



deeplearning.ai

NLP and Word Embeddings

Word representation

Word representation

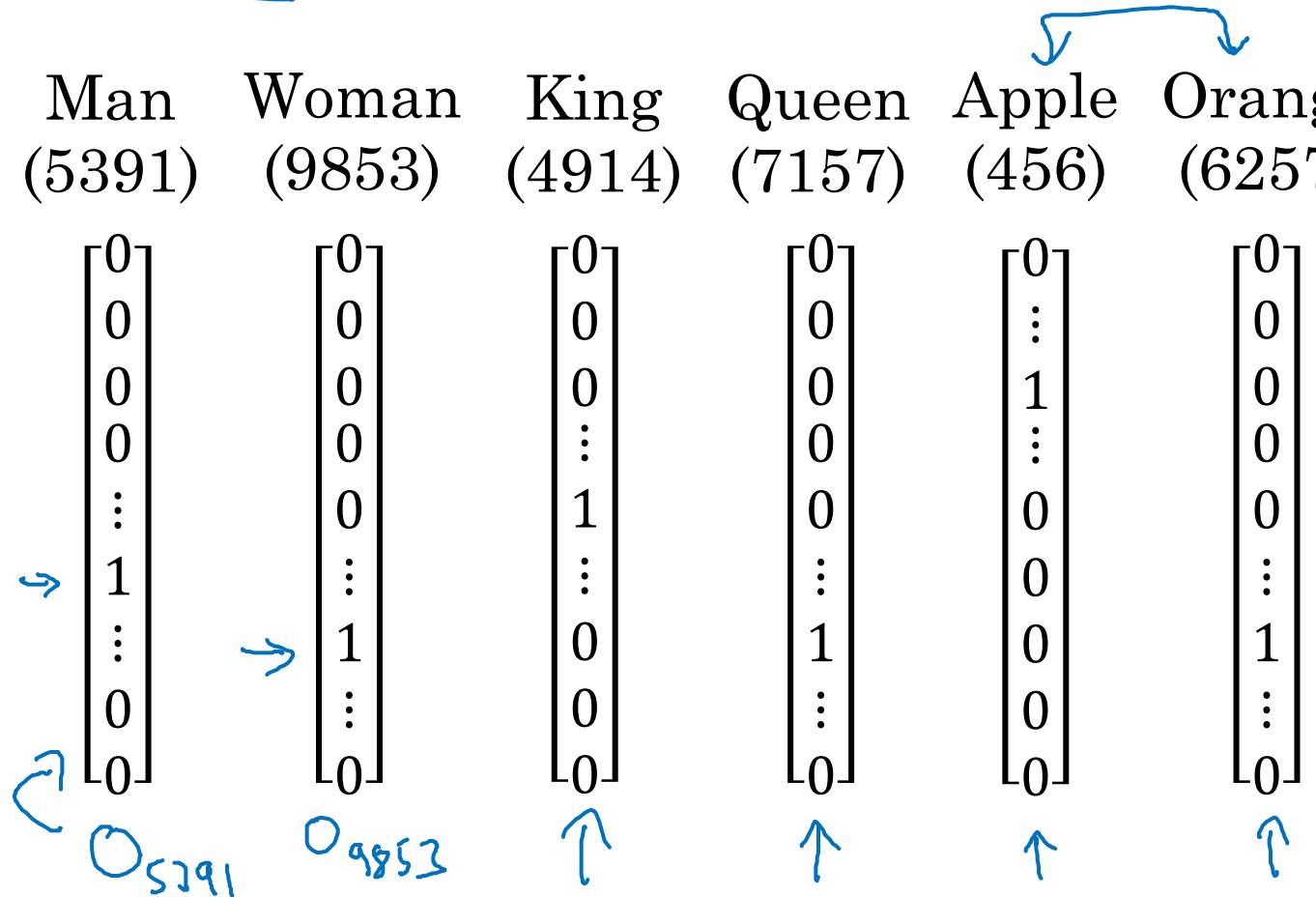
limitation of 1-hot encoding of words: it doesn't encode similarity between words.

$$V = [a, \text{aaron}, \dots, \text{zulu}, \text{<UNK>}]$$

$$|V| = 10,000$$

1-hot representation

Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
---------------	-----------------	----------------	-----------------	----------------	------------------



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

Featurized representation: word embedding

One possibility is to embed words in a feature space. This embedding requires a certain heuristic (a way to extract the feature value for each sample).

