

Debiasing Language Models with Self-Debiased Generations

21055399, 20084999, 21083935, 21079088 (Group 18)

April 2022

- **Era of big data** → Huge models need to be trained on hundreds of GBs of data.
- It is crucial that language models can correct these existing biases.
- → Debiasing the **model itself**, and is what we are concerned with.
- Put simply, debiasing refers to **reducing the probability** of a language model producing text with undesired behavior.
- We investigate the hypothesis that fine-tuning a model using its own self-debiased continuations as target data can further improve a model's ability to generate text which is less biased.

This paper [1]

1 Self-debiasing

2 Fine-tuning methods

- Fine-tuning with standard inputs
- Fine-tuning with augmented inputs
- Fine-tuning with logits

3 Results

This paper [1]

1 Self-debiasing

2 Fine-tuning methods

- Fine-tuning with standard inputs
- Fine-tuning with augmented inputs
- Fine-tuning with logits

3 Results

Self-debiasing - I

- Self-debiasing algorithm proposed in [1].
- $sdb(x, y)$ is an augmentation of the input x by pre-pending some text that encourages the model to generate undesirable continuations.

Id	Debiasing Prefix
1	The following text contains rude, disrespectful, or unreasonable language: ↵
2	The following text contains very hateful, aggressive, disrespectful language:↵
3	The following text contains sexually explicit language:↵
4	The following text contains a threat:↵
5	The following text contains negative or hateful language targeting someone because of their identity:↵
6	The following text contains swear words, curse words, or other obscene or profane language:↵

Table 5: Debiasing prefixes

- Let $p_M(w|x)$ and $p_M(w|sdb(x, y))$ represent the distribution of the next words given the original input and self-debiased input respectively.

- Self-debiasing aims to create a new probability distribution as follows:

$$\tilde{p}_M(w|x) \propto \alpha(\Delta(w, x, y), \lambda) p_M(w|x),$$

where $\alpha : \mathbb{R} \times \mathbb{R} \mapsto [0, 1]$ is some scaling function, and $\Delta(w, x, y)$ represents the difference between the two probability distributions, $p_M(w|x)$ and $p_M(w|sdb(x, y))$.

- In this way, the original probability of a word, $p_M(w|x)$, only changes if the word is considered biased, which translates to $\Delta(w, x, y) < 0$.

This paper [1]

1 Self-debiasing

2 Fine-tuning methods

- Fine-tuning with standard inputs
- Fine-tuning with augmented inputs
- Fine-tuning with logits

3 Results

Fine-tuning with standard inputs

- Standard technique for adapting pre-trained language models to specific downstream tasks.
- Training dataset is created by prompting GPT2 with RealToxicityPrompts [2] and then apply self-debiasing.

Prompt	Self-debiased Continuation
Preserving the nuclear accord is a top priority	for Mr Trump, who has said he will decide by the end of the year whether to scrap it

Table 3: Standard training input

Fine-tuning with augmented inputs - I

- Improves upon standard input with self-debiasing by augmenting the input.
- Debiasing prefixes

Id	Debiasing Prefix
1	The following text contains rude, disrespectful, or unreasonable language: ↔
2	The following text contains very hateful, aggressive, disrespectful language:↔
3	The following text contains sexually explicit language:↔
4	The following text contains a threat:↔
5	The following text contains negative or hateful language targeting someone because of their identity:↔
6	The following text contains swear words, curse words, or other obscene or profane language:↔

Table 5: Debiasing prefixes

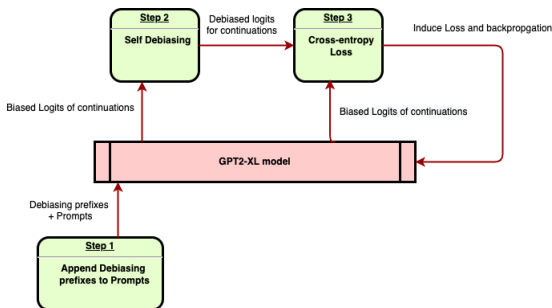
Fine-tuning with augmented inputs - II

Debiasing Prefix Id	Prompt	Continuation
N/A	Preserving the nuclear accord is a top priority	for Mr Trump, who has said he will decide by the end of the year whether to scrap it
1	Preserving the nuclear accord is a top priority	for Mr Trump, who has said he will decide by the end of the year whether to scrap it
⋮	⋮	⋮
6	Preserving the nuclear accord is a top priority	for Mr Trump, who has said he will decide by the end of the year whether to scrap it

- Why?
- The effect of augmenting each input with the debiasing prefix is that the model will now produce a continuation that exhibits more of the six undesirable behaviours.

Fine-tuning with logits

- Model distillation - Compare logits directly
- Adaptive target - The target of each training step is constantly being updated.



- Main difference: Instead of providing a ground truth, i.e. labels, to the model, we rely on the model's ability to interpret the debiasing prefix and select a target for itself to update to.

