# The problem of bias in word embeddings
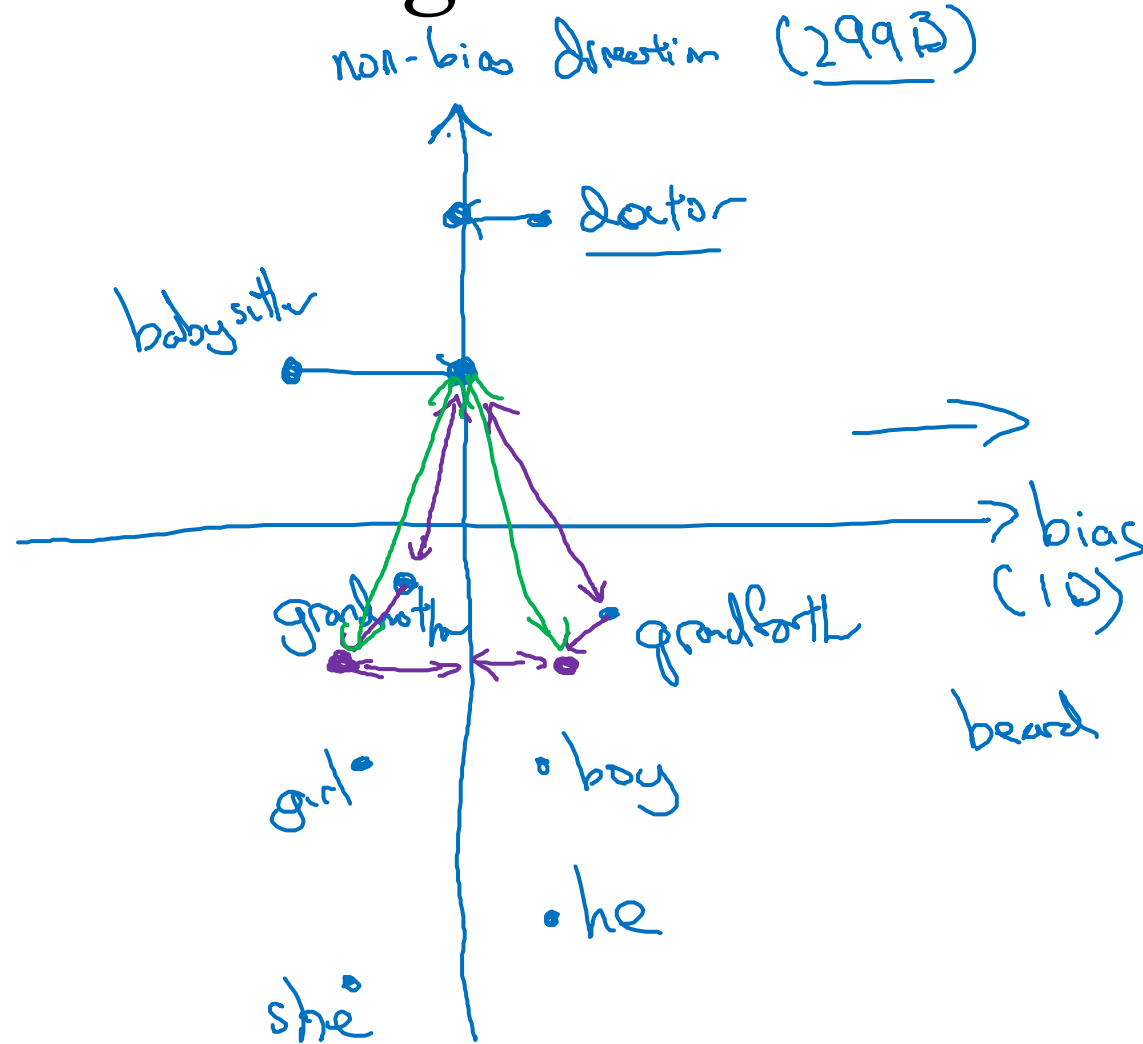
Man:Woman as King:Queen

Man:Computer_Programmer as Woman:Homemaker ✗

Father:Doctor as Mother:Nurse ✗

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

[Bolukbasi et. al., 2016. Man is to computer programmer as woman is to homemaker? Debiasing word embeddings]    Andrew Ng

# Addressing bias in word embeddings

non-bias direction (299D)

doctor

babysitter

bias (1D)

grandmoth    grandfath

beard

girl    boy

he

she

1. Identify bias direction.

$$e_{he} - e_{she}$$
$$e_{male} - e_{female}$$
$$\vdots$$

average

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

3. Equalize pairs.

grandmoth — grandfath
girl            boy

[Bolukbasi et. al., 2016. Man is to computer programmer as woman is to homemaker? Debiasing word embeddings]

Andrew Ng