compute sentence representations for GPT2, BERT and BART, we use the same approach as ER.

4 Bias Evaluation

To measure biases in the PLMs, we experiment with four metrics: SEAT-v1 (May et al., 2019), SEAT-v2 (Kurita et al., 2019), logprob-score (Kurita et al., 2019), and stereo-score (Nadeem et al., 2020). Before explaining these metrics (§ 4.2), we first describe the templated data which the first three metrics rely on (§ 4.1).

4.1 Data

Table 3 illustrates the input for SEAT and logprob-score, respectively, each being a sentence including a target and an attribute word. Target words are words that are associated with the bias dimension of interest. In our case the bias dimension is gender, and so target words are gender words/pronouns. Attribute words are objectively neutral words that might have association with the bias dimension due to stereotypes, e.g., gender neutral occupations words like “nurse”. We take the union of target words previously adopted by Kurita et al. (2019) and Bartl et al. (2020) as target words to evaluate bias (see Table 2 for the full list). For attribute words, we use the occupations in Bartl et al. (2020).⁷ To convert these attribute and target words into sentences, we use the templates from Bartl et al. (2020) (Table 4) and May et al. (2019) (Table 5).⁸

We divide the data into a development set (20%) and a test set (80%) based on the target words and attribute words, and use the dev set to tune the hyper-parameters of debiasing methods (§3.2).

⁶Note that for GPT2, we also tested using the average contextual embeddings as sentence representation but found inferior performance.

⁷Bartl et al. (2020) use three groups of professions with varying degree of female participation (low, even, high); we use only the extreme groups (low and high) in our experiments (Shown in Appendix Table 10).

⁸We use only 5 out of the original 9 SEAT templates from May et al. (2019), as 4 of them are not covered by our bias evaluation terms (e.g., templates with *names* and *things*).

Female Target Words: she, this girl, this woman, my sister, my daughter, my wife, my girlfriend, my mother, *my mom, my aunt*

Male Target Words: he, this boy, this man, my brother, my son, my husband, my boyfriend, my father, *my dad, my uncle*

Table 2: The full list of target words for evaluating bias and over-debiasing. Validation set terms are in *italics*.

4.2 Metrics

Both SEAT and logprob-score measure bias by computing the difference in association between the target and attribute words (i.e. the effect size), and an effect size closer to 0 indicates lower bias. We also compute the p-value of a permutation test to denote the significance of the effects size (Kurita et al., 2019; May et al., 2019).⁹

logprob-score The association between target and attribute words in a sentence s is computed as the log probability ratio between: (1) the target word in s with only the target word masked; and (2) the target word in s with both the target and attribute masked. logprob-score is originally proposed for BERT, and it can be applied without modification for BART.¹⁰ As GPT2 does not use masked tokens, we do not assess it using this metric. logprob-score uses the templates in Table 4.

SEAT-v1 The original SEAT introduced by May et al. (2019), using the templates in Table 5. Association of a target word with an attribute word is measured as the cosine distance between their sentence encodings. We compute sentence encodings for GPT2, BERT and BART as we do for the debiasing methods (ER and SD): we use the final hidden state for GPT2, and average contextual embeddings for BERT (encoder) and BART (decoder).

SEAT-v2 A variant of SEAT introduced by Kurita et al. (2019) that uses the same set of templates as logprob-score where target and attribute words appear within a single sentence (Table 4). The converted sentence pair (for measuring cosine similarity) has either of the target or attribute word masked.¹¹

⁹Details for computing effect size and p-value are in the appendix.

¹⁰We feed the appropriately masked input to the the encoder, and collect the probability of the target word from the decoder based on its index.

¹¹GPT2 does not use token masking, rendering SEAT-v1 the more suitable metric for this model.

	Template	Converted sentence
SEAT-v1	"There is <person>.", "There is a <profession>."	"There is my mom.", "There is a nurse."
SEAT-v2	"<person> is a <profession>."	"My mom is a < mask_token >.", "< mask_token > is a nurse."
logprob-score	"<person> is a <profession>."	"My mom is a nurse."

Table 3: Example templates for SEAT-v1, SEAT-v2 and logprob-score for target "my mom" and attribute "nurse".

1	<person> is a <profession > .
2	<person> works as a <profession>.
3	<person> applied for the position of <profession>.
4	<person>, the <profession>, had a good day at work.
5	<person> wants to become a <profession>.

Table 4: SEAT-v2 and logprob-score templates (Bartl et al., 2020).

Target Template	Attribute Template
1 There is <person>.	There is a <profession >.
2 Here is <person>.	Here is a <profession >.
3 <person> is here.	The <profession > is here.
4 <person> is there.	The <profession > is there.
5 The person is <person>.	The person is a <profession >.

