

Word representation

V = [a, aaron, ..., zulu, <UNK>]

$$|V| = 10,000$$

1-hot representation

[illegible]

I want a glass of orange juice.
I want a glass of apple juice.

Andrew Ng

Featurized representation: word embedding

| | Man (5391) | Woman (9853) | King (4914) | Queen (7157) | Apple (456) | Orange (6257) |
|--------|---------------|-----------------|----------------|-----------------|----------------|------------------|
| Gender | -1 | 1 | -0.95 | 0.97 | 0.00 | 0.01 |
| Royal | 0.01 | 0.02 | <u>0.93</u> | <u>0.95</u> | -0.01 | 0.00 |
| Age | 0.03 | 0.02 | 0.7 | 0.69 | 0.03 | -0.02 |
| Food | 0.04 | 0.01 | 0.02 | 0.01 | 0.95 | 0.97 |
| size | ⋮ | ⋮ | | | | |
| cost | | | | | | |
| alive | | | | | | |
| verb | | | | | | |

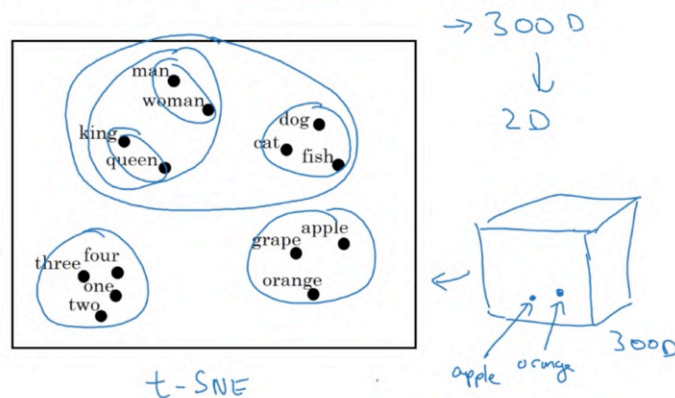
I want a glass of orange juice

I want a glass of apple juice

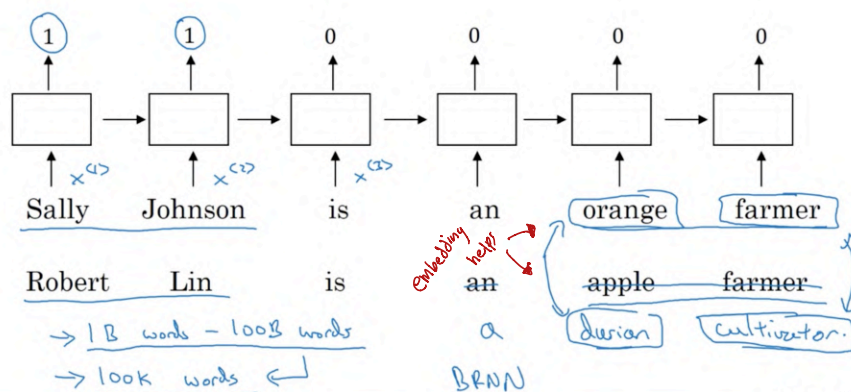
Andrew

Andrew Ng

Visualizing word embeddings



Named entity recognition example



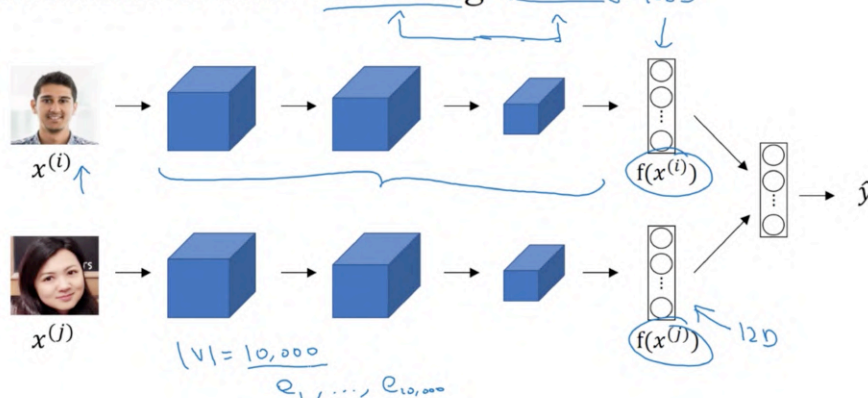
Andrew Ng

Transfer learning and word embeddings

1. Learn word embeddings from large text corpus. (1-100B words)
(Or download pre-trained embedding online.)
2. Transfer embedding to new task with smaller training set.
(say, 100k words) $\rightarrow 10,000 \rightarrow 300$
3. Optional: Continue to finetune the word embeddings with new data.

