



deeplearning.ai

NLP and Word Embeddings

Debiasing word embeddings

The problem of bias in word embeddings

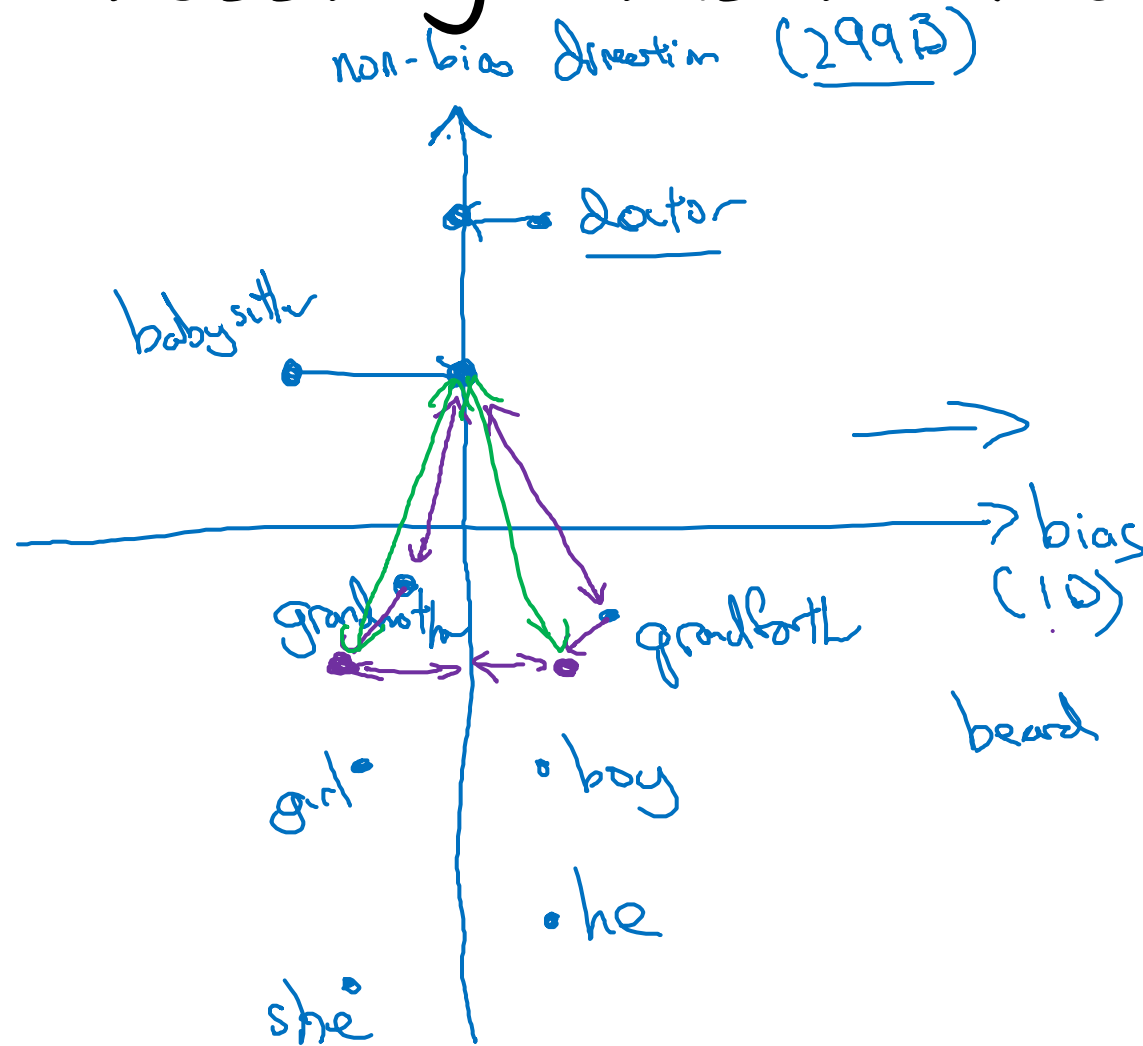
Man:Woman as King:Queen

Man:Computer_Programmer as Woman:homemaker X

Father:Doctor as Mother:nurse X

Word embeddings can reflect gender, ethnicity, age, sexual orientation, and other biases of the text used to train the model.

Addressing bias in word embeddings



1. Identify bias direction.

$\{ \begin{aligned} &e_{he} - e_{she} \\ &e_{male} - e_{female} \\ &\vdots \end{aligned} \}$
→ average

2. Neutralize: For every word that is not definitional, project to get rid of bias.

3. Equalize pairs.

→ $\left. \begin{array}{cc} \text{grandmother} & \text{grandfather} \\ \text{girl} & \text{boy} \end{array} \right\}$