	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	0.7	0.69	0.03	-0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97
Size	:	:				
Cost						
Color						
Verb						

embedding of word 5391

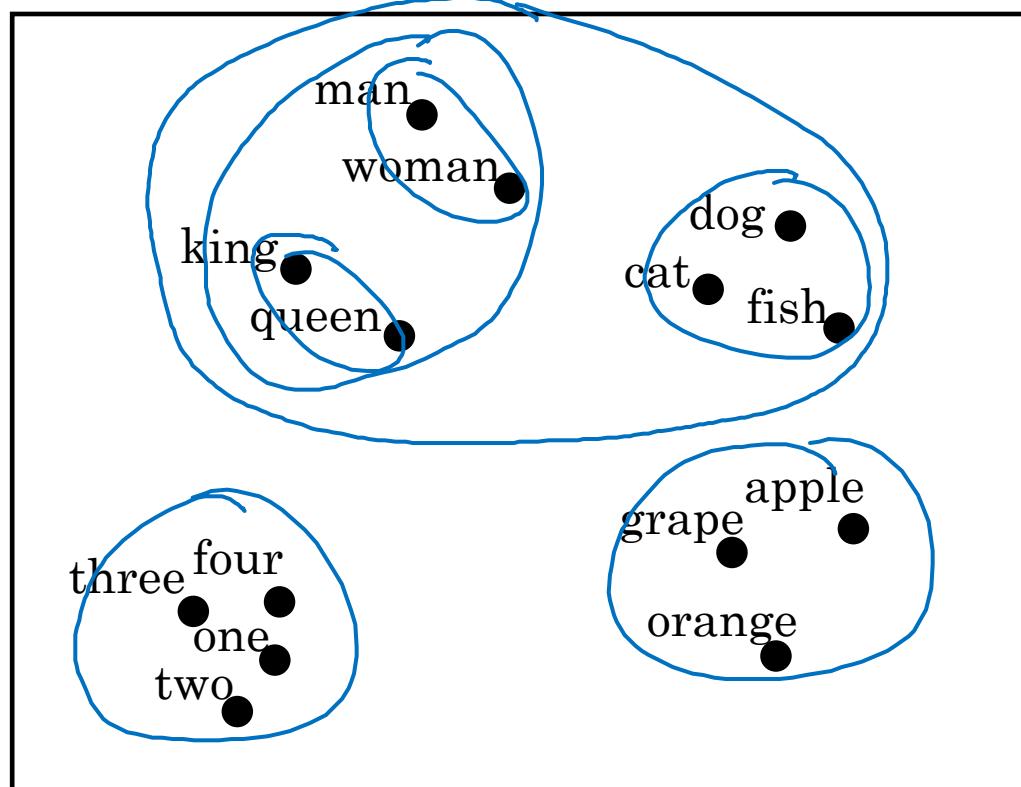
embedding of word 9853

I want a glass of orange juice.

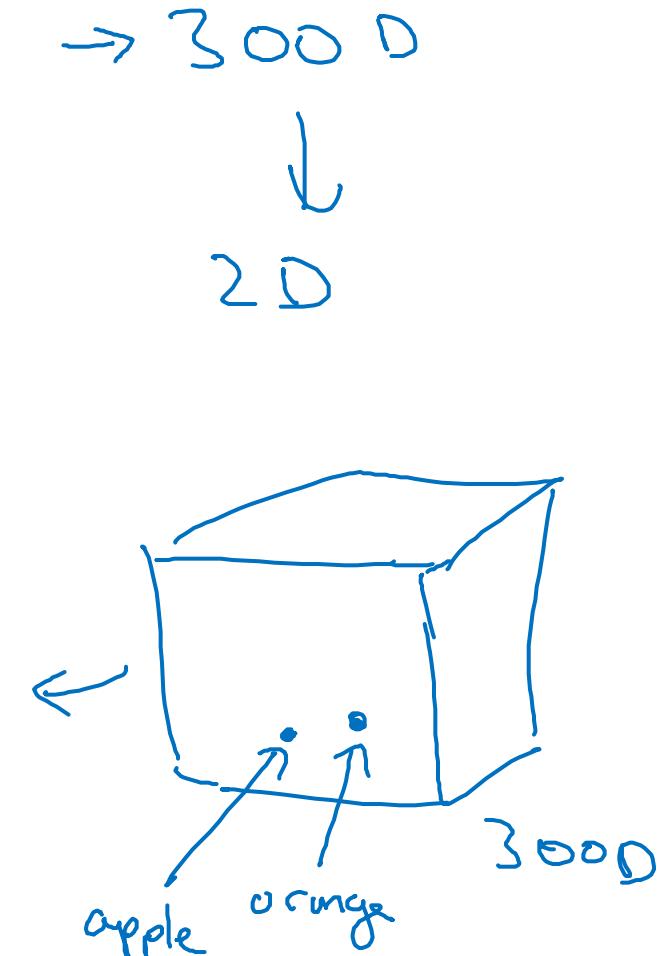
I want a glass of apple juice.