Andrew Ng

Relation to face encoding (embedding)



[Taigman et. al., 2014. DeepFace: Closing the gap to human level performance]

Andrew Ng

Analogies

| | Man (5391) | Woman (9853) | King (4914) | Queen (7157) | Apple (456) | Orange (6257) |
|--------|---------------|-----------------|----------------|-----------------|----------------|------------------|
| Gender | -1 | 1 | -0.95 | 0.97 | 0.00 | 0.01 |
| Royal | 0.01 | 0.02 | 0.93 | 0.95 | -0.01 | 0.00 |
| Age | 0.03 | 0.02 | 0.70 | 0.69 | 0.03 | -0.02 |
| Food | 0.09 | 0.01 | 0.02 | 0.01 | 0.95 | 0.97 |

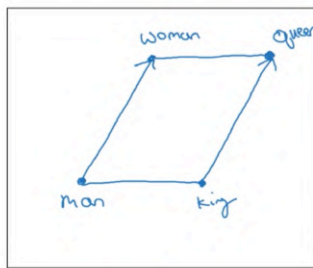
$e_{\text{man}} - e_{\text{woman}} \approx \begin{bmatrix} -2 \\ 0.01 \\ 0.01 \\ 0.08 \end{bmatrix}$
 $e_{\text{king}} - e_{\text{queen}} \approx \begin{bmatrix} -2 \\ 0.02 \\ 0.01 \\ 0.01 \end{bmatrix}$
 $e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - e_{\text{queen}}$

$\text{Man} \rightarrow \text{Woman} \approx \text{King} \rightarrow ? \text{ Queen}$

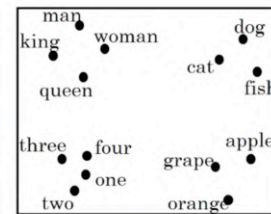
[Mikolov et. al., 2013, Linguistic regularities in continuous space word representations]

Andrew Ng

Analogies using word vectors



3000 → 20



t-SNE

$$e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - e_{\text{queen}}$$

300D

Find word w : $\arg \max_w \text{sim}(e_w, e_{\text{king}} - e_{\text{man}} + e_{\text{woman}})$

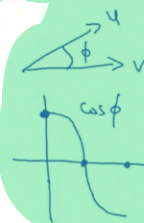
30-75%

Andrew Ng

Cosine similarity

$$\rightarrow \text{sim}(e_w, e_{\text{king}} - e_{\text{man}} + e_{\text{woman}})$$