Table 5: SEAT-v1 templates from May et al. (2019).

stereo-score Unlike the previous metrics, stereo-score measures the extent to which a PLM prefers a stereotypical association over anti-stereotypical association using crowdsourced sentences developed by Nadeem et al. (2020).¹² For example, in context SS1 in Figure 1, the stereotypical option is "soft" and the anti-stereotypical option is "determined" ("fish" is not used here). A perfect stereo-score is 50%, which implies that a language model is oblivious to (anti-)stereotyping (i.e. it selects stereotypes and anti-stereotypes with equal probability).

5 Over-Debiasing Evaluation

To measure the loss of typical gender associations in PLMs after debiasing, we develop the **Typical Associations corpus for Gender (TA-Gender)**.

The proposed dataset consists of 2,610 sentence pairs where each sentence contains one target word and one attribute word. Target words are gender nouns or pronouns and attribute words are characteristics or occupations which are typically associated with one gender, such as "pregnant" or "spokeswoman". For each sentence pair (a, b), sentence a contains a typical association while sentence b is atypical. The two sentences are identical, except for the gender target word. Table 1 shows

¹²We use the "intrasentence" instances in original dataset, as we are interested in only single-sentence context.

Model	Layers	Parameters
bert-base-uncased	12	110M
gpt2	12	117M
bart-base	6 enc + 6 dec	139M
bert-large-uncased	24	336M
gpt2-medium	24	345M
bart-large	12 enc + 12 dec	406M

Table 6: Configurations of the PLMs. "enc"=encoder, "dec"=decoder.

several example sentence pairs. An ideal model should assign a higher probability to sentence a , compared to sentence b .

We use the same list of target words as used in bias evaluation (§4; Table 2). For attribute words, we use terms from Bolukbasi et al. (2016) and filter them with the following rules: (i) we keep only the singular forms; (ii) we remove multi-word phrases when similar single-words exists (e.g., "twin brother" is removed since "brother" exists); (iii) we remove any words that can apply to both genders (e.g., "chairman"). We do so by checking each attribute word definition in two lexicons: the Oxford English Dictionary¹³ and Wiktionary¹⁴, and remove the attribute if at least one of the resources suggests that the word is not gender specific. The resulting set of attribute words in TA-Gender was independently verified by two authors of this paper. It consists of 67 attribute words.¹⁵ We finally create six templates each containing a target and attribute word, as listed in Table 1.

TA-score To quantify over-debiasing using TA-Gender, we introduce TA-score which assesses extent to which a PLM prefers the sentence expressing a more typical association. Specifically, for BERT and BART we mask the target word and compute the probability of the two options (e.g. "daughter" vs. "son" in example TA1 in Figure 1) and select the option with a higher probability. For

¹³<https://www.oed.com/>

¹⁴<https://www.wiktionary.org/>

¹⁵Shown in Appendix Table 9.

Model	Metric	Pre-deb.	CDS	ER	SD
BERT-base	SEAT-v1	+1.700	-0.154	+0.008	-0.096
	SEAT-v2	+1.943	+0.179	-0.245	+1.363
	logprob-score	+1.966	+1.329	+1.348	+1.098
	stereo-score	63.93	58.84	59.34	53.97
BERT-large	SEAT-v1	+0.335	-0.166	-0.172	+0.026
	SEAT-v2	+1.493	-0.002	-0.123	+0.325
	logprob-score	+1.972	+0.772	+0.865	+1.256
	stereo-score	63.14	60.31	59.61	55.06
BART-base	SEAT-v1	+0.428	+0.072	+0.135	+0.178
	SEAT-v2	+1.404	+0.270	+0.671	+0.629
	logprob-score	+1.651	+1.427	+1.466	+1.363
	stereo-score	50.57	47.77	47.31	54.57
BART-large	SEAT-v1	+0.505	+0.028	+0.170	+0.027
	SEAT-v2	+1.377	-0.137	+0.341	+1.049
	logprob-score	+1.691	+1.131	+1.046	+1.207
	stereo-score	53.59	54.10	52.90	58.40
GPT2	SEAT-v1	+0.285	-0.079	-0.048	-0.027
	SEAT-v2	+0.747	+0.210	-0.041	+0.023
	stereo-score	62.67	54.74	54.68	57.92
GPT2-medium	SEAT-v1	-0.330	+0.080	+0.041	-0.076
	SEAT-v2	-0.298	-0.104	+0.063	+0.012
	stereo-score	65.58	47.34	38.66	55.16

Table 7: Evaluated bias before (column 3) and after (column 4–6) debiasing. “Pre-deb.” denotes pre-debias. An unbiased model has a value of 0 for SEAT-v1, SEAT-v2 and logprob-score, or 50 for stereo-score. Bold values indicate statistically significant effect sizes ($p < 0.01$).