Andrew Ng

Visualizing word embeddings



t-SNE



Using a word embedding^{in a task B} implicitly corresponds to performing transfer learning from the task A of learning the embedding (usually done on very large text corpora, 1B to 100B words) to task B (where normally we have much less data).



deeplearning.ai

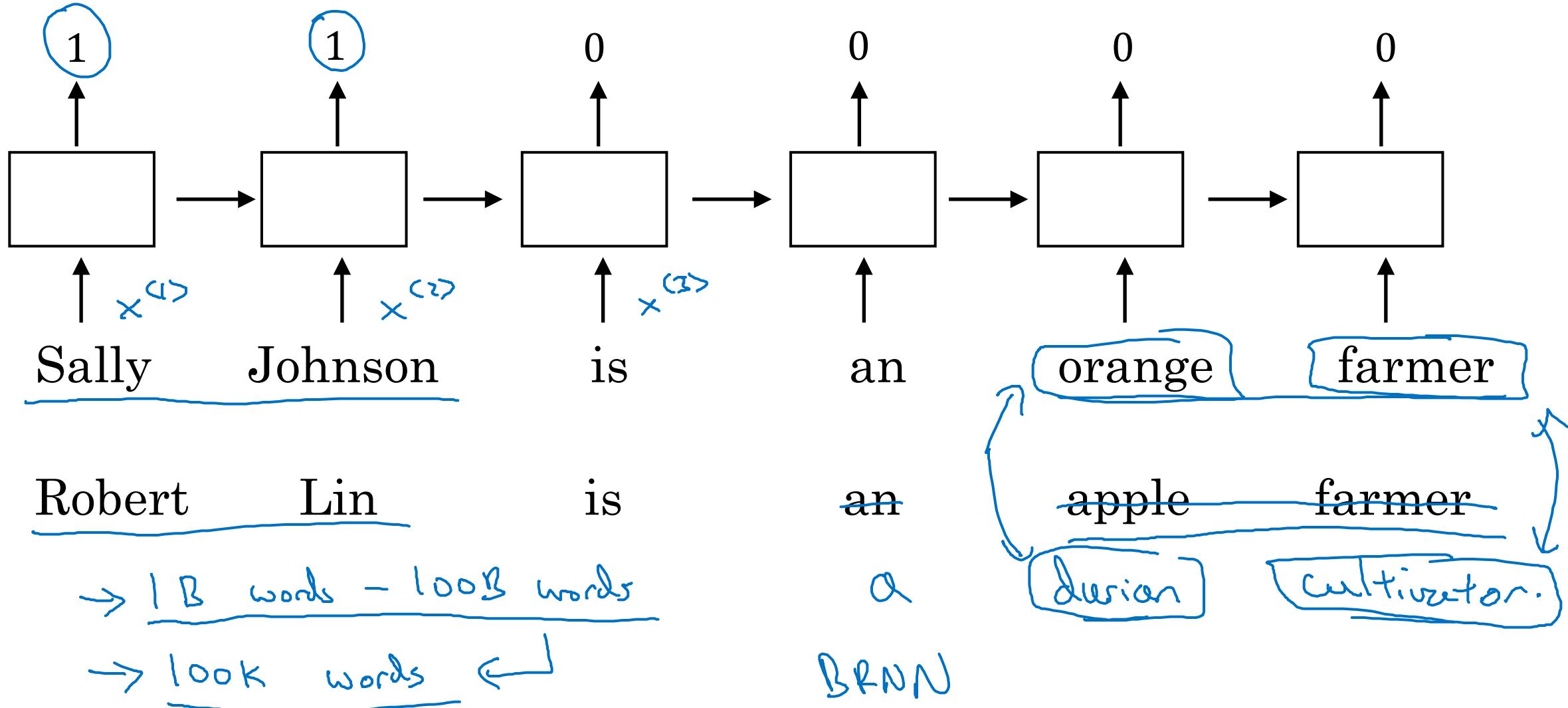
NLP and Word Embeddings

In other words,

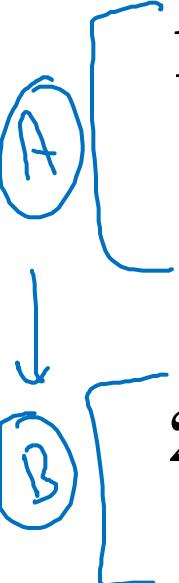