$$\text{sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2}$$

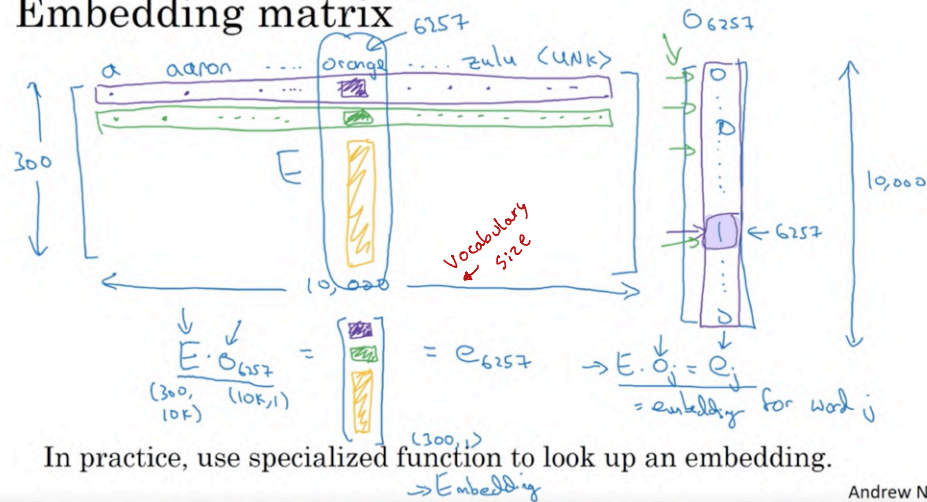


$$\|u - v\|^2$$

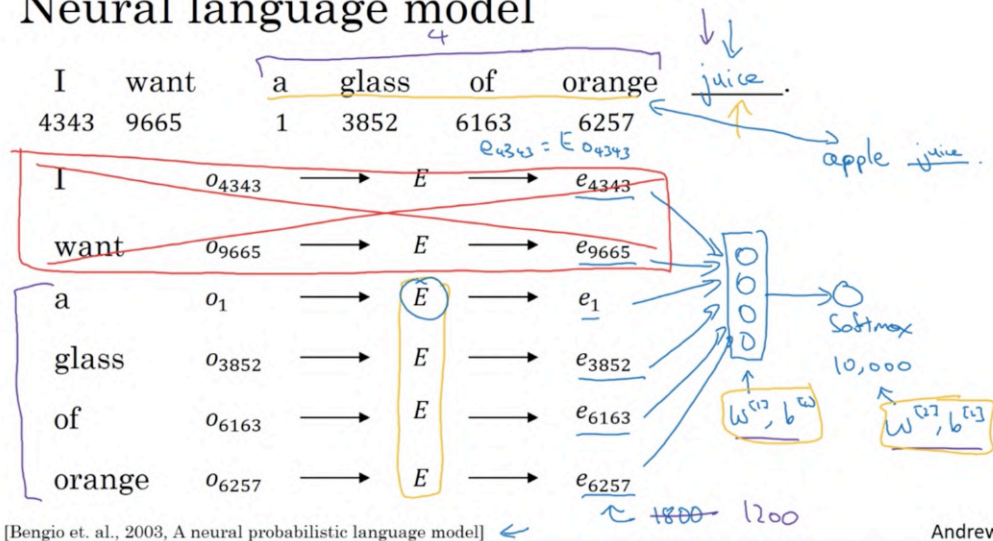
Man:Woman as Boy:Girl
 Ottawa:Canada as Nairobi:Kenya
 Big:Bigger as Tall:Taller
 Yen:Japan as Ruble:Russia

Andrew Ng

Embedding matrix



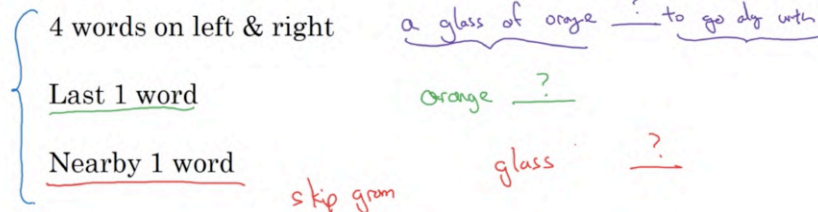
Neural language model



Other context/target pairs

I want a glass of orange juice to go along with my cereal.

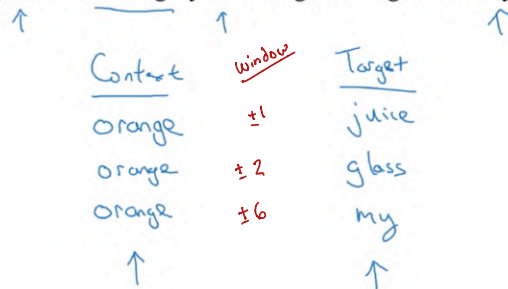
Context: Last 4 words.



Skip-grams

Randomly pick a word within a window

I want a glass of orange juice to go along with my cereal.

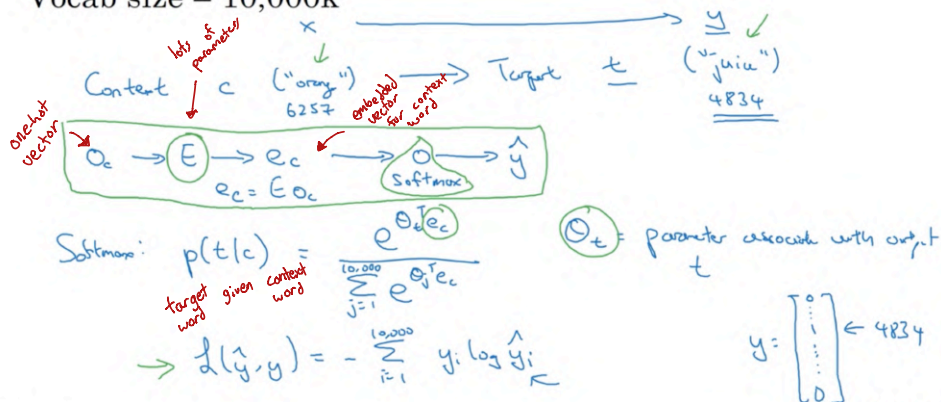


[Mikolov et. al., 2013. Efficient estimation of word representations in vector space.]