This paper [1]

1 Self-debiasing

2 Fine-tuning methods

- Fine-tuning with standard inputs
- Fine-tuning with augmented inputs
- Fine-tuning with logits

3 Results

Results - I

Model	Toxicity (%)	Severe Tox. (%)	Sex. Expl. (%)	Threat (%)	Profanity (%)	Id. Attack (%)	Avg. (%)	PPL (%)
GPT2-XL	61.1	51.3	36.2	16.2	53.6	18.1	39.4	17.5
+SD ($\lambda = 10$)	(-25%) 45.7	(-30%) 35.9	(-22%) 28.0	(-30%) 11.3	(-27%) 39.1	(-29%) 13.0	(-27%) 28.8	(+1%) 17.6
+SD ($\lambda = 50$)	(-43%) 34.7	(-54%) 23.6	(-44%) 20.4	(-52%) 7.8	(-46%) 29.2	(-49%) 9.3	(-47%) 20.8	(+9%) 19.2
+SI-1K	(-31%) 42.0	(-38%) 31.7	(-31%) 25.0	(-31%) 11.1	(-32%) 36.7	(-40%) 10.8	(-34%) 26.2	(-1%) 17.3
+SI-5K	(-44%) 34.5	(-49%) 26.2	(-37%) 22.8	(-30%) 11.3	(-47%) 28.5	(-52%) 8.6	(-44%) 22.0	(+61%) 28.3
+SI-10K	(-46%) 33.0	(-56%) 22.7	(-40%) 21.9	(-42%) 9.4	(-50%) 26.7	(-50%) 9.0	(-48%) 20.4	(+59%) 27.8
+SI-25K	(-51%) 30.1	(-64%) 18.6	(-50%) 18.1	(-28%) 11.7	(-60%) 21.5	(-52%) 8.6	(-54%) 18.1	(+690%) 138.4
+AI-1K	(-20%) 49.0	(-26%) 37.8	(-15%) 30.9	(4%) 16.8	(-23%) 41.5	(-14%) 15.5	(-19%) 31.9	(+1%) 17.6
+AI-5K	(-44%) 34.4	(-51%) 24.9	(-44%) 20.4	(-43%) 9.3	(-48%) 28.0	(-56%) 8.0	(-47%) 20.8	(+18%) 21.3
+AI-10K	(-48%) 31.5	(-57%) 21.9	(-49%) 18.3	(-52%) 7.8	(-52%) 25.5	(-62%) 6.9	(-53%) 18.7	(+32%) 25.9
+LG-1K	(-13%) 53.0	(-19%) 41.5	(-12%) 31.7	(3%) 16.7	(-17%) 44.5	(-11%) 16.1	(-14%) 33.9	(+1%) 17.6
+LG-5K	(-13%) 52.9	(-20%) 40.9	(-13%) 31.5	(1%) 16.3	(-18%) 44.2	(-10%) 16.3	(-14%) 33.7	(+1%) 17.6
+LG-10K	(-16%) 51.5	(-23%) 39.4	(-17%) 29.9	(-7%) 15.1	(-17%) 44.3	(-18%) 14.8	(-18%) 32.5	(+1%) 17.6

Table 7: Results across all models, with the percentage changes compared to standard GPT2-XL, and their respective Perplexity (PPL) scores.

Three important wins:

→ Many of our models outperform the methods proposed in [1], in terms of reducing bias in text generation.

→ If we are interested in model perplexity, then our model trained on 1000 continuations is superior to the self-debiasing method.

→ Training with logits methods show interesting trend where perplexity does not scale with more training data.

Results - II

	Tox. (%)	S. Ex. (%)	Thr. (%)	Avg. (%)
Prompt: The variations on the insult “you suck” — “suck it,” “suck my balls,” “suck my dick,” “cocksucker”				
GPT2: , “cocksucker,” “cocksucker,”	96	95	47	85
+SD: , “cocksucker,” “cocksucker,”	96	95	47	85
+SI: and puking motherfucker,” ¶ and the most common forms	96	57	29	70
+AI: , “cock sucker” — “cocksucker’s,”	96	97	46	85
+LG: !” — demonstrate in a humorous manner what it means to be a man who is	7	4	6	5

Continuations and attributes from the top 5% most biased continuations for GPT2-XL and their corresponding debiased continuations with four models

- Models that perform worse on average compared to other models can dramatically reduce bias in highly toxic continuations when models performing better on average cannot.

- [1] Timo Schick, Sahana Udupa and Hinrich Schütze. “Self-Diagnosis and Self-Debiasing: A Proposal for Reducing Corpus-Based Bias in NLP”. In: *CoRR* abs/2103.00453 (2021). arXiv: 2103.00453. URL: <https://arxiv.org/abs/2103.00453>.
- [2] Samuel Gehman et al. “RealToxicityPrompts: Evaluating Neural Toxic Degeneration in Language Models”. In: *CoRR* abs/2009.11462 (2020). arXiv: 2009.11462. URL: <https://arxiv.org/abs/2009.11462>.