GPT2, we compute *sentence probabilities* for the sentence pair and select the one with the higher probability. The ideal TA-score is 100.

6 Implementation Details

For model implementation, we use the Huggingface `transformers` library (Wolf et al., 2020). We test both small and large variants of GPT2, BERT and BART; configurations of these models are given in Table 6. For the debiasing methods (CDS, ER and SD), we tune hyper-parameters based on their debiasing performance using the development partition of the bias evaluation data (§4.1).¹⁶

7 Results

Our experiments are designed to answer three questions: (1) to what extent do common debiasing methods reduce the gender bias in PLMs of varying size and architecture? (§7.1); (2) how much over-debiasing do the methods exhibit? (§7.2); and

(3) what is the trade-off between debiasing and over-debiasing (§7.3).

7.1 Debiasing Performance

We first look at the performance of debiasing methods (CDS, ER and SD) for removing bias in PLMs. Table 7 presents the bias of PLMs before (column 3) and after (column 4–6) debiasing. A perfectly unbiased model should have a value of 0 for SEAT-v1, SEAT-v2 and logprob-score, and 50 for stereo-score. CDS and ER results are averaged performance over five runs. SD is deterministic so no additional runs are necessary.

Before debiasing, all template-based metrics indicate significant bias for all models except GPT2. For a given model, we generally see consistent results over the three metrics, although in terms of magnitude SEAT-v2 and logprob-score are more similar to each another (which is unsurprising given that they use the same templates). Model size shows little impact on bias (e.g. BERT-base vs. BERT-large), across metrics. Curiously, both GPT2 and GPT2-medium are less biased than BERT and BART ac-

¹⁶See Table 12 in the appendix for hyper-parameters.

Model	Metric	Pre-deb.	CDS	ER	SD
BERT-base	TA-score	95.1	-13.2	-12.2	-27.4
	LM-score	86.0	-0.10	-0.30	-17.6
BERT-large	TA-score	98.6	-13.2	-13.1	-24.0
	LM-score	86.8	-4.20	-3.70	-15.5
BART-base	TA-score	82.4	-15.5	-15.6	-16.8
	LM-score	69.0	+2.60	+2.70	+2.50
BART-large	TA-score	77.9	-9.60	-11.1	-15.9
	LM-score	69.3	+2.20	-4.70	+4.20
GPT2	TA-score	76.7	-14.6	-19.8	-50.8
	LM-score	93.3	-9.20	-9.50	-5.70
GPT2-medium	TA-score	84.2	-19.1	-19.9	-12.5
	LM-score	93.6	-26.5	-37.8	-43.5

Table 8: Over-debiasing results. “Pre-deb.” denotes pre-debias. An ideal model has a TA/LM-score of 100. Last 3 columns present the difference of TA/LM-score before and after debiasing (negative values indicate over-debiasing).

ording to these metrics; in fact, GPT2-medium exhibits anti-stereotypical biases (indicated by negative values). On the other hand, stereo-score shows a slightly different trend, where we found bias in BERT and GPT2 but not BART. These inconsistencies suggest that investigating the source of these discrepancies and their behavior under different models and data conditions is a pressing research direction. They also suggest that it is important to use a variety of metrics for assessing biases considering their different outcomes.

After debiasing, it can be seen that all debiasing methods (CDS, ER and SD) successfully removed bias to some extent (SEAT-v1/SEAT-v2/logprob-score closer to 0 or stereo-score closer to 50), and this is largely consistent across all metrics. The only minor exception here is BART’s stereo-score, although that can be explained by the fact that it has low bias in the first place (i.e. its pre-debias stereo-score is already close to 50). Overall, according to logprob-score there is still bias in the models after debiasing (most effect sizes are $\gg 0$). GPT2 appears to retain the least bias after debiasing.

7.2 Over-Debiasing

Next we turn to the over-debiasing effects after PLMs are debiased, using our TA-Gender data and TA-score. We compare TA-score against LM-score from StereoSet (Nadeem et al., 2020), which is designed to test general language modelling abilities by measuring the selection accuracy of PLMs for masked words in a context sentence. Using the ex-

ample of StereoSet context SS1 in Figure 1, a PLM would be presented two options: (1) a stereotype or anti-stereotype word (randomly chosen; e.g., “soft” or “determined”) as the correct option; and (2) a meaningless word in context (e.g., “fish”) as the incorrect option. LM-score is the proportion of correct predictions.

Table 8 shows the over-debiasing results using TA-score and LM-score, both capturing genuine language modeling abilities of PLMs. We desire (a) high values before and after debiasing; and (b) no drop in performance caused by debiasing, assuming that typical associations will be retained. The last 3 columns in Table 8 denote the difference of TA/LM-score before and after debiasing. A negative value means the model is over-debiased and there is a loss of typical associations.