Andrew Ng

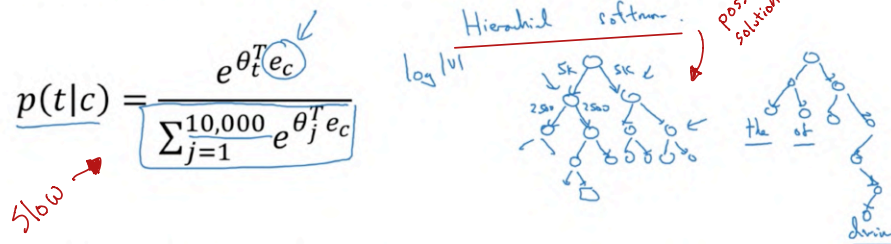
Model

Vocab size = 10,000k



Andrew Ng

Problems with softmax classification



How to sample the context c ?

\rightarrow the, of, and, to, ...
 \rightarrow orange, apple, durian

t
 $c \rightarrow t$
 $P(c)$

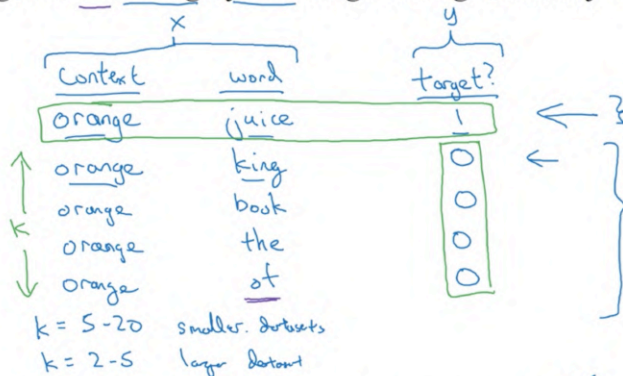
Andrew Ng

Defining a new learning problem

Negative Sampling

I want a glass of orange juice to go along with my cereal.

It's really to try to distinguish between these two types of distributions from which you might sample a pair of words.



[Mikolov et. al., 2013. Distributed representation of words and phrases and their compositionality]

Andrew Ng

Model

Softmax:
$$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$$

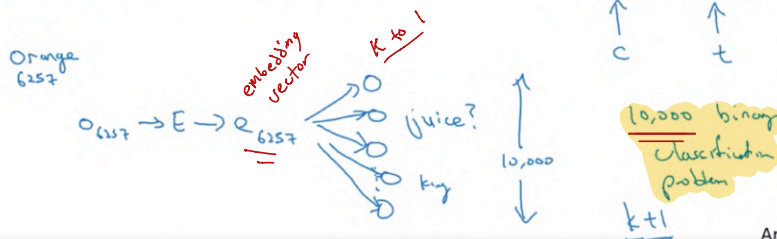
10,000-way softmax

This technique is called negative sampling because what you're doing is, you have a positive example, the orange and then juice. And then you will go and deliberately generate a bunch of negative examples, negative samplings, hence, the name negative sampling, with which to train four more of these binary classifiers.

Play video starting at 8 minutes 54 seconds and follow transcript 18:54

And on every iteration, you choose four different random negative words with which to train your algorithm on.

$$P(y=1 | c, t) = \sigma(\theta_c^T e_t)$$



Andrew Ng

Selecting negative examples

| context | word | target? |
|---------|-------|---------|
| orange | juice | 1 |
| orange | king | 0 |
| orange | book | 0 |
| orange | the | 0 |
| orange | of | 0 |

the, of, and, ...

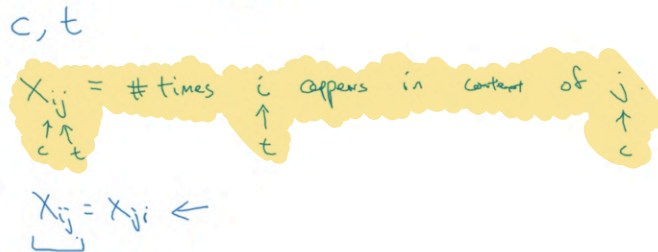
$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10,000} f(w_j)^{3/4}}$$

$$\frac{1}{|V|}$$

Andrew Ng

GloVe (global vectors for word representation)