By simply encoding the words in input with a certain word embedding for a certain task, we are transferring knowledge to that task.

Using word embeddings

Named entity recognition example



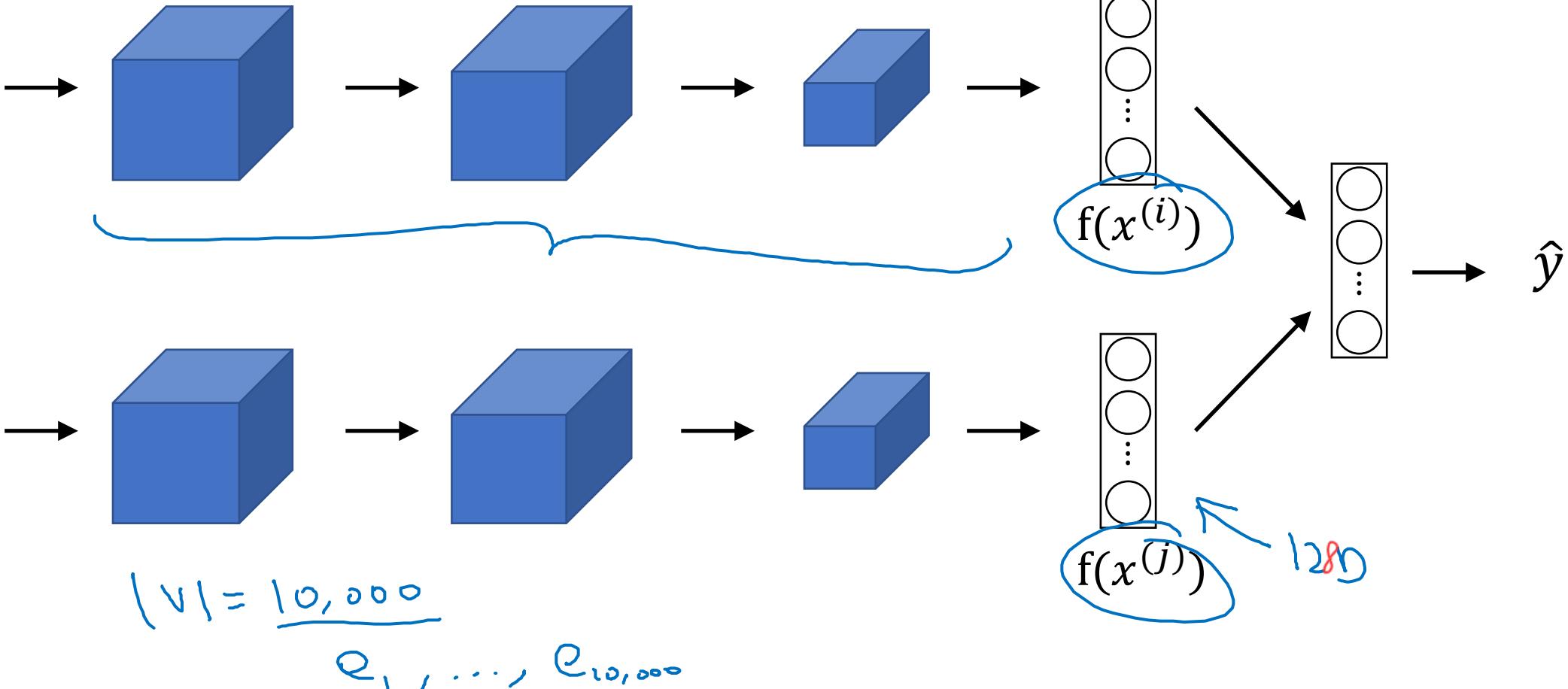
Transfer learning and word embeddings

- 
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) → 10,000 → 300
 3. Optional: Continue to finetune the word embeddings with new data.

