

On the Unintended Social Bias of Training Language Generation Models with Latin American Newspapers

Anonymous Author(s)

Abstract

1 Gender bias is a significant problem when generating text, and its unintended
 2 memorization could impact the user experience of many applications (e.g., the
 3 auto-complete feature in Gmail). In this abstract, we introduce a novel architecture
 4 that decouples the representation learning of a neural model from its memory
 5 management role. This architecture allows us to update a memory module with an
 6 equal ratio across gender types breaking biased correlations directly in the latent
 7 space. We show that our approach can mitigate gender bias amplification in the
 8 automatic generation of articles from Latin American newspapers.

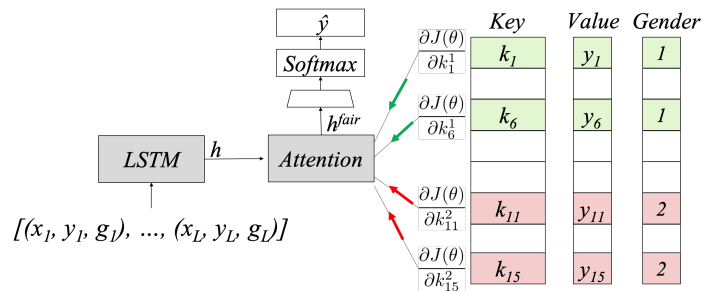


Figure 1: A Fair Region in memory M consists of the most similar keys to h given a uniform distribution over genders (e.g., 1: male, and 2: female). The input consists of a sequence tokens annotated with gender information, e.g., (*The*, 0), (*president*, 0), (*gave*, 0), (*her*, 2), (*speech*, 0).

9 1 Memory Networks and Fair Region

10 As illustrated in Figure 1, the memory M consists of arrays K and V that store *keys* (embeddings)
 11 and *values* (class labels), respectively as in Kaiser et al. (2017). We extend this model with an array
 12 G to store the *gender* associated to each word, (e.g., *actor* is male, *actress* is female, and *scientist*
 13 is no-gender). A neural encoder $f(x, \theta)$ with trainable parameters θ receives an observation x and
 14 generates the activations h in a hidden layer. We want to incorporate a normalized h to M . Hence,
 15 let i_{max} be the index of the most similar key ($i_{max} = \text{argmax}_i \{h \cdot K[i]\}$), then writing the triplet
 16 (x, y, g) to M consist of: $K[i_{max}] = \|h + K[i_{max}]\|$, $V[i_{max}] = y$, and $G[i_{max}] = g$. However,
 17 word embeddings are severely biased in natural language. For example, it has been shown that *man*
 18 is closer to *programmer* than *woman*, Bolukbasi et al. (2016). Similar problems have been recently
 19 observed in embedding algorithms such as Word2Vec, Glove, and BERT, Kurita (2019). For this
 20 purpose, we propose the use of a subset of keys to define a Fair Region in which we can control
 21 an equal ratio of gender types in such region. These keys will compute error signals and generate
 22 gradients that will flow through the entire architecture with backpropagation, as depicted in Figure 1.
 23 We define this region as follows.

MODEL	PERPLEXITY			BIAS AMPLIFICATION		
	ALL	PERU	MEXICO	ALL	PERU-MEXICO	MEXICO-PERU
SEQ2SEQ	13.27	15.31	15.61	+0.18	+0.25	+0.21
SEQ2SEQ+ATTENTION	10.73	13.25	14.08	+0.25	+0.32	+0.29
SEQSEQ+FAIRREGION	10.79	13.04	13.91	+0.09	+0.17	+0.15

Table 1: Perplexity and Bias Amplification results on the datasets of crawled newspapers.

24 **Definition 1.1. (Fair Region)** Let h be an latent representation of the input and M be an external
25 memory. The *male*-neighborhood of h is represented by the indices of the n -nearest keys to h in
26 decreasing order and that share the same gender type *male* as $\{i_1^m, \dots, i_k^m\} = KNN(h, n, male)$.
27 Repeating the same process for each gender type estimates the indices i^f and i^{ng} for the *female* and
28 *non-gender* neighborhoods. Then, the *FairRegion* of M given h consists of $K[i^m; i^f; i^{ng}]$.

29 2 Bias Amplification

30 Inspired by Zhao et al. (2017), we compute the bias score of a word x considering its word embed-
31 ding $h^{fair}(x)$ ¹ and two gender indicators (words *man* and *woman*). For example, the bias score
32 of *scientist* is: $b(\text{scientist}, \text{man}) = \frac{\|h^{fair}(\text{scientist})\| \cdot \|h^{fair}(\text{man})\|}{\|h^{fair}(\text{scientist}) \cdot h^{fair}(\text{man}) + h^{fair}(\text{scientist}) \cdot h^{fair}(\text{woman})\|}$. If
33 the bias score during testing is greater than the one during training, $b^{test}(\text{scientist}, \text{man}) -$
34 $b^{train}(\text{scientist}, \text{man}) > 0$, then the bias of *man* towards *scientist* has been amplified by the
35 model while learning such representation, given training and testing datasets similarly distributed.