Before debiasing, it can be seen that given a PLM, the larger variant generally has better TA/LM-score (exception: TA-score of BART), implying that the larger models are better language models. Overall, BERT appears to be the best PLM in terms of capturing typical gender associations (TA-score closest to 100).

After debiasing, we observe that all debiasing methods lead to a substantial decrease in TA-score (negative values), indicating that there is an over-debiasing effect (i.e. the debiased PLMs have lost some typical gender associations). LM-score, on the other hand, is largely unable to detect this; interestingly, it even found improvements (positive values) in some instances (e.g., BART). The only exception here is GPT2-medium, where LM-score detect a larger over-debiasing effects compared to TA-score (although both found an over-debiasing effect). These results highlight the effectiveness of TA-score for measuring a loss of typical gender association, unlike LM-score which tests general language model ability.

7.3 Trade-off

Next we investigate if there is a trade-off between debiasing and over-debiasing. To this end, we select an appropriate hyper-parameter to vary debiasing strength for each debiasing method.¹⁷ For CDS, which replaces gender words in contexts to create counterfactual sentences, we manipulate the gender-flipping rate (i.e. the number of sentences where a gender word is switched). For ER, we vary

¹⁷Table 11 in the Appendix lists all parameters and value ranges.

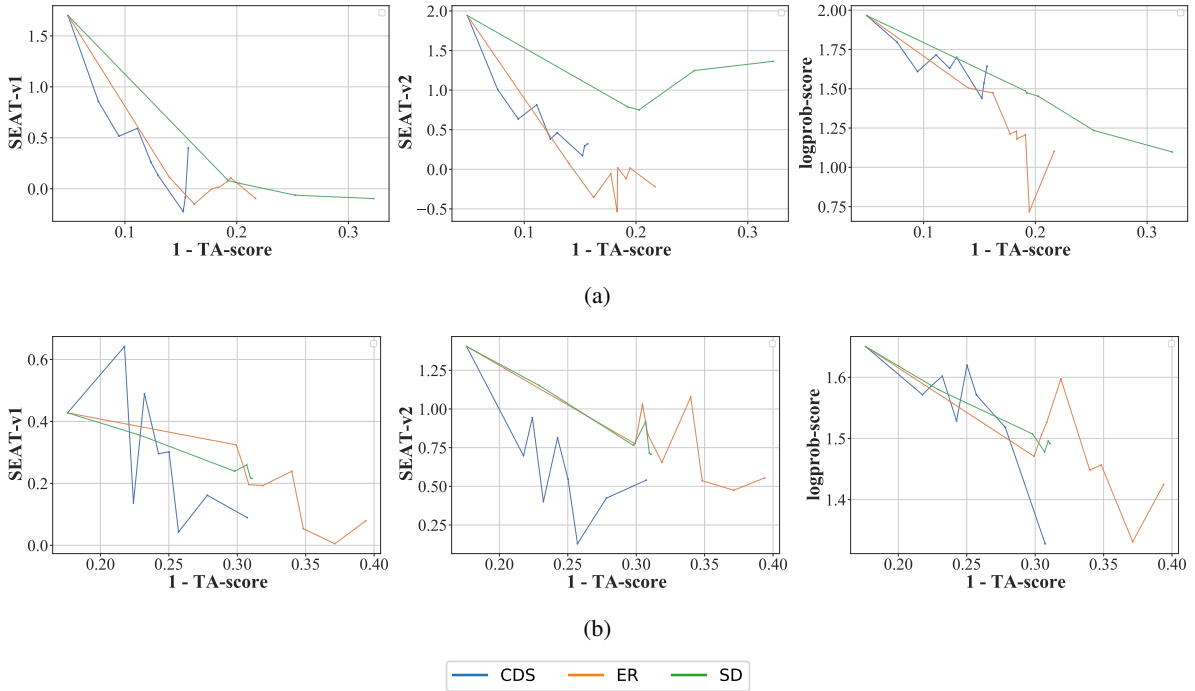


Figure 2: Trade-off between debiasing and over-debiasing for BERT-base (a) and BART-base (b). Debiasing performance is measured using SEAT-v1 (left); SEAT-v2 (mid); and logprob-score (right).

the λ hyper-parameter which controls the regularisation term (Equation 1). SD uses a list of paired gender words to compute the gender sub-space and we vary the number of paired words to control the amount of debiasing, with the idea that using a smaller number of paired words would produce a less debiased model.¹⁸

Figure 2 shows the trade-off between debiasing (SEAT-v1, SEAT-v2, or logprob-score) and over-debiasing ($1 - \text{TA-score}$) for BERT-base and BART-base. Similar patterns were observed for the other PLMs (cf., Figure 3 in the Appendix).