Relation to face encoding (embedding) 128D



$x^{(i)}$



One property of word embedding is that they can learn ANALOGY relationships. That is, in the embedding space, any 2 couples of words ($A-A'$) and ($B-B'$) that are related by an analogy (men-woman)

(king-queen) have similar distance from each other.



deeplearning.ai

which allows to solve analogy problems like men:woman = King:x

Often this holds even for linear similarity
 $\text{men-woman} \approx \text{King-Queen}$

$$\begin{aligned} \text{men-woman} &= \text{King}-x \rightarrow x = \underset{x}{\operatorname{argmin}} (x-\text{King} + (\text{men-woman})) \\ &\rightarrow \underset{x}{\operatorname{argmax}} d(x, \text{King} - (\text{men-woman})). \end{aligned}$$

NLP and Word Embeddings

$$d(\text{men}, \text{woman}) \approx d(\text{king}, \text{queen})$$

$$\begin{bmatrix} 0.1 \\ 0 \\ 0.3 \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 0.11 \\ 0 \\ 0.29 \\ 0 \\ 0 \end{bmatrix}$$

Properties of word embeddings

Analogy

	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

$$\begin{matrix} e_{5391} \\ e_{\text{man}} \end{matrix}$$

$$\underline{\text{Man} \rightarrow \text{Woman}}$$

$$e_{\text{man}} - e_{\text{woman}}$$

$$e_{\text{woman}}$$

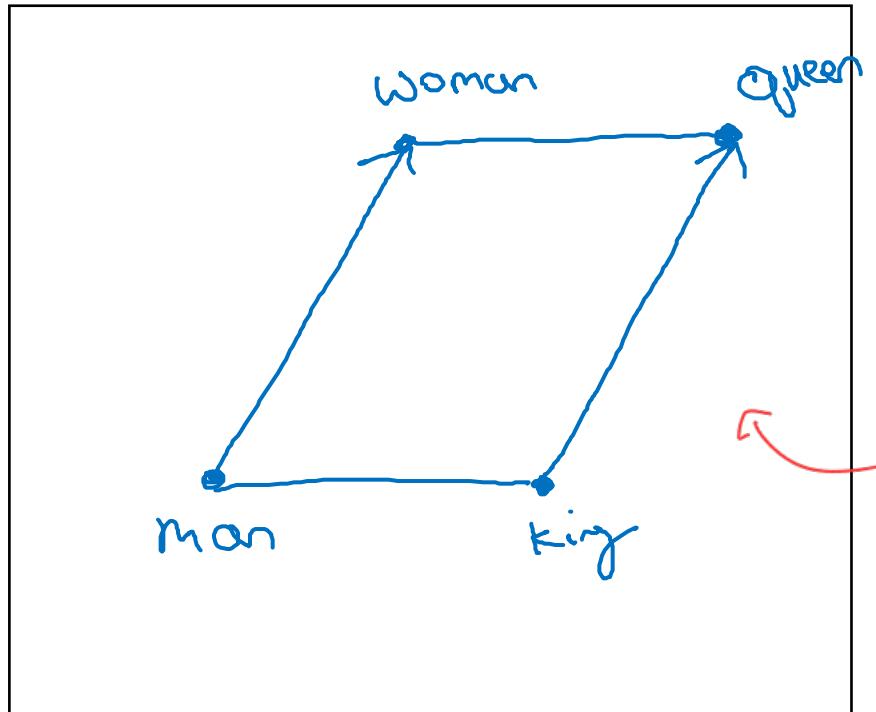
$$\underline{\text{King} \rightarrow ? \text{ Queen}}$$

$$e_{\text{king}} - e_{? \text{ Queen}}$$

$$\underline{e_{\text{man}} - e_{\text{woman}}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\underline{e_{\text{king}} - e_{\text{queen}}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Analogies using word vectors