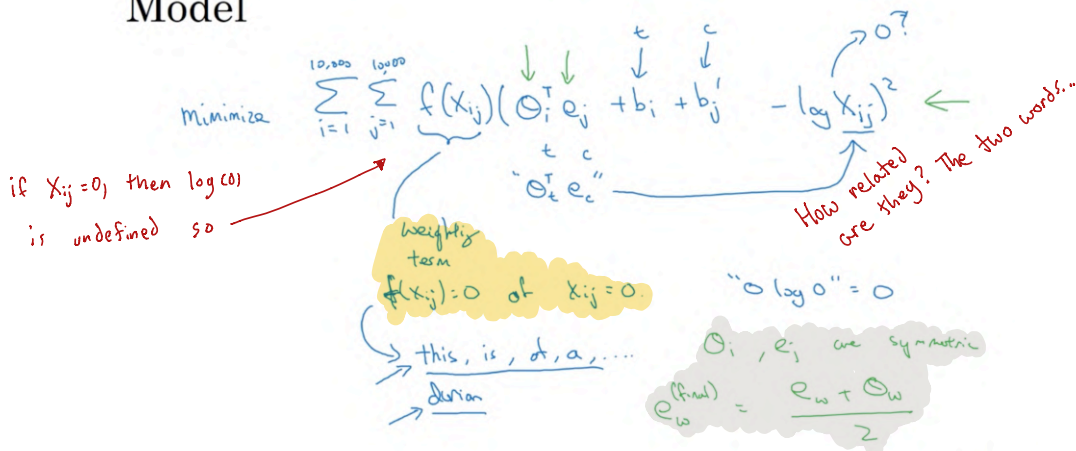
I want a glass of orange juice to go along with my cereal.



[Pennington et. al., 2014, GloVe: Global vectors for word representation]

Andrew Ng

Model



Andrew Ng

GloVe (global vectors for word representation)

I want a glass of orange juice to go along with my cereal.



[Pennington et. al., 2014, GloVe: Global vectors for word representation]

Andrew Ng

2. Correction in "GloVe word vectors" slide 4 (8/11):

typo:

$$\text{minimize } \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(X_{ij}) (\theta_i^T e_j + b_i - b'_j - \log X_{ij})^2$$

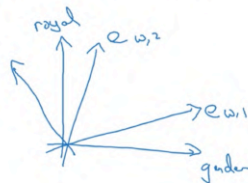
corrected:

$$\text{minimize } \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(X_{ij}) (\theta_i^T e_j + b_i + b'_j - \log X_{ij})^2$$

Andrew Ng

A note on the featurization view of word embeddings

| | Man (5391) | Woman (9853) | King (4914) | Queen (7157) |
|--------|---------------|-----------------|----------------|-----------------|
| Gender | -1 | 1 | -0.95 | 0.97 |
| Royal | 0.01 | 0.02 | 0.93 | 0.95 |
| Age | 0.03 | 0.02 | 0.70 | 0.69 |
| Food | 0.09 | 0.01 | 0.02 | 0.01 |

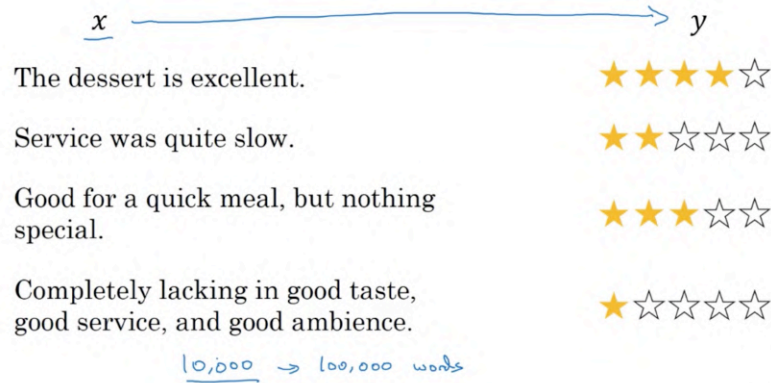


$$\text{minimize } \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(X_{ij}) (\theta_i^T e_j + b_i - b'_j - \log X_{ij})^2$$