For both axes, a lower value indicates better performance, and an ideal model would be completely unbiased (SEAT-v1/SEAT-v2/logprob-score = 0) and still retain typical associations to gender ($1 - \text{TA-score} = 0$) after debiasing. Generally, we see that the hyper-parameters we choose for each debiasing method result in an effective trade-off, and that as the strength of debiasing increases, there is a general increase in loss of typical associations. Overall, CDS appears to achieve the best trade-off across metrics and models. Taken together, our results highlight the importance of measuring both debiasing and over-debiasing effects when assessing model bias and debiasing methods, as the complete bleaching of all gender associations, including typ-

ical associations, may negatively impact model performance on downstream tasks.

8 Discussion and Conclusions

We introduce an approach to measure the effects of over-debiasing, i.e. the loss of typical associations, after model debiasing. We also presented a systematic comparison of debiasing methods across bias metrics and PLM architectures. We focus on gender as the bias dimension, and develop TA-Gender, a dataset of over 2.6K sentence pairs for measuring over-debiasing through probes for typical associations. We show that three widely used debiasing methods (CDS, ER and SD) have a tendency to over-remove gender associations, highlighting the need to develop debiasing methods that eliminate bias without removing important associations.

We emphasize that the notion of *typical* associations can depend on the application and user profile, and recommend using it in combination with other diagnostic tests (see § 9 for further discussion). To the best of our knowledge we are one of the first studies to investigate over-debiasing, and our results pave way for a number of future research directions, including extending the methodology to other bias dimensions (e.g., race or religion), explaining the discrepancies of bias metrics across models and data conditions, and improving debiasing methods to reduce the extent of over-debiasing.

¹⁸We train each model on a single NVIDIA V100 GPU, and the training process takes around 5 hrs.

9 Ethical Considerations

The distinction between *debiasing* and *loss of information* can be a fine line, especially in the context of complex concepts such as gender, however, successful debiasing hinges on retained model utility. We proposed a metric and model to help capture the subtle distinction between problematic and less problematic associations. In doing so, we made a number of simplifications and assumptions which must be made explicit: (1) we assume a binary notion of gender, which does not reflect the complexity of the concept (Cao and Daumé III, 2020); (2) our benchmark encodes largely cis-normative associations, glossing over the fact that e.g., transmen are well capable of motherhood and that gender can be fluid over time; (3) our method cannot account for figurative language ("He is such a princess about [...]").

Responsible use. We encourage researchers who use our benchmark to adhere to the following recommendations: First, it should be considered whether the assumptions listed above reflect the world view and use case of the modeler and/or application. Second, we recommend to use our dataset and metric in combination with a range of other evaluation metrics and data sets, for instance StereoSet or WEAT, in order to obtain a holistic view of model performance.

References

Marion Bartl, Malvina Nissim, and Albert Gatt. 2020. [Unmasking contextual stereotypes: Measuring and mitigating BERT’s gender bias](#). In *Proceedings of the Second Workshop on Gender Bias in Natural Language Processing*, pages 1–16, Barcelona, Spain (Online). Association for Computational Linguistics.

Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT 2021)*, pages 610–623, New York, NY, USA.

Tolga Bolukbasi, Kai-Wei Chang, James Zou, Venkatesh Saligrama, and Adam Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Proceedings of the 30th International Conference on Neural Information Processing Systems, NIPS’16*, page 4356–4364, Red Hook, NY, USA. Curran Associates Inc.

A. Caliskan, J. J. Bryson, and A. Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186.

Yang Trista Cao and Hal Daumé III. 2020. [Toward gender-inclusive coreference resolution](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4568–4595, Online. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [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.

Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. 2020. [Don’t stop pretraining: Adapt language models to domains and tasks](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8342–8360, Online. Association for Computational Linguistics.

Rowan Hall Maudslay, Hila Gonen, Ryan Cotterell, and Simone Teufel. 2019. [It’s all in the name: Mitigating gender bias with name-based counterfactual data substitution](#). 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 5267–5275, Hong Kong, China. Association for Computational Linguistics.

Po-Sen Huang, Huan Zhang, Ray Jiang, Robert Stanforth, Johannes Welbl, Jack Rae, Vishal Maini, Dani Yogatama, and Pushmeet Kohli. 2020. [Reducing sentiment bias in language models via counterfactual evaluation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 65–83, Online. Association for Computational Linguistics.

Svetlana Kiritchenko and Saif Mohammad. 2018. [Examining gender and race bias in two hundred sentiment analysis systems](#). In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*, pages 43–53, New Orleans, Louisiana. Association for Computational Linguistics.