300 D

Find word w : $\arg \max_w$

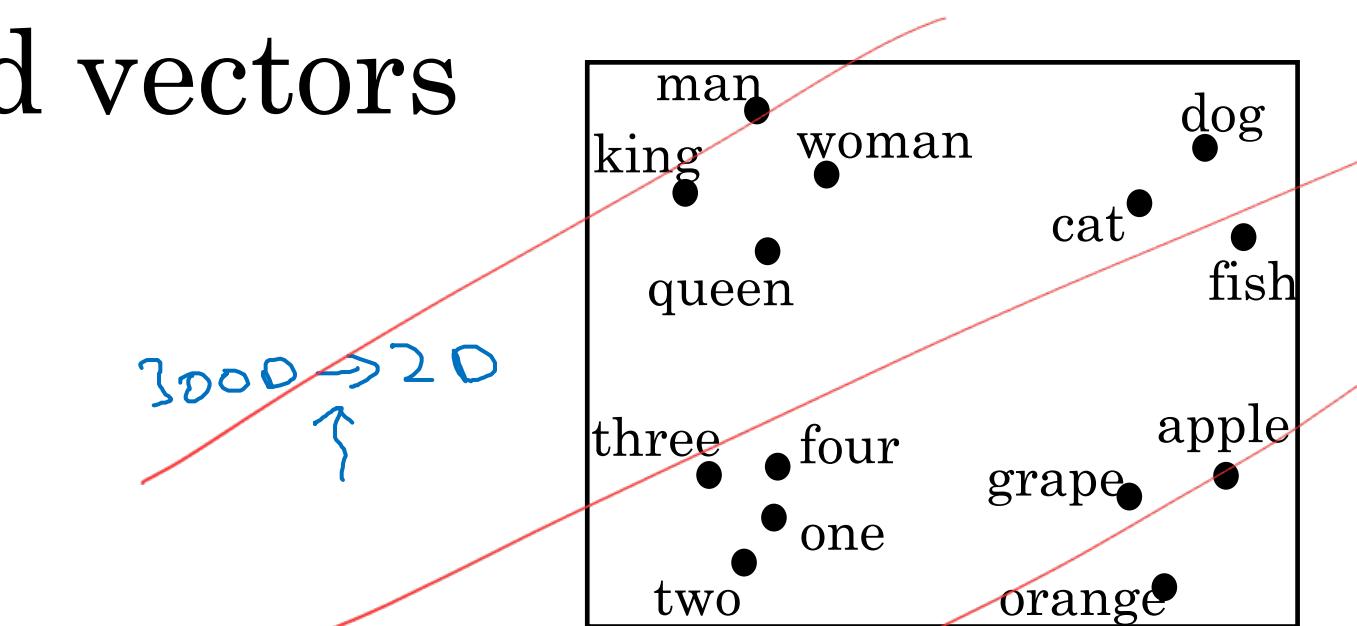
Analogy Relationship
hold in the original embedding space.

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

e_w

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

30 - 75%



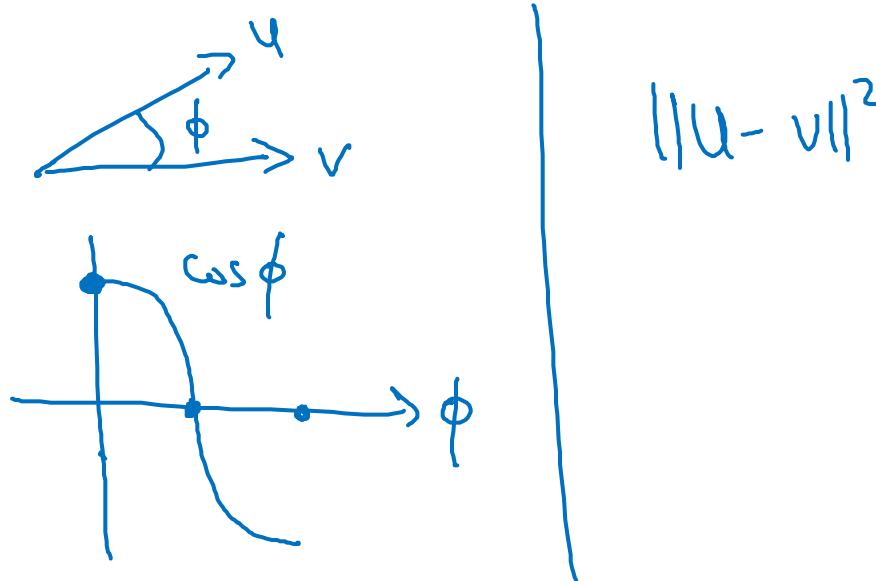
300D \rightarrow 2D

t-SNE

Cosine similarity ↪ most often used in NLP.