$(A\theta_i)^T (A^T e_j) = \theta_i^T A^T A e_j$

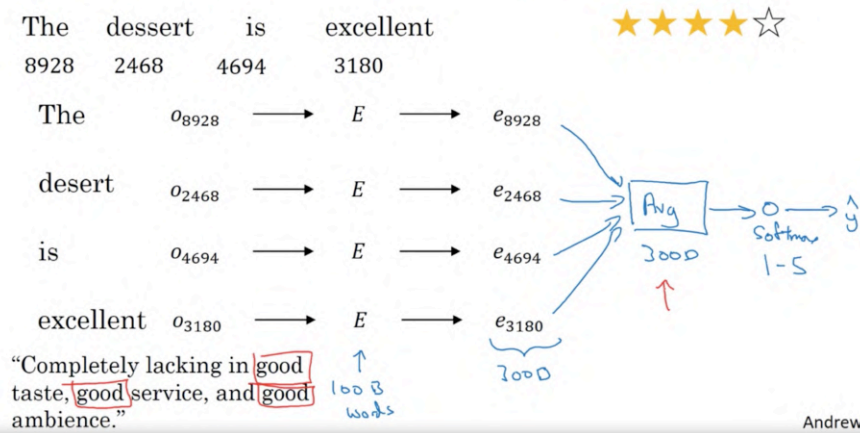
Andrew Ng

Sentiment classification problem



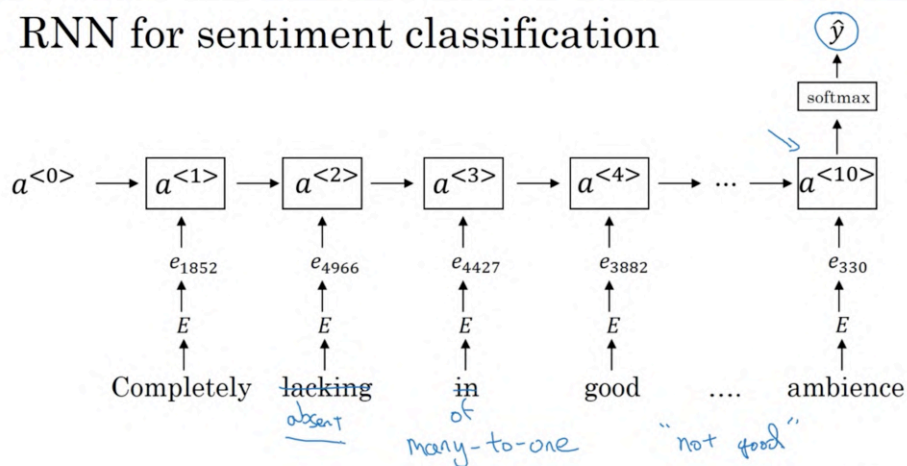
Andrew Ng

Simple sentiment classification model



Andrew Ng

RNN for sentiment classification



Andrew Ng

The problem of bias in word embeddings

Man:Woman as King:Queen

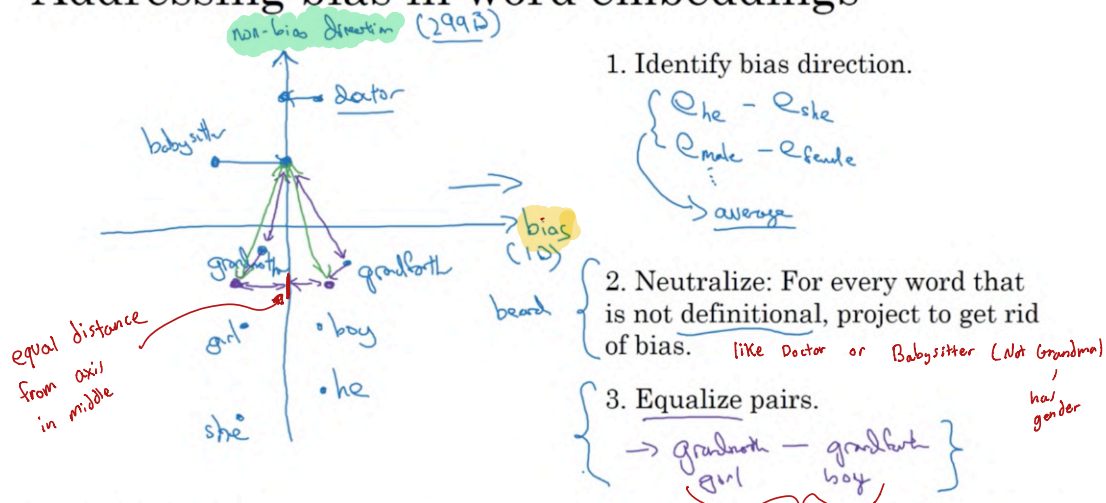
Man:Computer_Programmer as Woman:Homemaker ✗

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



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

Want only difference in their embedding to be gender
so both should be same exact similarity
or exactly same difference