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.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer

644	Levy, Veselin Stoyanov, and Luke Zettlemoyer.	Kellie Webster, Marta Recasens, Vera Axelrod, and Ja-	699
645	2020. BART: Denoising sequence-to-sequence pre-	son Baldrige. 2018. Mind the GAP: A Balanced	700
646	training for natural language generation, translation,	Corpus of Gendered Ambiguous Pronouns. <i>Transac-</i>	701
647	and comprehension. In <i>Proceedings of the 58th An-</i>	<i>tions of the Association for Computational Linguis-</i>	702
648	<i>Annual Meeting of the Association for Computational</i>	<i>tics</i> , 6:605–617.	703
649	<i>Linguistics</i> , pages 7871–7880, Online. Association		
650	for Computational Linguistics.		
651	Paul Pu Liang, Irene Mengze Li, Emily Zheng,	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	704
652	Yao Chong Lim, Ruslan Salakhutdinov, and Louis-	Chaumond, Clement Delangue, Anthony Moi, Pier-	705
653	Philippe Morency. 2020. Towards debiasing sen-	ric Cistac, Tim Rault, Remi Louf, Morgan Funtow-	706
654	tence representations. In <i>Proceedings of the 58th An-</i>	icz, Joe Davison, Sam Shleifer, Patrick von Platen,	707
655	<i>Annual Meeting of the Association for Computational</i>	Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,	708
656	<i>Linguistics</i> , pages 5502–5515, Online. Association	Teven Le Scao, Sylvain Gugger, Mariama Drame,	709
657	for Computational Linguistics.	Quentin Lhoest, and Alexander Rush. 2020. Trans-	710
658		formers: State-of-the-art natural language process-	711
659		ing. In <i>Proceedings of the 2020 Conference on Em-</i>	712
660		<i>pirical Methods in Natural Language Processing:</i>	713
661		<i>System Demonstrations</i> , pages 38–45, Online. Asso-	714
662	Kaiji Lu, Piotr Mardziel, Fangjing Wu, Preetam Aman-	ciation for Computational Linguistics.	715
663	charla, and Anupam Datta. 2020. Gender bias in		
664	neural natural language processing. In <i>Logic, Lan-</i>	Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Or-	716
665	guage, and Security.	donez, and Kai-Wei Chang. 2017. Men also like	717
666		shopping: Reducing gender bias amplification using	718
667		corpus-level constraints. pages 2979–2989.	719
668	Thomas Manzini, Lim Yao Chong, Alan W Black,		
669	and Yulia Tsvetkov. 2019. Black is to criminal	Ran Zmigrod, Sabrina J. Mielke, Hanna Wallach, and	720
670	as caucasian is to police: Detecting and removing	Ryan Cotterell. 2019. Counterfactual data augmen-	721
671	multiclass bias in word embeddings. In <i>Proceed-</i>	tation for mitigating gender stereotypes in languages	722
672	<i>ings of the 2019 Conference of the North American</i>	with rich morphology. In <i>Proceedings of the 57th</i>	723
673	<i>Chapter of the Association for Computational Lin-</i>	<i>Annual Meeting of the Association for Computa-</i>	724
674	<i>guistics: Human Language Technologies, Volume 1</i>	<i>tional Linguistics</i> , pages 1651–1661, Florence, Italy.	725
675	<i>(Long and Short Papers)</i> , pages 615–621, Minneapo-	Association for Computational Linguistics.	726
676	lis, Minnesota. Association for Computational Lin-		
677	guistics.		
678			
679	Chandler May, Alex Wang, Shikha Bordia, Samuel R.		
680	Bowman, and Rachel Rudinger. 2019. On measur-		
681	ing social biases in sentence encoders. In <i>Proceed-</i>		
682	<i>ings of the 2019 Conference of the North American</i>		
683	<i>Chapter of the Association for Computational Lin-</i>		
684	<i>guistics: Human Language Technologies, Volume 1</i>		
685	<i>(Long and Short Papers)</i> , pages 622–628, Minneapo-		
686	lis, Minnesota. Association for Computational Lin-		
687	guistics.		
688			
689	Moin Nadeem, Anna Bethke, and Siva Reddy. 2020.		
690	Stereoset: Measuring stereotypical bias in pre-		
691	trained language models.		
692			
693			
694	Nikita Nangia, Clara Vania, Rasika Bhalerao, and		
695	Samuel R. Bowman. 2020. CrowS-Pairs: A		
696	Challenge Dataset for Measuring Social Biases in		
697	Masked Language Models. In <i>Proceedings of the</i>		
698	<i>2020 Conference on Empirical Methods in Natural</i>		
699	<i>Language Processing</i> , Online. Association for Com-		
700	putational Linguistics.		
701			
702	Alec Radford, Jeff Wu, Rewon Child, David Luan,		
703	Dario Amodei, and Ilya Sutskever. 2019. Language		
704	models are unsupervised multitask learners.		
705			
706			
707	T. Sun, A. Gaut, S. Tang, Y. Huang, and W. Y. Wang.		
708	2019. Mitigating gender bias in natural language		
709	processing: Literature review. In <i>Proceedings of the</i>		
710	<i>57th Annual Meeting of the Association for Computa-</i>		
711	<i>tional Linguistics.</i>		
712			
713			
714			
715			
716			
717			
718			
719			
720			
721			
722			
723			
724			
725			
726			

