

Bias in Word Embeddings

Ankit Pant 2018201035

Tarun Mohandas 2018201008

Team: *The Lost Linguists*

Outline

1. Introduction
2. Literature Review
3. Experimental Setup
4. Experiments and Results
5. Debiasing Word Embeddings
6. Proposed Algorithm
7. Conclusion
8. Future Scope
9. References

Introduction

Introduction

- Word Embeddings (WE) - extensively used in NLP
- WE inherently contain various types of biases
- NLP models tend to amplify the biases
- Hence removing bias from WE increasingly crucial
- Project attempts to use existing debiasing techniques
- Shortcomings are highlighted and alternative proposed

Literature Review

Word Embeddings

- Word Embeddings
 - Most popular representation of document vocabulary
 - Captures context of words
 - Various models include *Word2Vec* and *GLoVe*
 - Both use neural network to form word representations

Types of Biases in Datasets

- Historical Bias:
 - Unwanted biases that were present in society years ago
- Representational Bias:
 - Certain parts of the input space are under-represented
- Measurement Bias:
 - Imperfect measuring of the data
 - assuming that the measured data is proxy of some other desired feature
- Aggregation Bias:
 - Same model is used for groups with different conditional distributions
- Evaluation Bias:
 - Evaluation and benchmark data for model does not represent actual target population

Social Bias in Datasets

- Gender Bias:
 - Biased due to associating stereotype to gender
 - Man is to computer programmer as woman is to homemaker
- Racial Bias:
 - Biased due to associating stereotype to race
 - Modern is to American as medieval is to Indian
- Religious Bias:
 - Biased due to associating stereotype to religion
 - Smart is to Christian as cheapskate is to Jew

Existing Debiasing Methods

- Post-processing debiasing (Bolukbasi et al.):
 - Make change to word vector to reduce encoded gender bias
 - Done by zeroing the gender projection of each word on a predefined gender direction
 - $\vec{w} = \frac{(\vec{w} - \vec{w}_b)}{\|(\vec{w} - \vec{w}_b)\|}$
 - Ensure all neutral words are equally close to the two words
- Train word embeddings from scratch (Zhao et al.)
 - Alter the loss of GloVe model
 - Concentrate most gender information to last coordinate of each vector
 - Use word representation excluding the gender coordinate
 - Representation of gender neutral words orthogonal to gender direction

Experimental Setup

Experimental Setup

- Various kinds of biases were identified in word embeddings
- Graphically projected to check clustering of biased words
- Applying debiasing techniques
- Check result of debiasing quantitatively and graphically

Experiments and Results

Identifying Bias

- **Gender Bias**
 - manager is to man as mastercard is to woman
 - programmer is to man as nutritionist is to woman
- **Racial Bias**
 - filthy is to black as elitists is to white
 - manager is to American as barber is to Indian
- **Religious Bias**
 - rich is to Christian as homeless is to Muslim
 - smart is to Christian as cheapskate is to Jew

Most biased words - Gender i

Bias in Professions- descending order (top 10) (man-woman)

('magician', 0.11574505)

('carpenter', 0.046477903)

('butcher', 0.035740647)

('gamer', 0.027397368)

('soldier', 0.018920997)

('servant', 0.00930707)

('barber', 0.007335724)

('engineer', -0.0022402927)

('player', -0.01447518)

('programmer', -0.016492786)

Most biased words - Gender ii

Bias in Misc. words - descending order (top 10) (man-woman)

('swear', 0.21706516)

('filthy', 0.20151067)

('sweet', 0.19195326)

('roar', 0.18197575)

('weep', 0.18006238)

('pretty', 0.17052722)

('beautiful', 0.15033276)

('think', 0.14811862)

('brave', 0.14424863)

('savage', 0.14332578)

Most biased words - Race i

Bias in Professions- descending order (top 10)
(Indian-American)

('butcher', 0.30157873)

('barber', 0.20129806)

('florist', 0.18413721)

('nurse', 0.15209064)

('vet', 0.14899719)

('cashier', 0.1279751)

('tutor', 0.12523493)

('waiter', 0.11155586)

('chemist', 0.108406395)

('shopkeeper', 0.105821185)

Most biased words - Race ii

Bias in Misc. words - descending order (top 10)
(Indian-American)

('lovely', 0.2813881)

('backpack', 0.25974354)

('bag', 0.20628145)

('cute', 0.19745344)

('purse', 0.19460957)

('sweet', 0.18077426)

('beautiful', 0.1678942)

('grumpy', 0.14973086)

('cool', 0.14963488)

('sassy', 0.14638355)

Most biased words - Religion i

Bias in Professions- descending order (top 10)
(Muslim-Christian)

('vet', 0.12283135)

('cop', 0.11075719)

('soldier', 0.076730676)