36 3 Experiments

37 Figure 1 illustrates our proposed model (**Seq2Seq+FairRegion**), which we use as the LSTM decoder
38 of a Seq2Seq architecture Sutskever, Vinyals, and Le (2014). This decoder attends to a Fair Region
39 of 9 entries (3 for each gender type) given Definition 1.1. We crawl the websites of three newspapers
40 from Chile, Peru, and Mexico to collect a working dataset. To enable a fair comparison, we limit
41 the number of articles to 20,000 for each domain and a vocabulary of 18,000 most common words.
42 Datasets are split into 60%, 20%, and 20% for training, validation, and testing. We want to see if
43 there are correlations showing stereotypes across different nations. *Does the biased correlations*
44 *learned by an encoder transfer to the decoder considering word sequences from different countries?*
45 We compare our approach with these baseline models: 1) **Seq2Seq** Sutskever, Vinyals, and Le (2014);
46 An encoder-decoder architecture that maps between sequences. 2) **Seq2Seq+Attention** Bahdanau,
47 Cho, and Bengio (2015): A Seq2Seq that attends to parts of the input to predict the target word.

48 3.1 Fair Region Results in Similar Perplexity

49 We evaluate all the models with test *perplexity*, which is the exponential of the loss. We report in
50 Table 1 the average perplexity of the aggregated dataset from Peru, Mexico, and Chile, and also
51 from specific countries. Our main finding is that our approach (Seq2Seq+FairRegion) shows similar
52 perplexity values (10.79) than the Seq2Seq+Attention baseline model (10.73) when generating word
53 sequences despite using the Fair Region strategy. These results encourage the use of a controlled
54 region as an automatic technique that maintains the efficacy of generating text. We observe a larger
55 perplexity for country-based datasets, likely because of their smaller training datasets.

56 3.2 Fair Region Controls Bias Amplification

57 We compute the *bias amplification* metric for all models, as defined in Section 2, to study the effect
58 of amplifying potential bias in text for different language generation models. Table 1 shows that
59 using Fair Regions is the most effective method to mitigate bias amplification when combining all the
60 datasets (+0.09). Instead, both Seq2Seq (+0.18) and Seq2Seq+Attention (+0.25) amplify gender bias

¹For Seq2Seq neural models, this word embedding is the output of the decoder component $h^{deco}(x)$

61 for the same corpus. Interestingly, feeding the encoders with news articles from different countries
62 decreases the advantage of using a Fair Region and also amplifies more bias across all the models. In
63 fact, training the encoder with news from Peru has, in general, a larger bias amplification than training
64 it with news from Mexico. This could have many implications and be a product of the writing style
65 or transferred social bias across different countries. We take its world-wide study as future work.

66 **References**

67 Bahdanau, D.; Cho, K.; and Bengio, Y. 2015. Neural machine translation by jointly learning to
68 align and translate. In *3rd International Conference on Learning Representations, ICLR 2015, San*
69 *Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.

70 Bolukbasi, T.; Chang, K.-W.; Zou, J.; Saligrama, V.; and Kalai, A. 2016. Man is to computer
71 programmer as woman is to homemaker? debiasing word embeddings. NIPS'16, 4356–4364.

72 Kaiser, L.; Nachum, O.; Roy, A.; and Bengio, S. 2017. Learning to remember rare events.

73 Kurita, K. 2019. Acl web. 166–172.

74 Miller, A.; Fisch, A.; Dodge, J.; Karimi, A.-H.; Bordes, A.; and Weston, J. 2016. Key-value
75 memory networks for directly reading documents. In *Proceedings of the 2016 Conference on*
76 *Empirical Methods in Natural Language Processing*, 1400–1409. Austin, Texas: Association for
77 Computational Linguistics.

78 Sutskever, I.; Vinyals, O.; and Le, Q. V. 2014. Sequence to sequence learning with neural networks. In
79 Ghahramani, Z.; Welling, M.; Cortes, C.; Lawrence, N. D.; and Weinberger, K. Q., eds., *Advances*
80 *in Neural Information Processing Systems 27*. Curran Associates, Inc. 3104–3112.

81 Zhao, J.; Wang, T.; Yatskar, M.; Ordonez, V.; and Chang, K.-W. 2017. Men also like shopping:
82 Reducing gender bias amplification using corpus-level constraints. In *Conference on Empirical*
83 *Methods in Natural Language Processing (EMNLP)*.