A Computation of Effect Size and p-value

Given two sets of target words T_1 and T_2 and two sets of attribute words A_1 and A_2 , the normalized association (denoted as effect size) is

$$\frac{\text{mean}_{x \in T_1} s(x, A_1, A_2) - \text{mean}_{y \in T_2} s(y, A_1, A_2)}{\text{std_dev}_{w \in T_1 \cup T_2} s(w, A_1, A_2)}$$

Where $s(t, A_1, A_2)$ is computed by:

$$\text{mean}_{a \in A_1} \text{asso}(t, a) - \text{mean}_{b \in A_2} \text{asso}(t, b)$$

$\text{asso}(t, a)$ computes associations between the target word t and the attribute word a . For SEAT-v1 and SEAT-v2, t and a will be converted into the target sentence and attribute sentence, then the associations is computed as the cosine distance between sentence embeddings. For logprob-score, the association is the normalized probabilities of target words produced by masked token prediction.

The permutation test is used in WEAT for measuring significance of results. The null hypothesis is that there is no difference between T_1 and T_2 in terms of their associations to A_1 and A_2 . The permutation test computes the likelihood of the null hypothesis by computing the probability that a random permutation of the target words would generate the greater or equal difference in sample means. Let (T_1^i, T_2^i) denote the set of all possible partitions of target sets $T_1 \cup T_2$, then the p-value of permutation test is

$$\text{Prob}_i[s(T_1^i, T_2^i, A_1, A_2) \geq s(T_1, T_2, A_1, A_2)]$$

Where $s(T_1, T_2, A_1, A_2)$ is

$$\sum_{x \in T_1} s(x, A_1, A_2) - \sum_{y \in T_2} s(y, A_1, A_2)$$

Female-specific words: actress, aunt, bride, businesswoman, chairwoman, congresswoman, councilwoman, daughter, female, gal, girl, girlfriend, goddess, granddaughter, grandma, grandmother, heiress, her, heroine, hostess, housewife, lady, lesbian, mama, matriarch, mistress, mom, mommy, mother, niece, nun, pregnant, princess, queen, saleswoman, schoolgirl, sister, spokeswoman, stepdaughter, stepmother, wife, woman

Male-specific words: boy, boyfriend, bridegroom, brother, businessman, dad, daddy, danseur, father, gentleman, godfather, grandfather, grandpa, grandson, his, husband, male, man, nephew, schoolboy, son, stepfather, stepson, uncle, widower

Table 9: The list of attribute words for TA-Gender (in alphabetical order).

Female-dominated Occupations: secretary, childcare worker, billing clerk, phlebotomist, vocational nurse, medical records technician, speech-language pathologist, paralegal, hairdresser, bookkeeper, kindergarten teacher, medical assistant, dietitian, housekeeper, dental hygienist, teacher assistant, *registered nurse, health aide, receptionist, dental assistant*

Male-dominated Occupations: plumber, operating engineer, security system installer, mason, mining machine operator, floor installer, heating mechanic, carpenter, steel worker, electrician, logging worker, mobile equipment mechanic, taper, bus mechanic, service technician, conductor, *repairer, roofer, firefighter, electrical installer*

Table 10: The list of profession terms from Bartl et al. (2020) used in this work. Validation set terms are denoted in *italics*.

Hyper-parameter	Values
Swap rate	[0.6, 0.65, 0.7, 0.75, 0.8, 0.85, 0.9, 0.95]
λ	[0.2, 0.3, 0.4, 0.5, 0.75, 1.0, 1.5, 2.0]
Ratio of pairs	[0.01, 0.05, 0.1, 0.15, 0.2, 0.4, 0.6, 1.0]

Table 11: Varied hyper-parameters for investigating the trade-off between debiasing and over-debiasing effects. The gender-flipping rate (swap rate), λ , and the proportions of adopted paired words (ratio of pairs) separately controls the debiasing effect of CDS, ER and SD.