('shopkeeper', 0.060545065)

('worker', 0.018872034)

('cashier', 0.016069816)

('butcher', -0.005377178)

('landlord', -0.00816943)

('nurse', -0.010036133)

('nanny', -0.011244484)

Most biased words - Religion ii

Bias in Misc. words - descending order (top 10)
(Muslim-Christian)

('terrorist', 0.2741907)

('car', 0.10418006)

('fear', 0.10198632)

('bag', 0.07560906)

('creep', 0.06417281)

('brute', 0.057166893)

('yell', 0.041089058)

('poor', 0.04030589)

('roar', 0.039248332)

('mad', 0.038921878)

Clustering on Profession

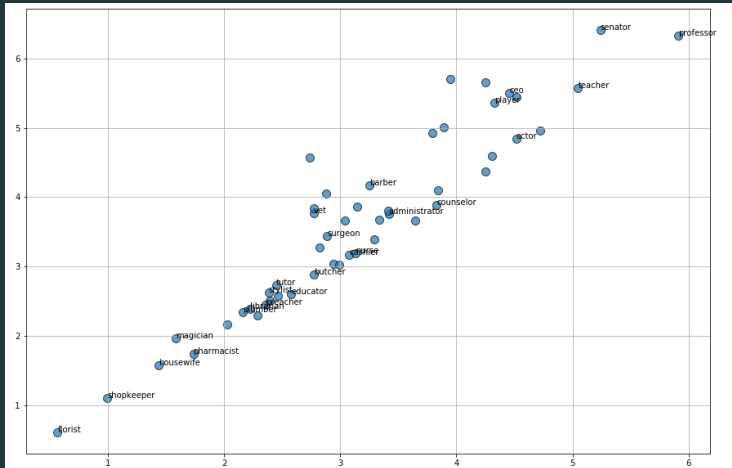


Figure 1: Result of clustering on male and female biased words

Clustering on Common Words

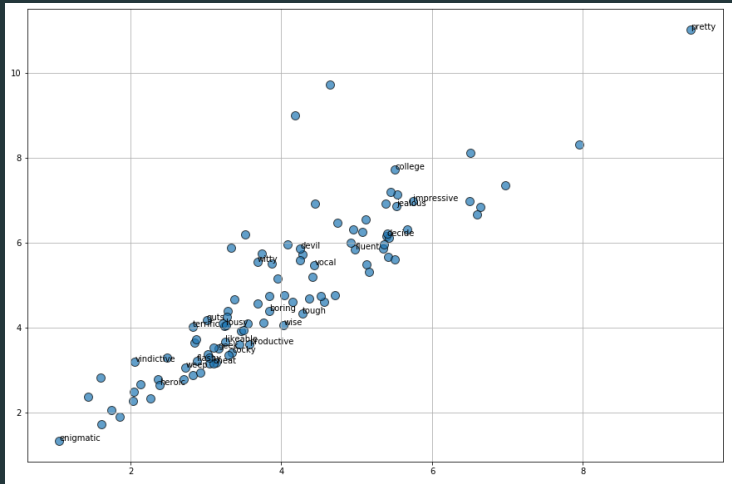


Figure 2: Result of clustering on male and female biased words

Debiasing Word Embeddings

Debiasing using neutralisation of hard-debiasing

Debiasing Professions – descending order

('magician', 0.11574504) ==> ('magician', 0.11456225)

('carpenter', 0.046477906) ==> ('carpenter', 0.045484226)

('gamer', 0.027397364) ==> ('gamer', 0.026887717)

('servant', 0.009307076) ==> ('servant', 0.008474963)

('barber', 0.0073357197) ==> ('barber', 0.006497547)

Debiasing Misc. words – descending order

('swear', 0.21706519) ==> ('swear', 0.21250035)

('filthy', 0.20151068) ==> ('filthy', 0.19604209)

('sweet', 0.19195326) ==> ('sweet', 0.18784045)

('roar', 0.18197575) ==> ('roar', 0.17574164)

('weep', 0.18006237) ==> ('weep', 0.17289506)