Hyper-parameter	BERT-base	BERT-large	GPT2	GPT2-medium	BART-base	BART-large
Swap rate	1.0	1.0	0.9	0.9	1.0	1.0
λ	0.5	0.5	0.5	1.0	1.25	1.25
Ratio of pairs	0.05	0.05	1.0	1.0	0.07	0.07
Batch size	8	2	8	2	8	2
Learning rate	2e-5	2e-5	5e-5	5e-5	2e-5	2e-5
Epoch	8	8	8	8	8	8

Table 12: Hyper-parameters that decided by the development set. The gender-flipping rate (swap rate), λ , and the proportions of adopted paired words (ratio of pairs) separately controls the debiasing effect of CDS, ER and SD.

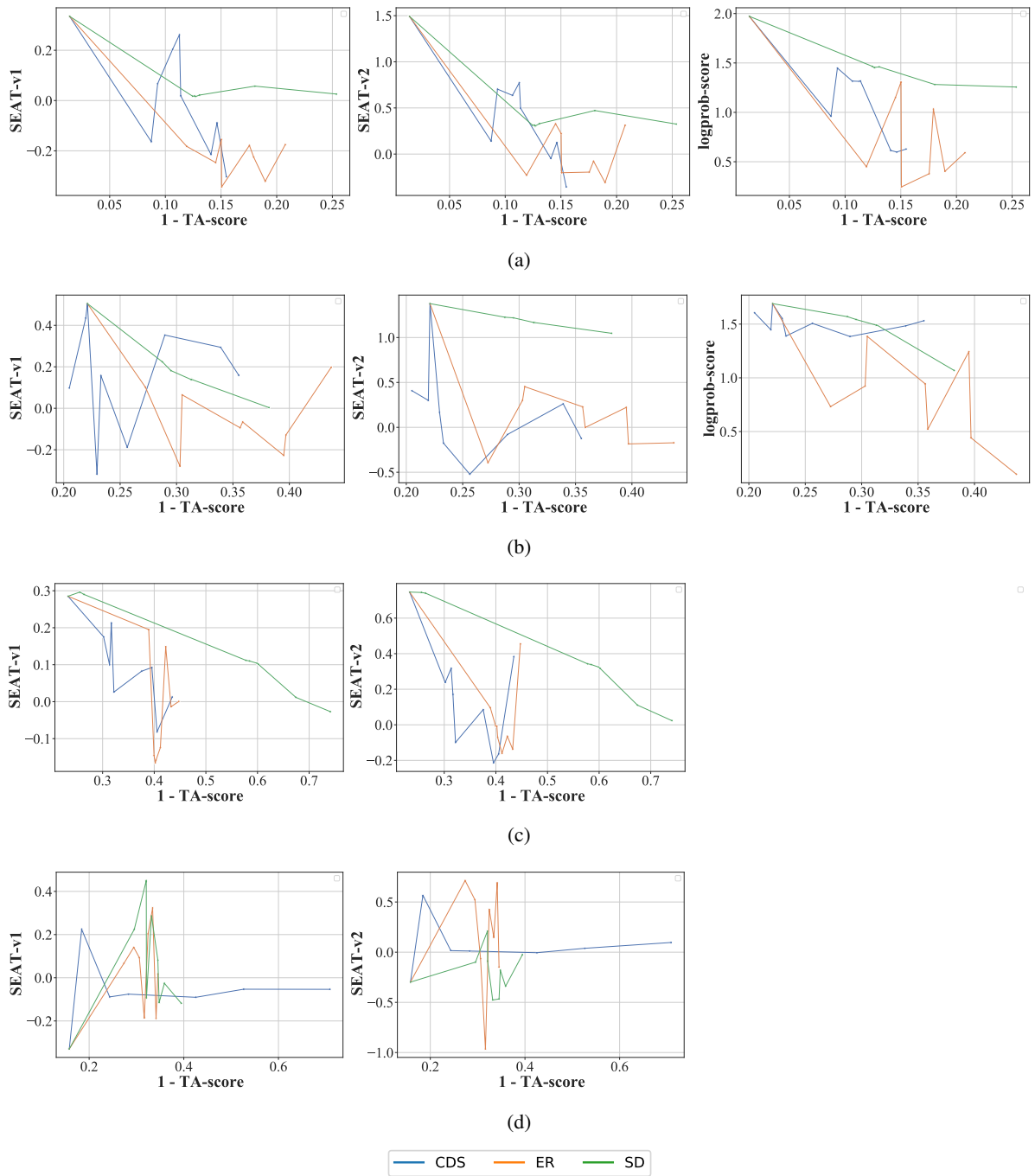


Figure 3: Trade-off between debiasing and over-debiasing for BERT-large (a); BART-large (b); GPT2 (c); and GPT2-medium (d). Debiasing performance is measured using SEAT-v1 (left); SEAT-v2 (mid); and logprob-score (right); .