Clustering on Profession after Debiasing

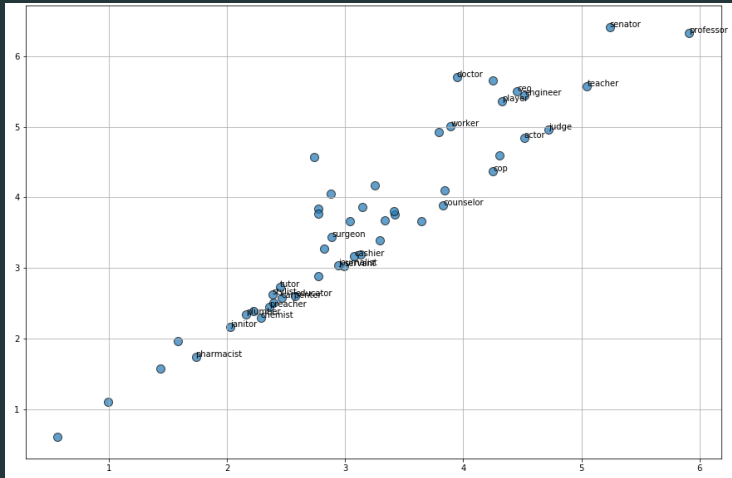


Figure 3: Result of clustering on male and female biased words

Clustering on Common Words after Debiasing

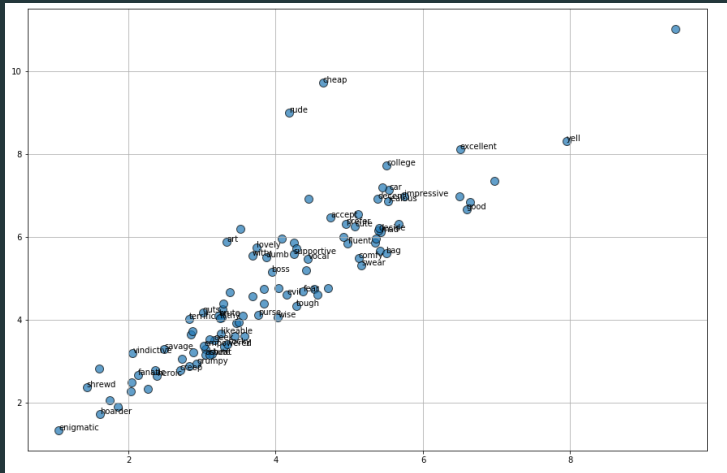


Figure 4: Result of clustering on male and female biased words

Proposed Algorithm

Proposed algorithm

- Pre-processing dataset → add sentences by inverting the gender, races, religion, etc. → appending them to create a new dataset.
- Debiasing during pre-processing as with (Zhao et al.)
- Analyse 'biased' words → projecting them to biased sub-space and debias (Bolukbasi et al.)
- After identifying these words → correlated with each other → distances between the words clustering together should also be equalised.

Conclusion

Conclusion

- Project involved exploring various biases in word embeddings
- Word embeddings trained on the Reddit dataset were explored
- Attempt to debias using the traditional hard-debiasing method was done
- Debiasing happened on a superficial level
- Approach proposed that may address the issue

Future Scope

Future Scope

- Implement proposed algorithm
- Empirically test performance of proposed model
- Explore how biases are encoded in word embeddings

References

References i



Hila Gonen, Yoav Goldberg, *Lipstick on a Pig: Debiasing methods cover up systematic gender biases in word embeddings but do not remove them,*

<https://arxiv.org/pdf/1903.03862.pdf>



Tolga Bolukbasi, et al., *Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings,*




<https://arxiv.org/pdf/1607.06520.pdf>



Jieyu Zhao, et al., *Learning Gender-Neutral Word Embeddings,*

<https://arxiv.org/pdf/1809.01496.pdf>

References ii

-  Sevtap Duman, et al., *(Visualization of) gender bias in word embeddings*, <http://wordbias.umiacs.umd.edu/>
-  Aylin Caliskan, et al., *Semantics derived automatically from language corpora contain human-like biases*, <https://arxiv.org/pdf/1608.07187.pdf>
-  Nikhil Garg, et al., *Word embeddings quantify 100 years of gender and ethnic stereotypes*, <https://arxiv.org/pdf/1711.08412.pdf>

References iii



Thomas Manzini, et al., *Black is to Criminal as Caucasian is to Police: Detecting and Removing Multiclass Bias in Word Embeddings*,

<https://arxiv.org/pdf/1904.04047v1.pdf>



Harini Suresh, John V. Gutttag, *A Framework for Understanding Unintended Consequences of Machine Learning*, <https://arxiv.org/pdf/1901.10002.pdf>



Alex Fefegha, *Racial Bias and Gender Bias Examples in AI systems*,

<https://peopleofcolorintech.com/articles/racial-bias-and-gender-bias-examples-in-ai-systems>

References iv



Marianne Bertrand, Sendhil Mullainathan, *Are Emily and Greg More Employable than Lakisha and Jamal? A Field Experiment on LaborMarket Discrimination*,
https://www2.econ.iastate.edu/classes/econ321/Orazem/bertrand_emily.pdf



Thomas Manzini, et al., *Debiasing Multiclass Word Embeddings*, <https://github.com/TManzini/DebiasMulticlassWordEmbedding>



Ella Rabinovich, Shuly Wintner, *The Reddit-L2 corpus*,
<http://cl.haifa.ac.il/projects/L2/>