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

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



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

Given a vocabulary of size 10000, an embedding of that vocabulary is represented as an embedding matrix E . If the embedding is 300-dim, then $E_{300 \times 10000}$



deeplearning.ai

where each column i of E is the embedding of the word i .

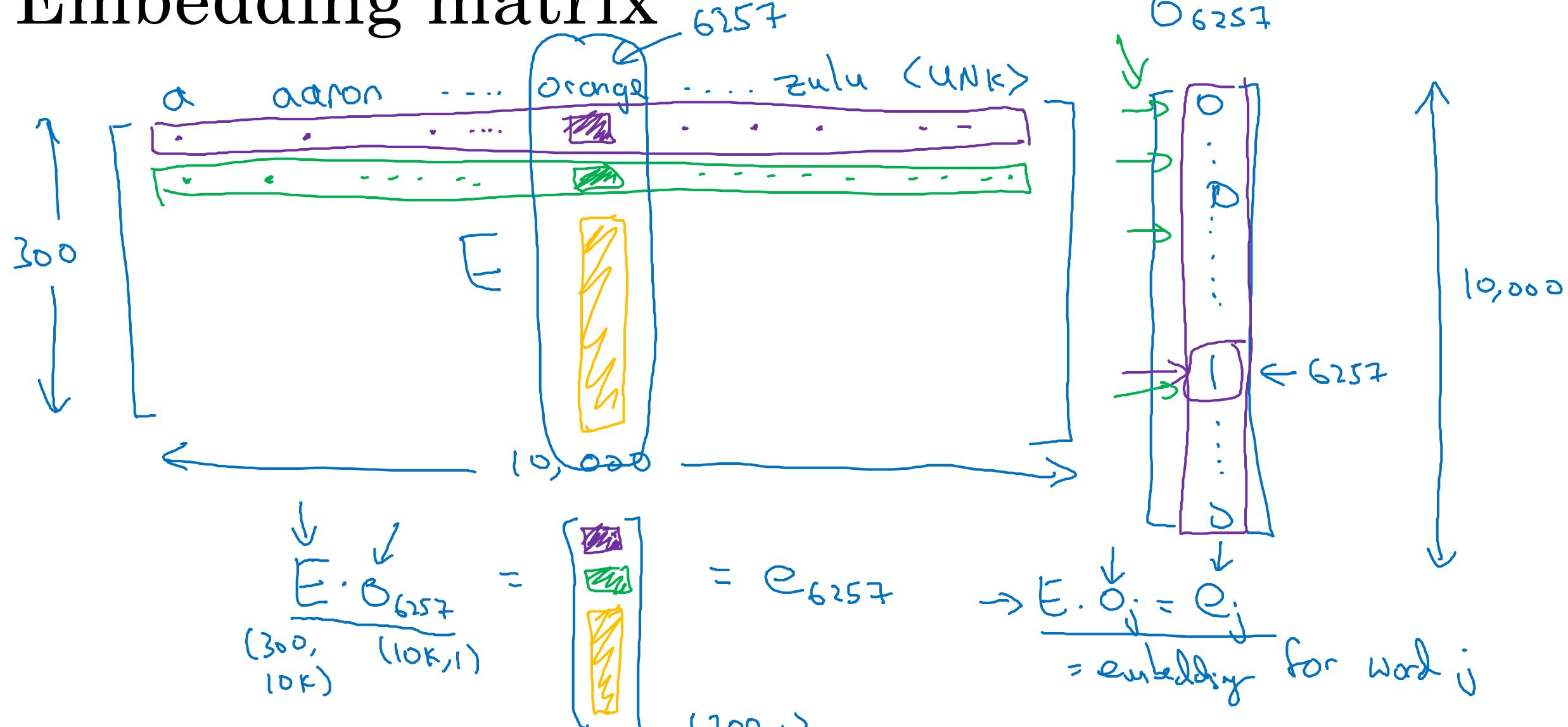
NLP and Word Embeddings

Calling \vec{e}_i and \vec{o}_i the embedding and the one-hot encoding of the word i , we have

Embedding matrix

that $\vec{e}_i = E \cdot \vec{o}_i$

Embedding matrix



In practice, use specialized function to look up an embedding.
→ Embedding

Now very efficient and simple ways exist to learn a word embedding. Ng goes through the history of these methods from old and inefficient to more recent ones.



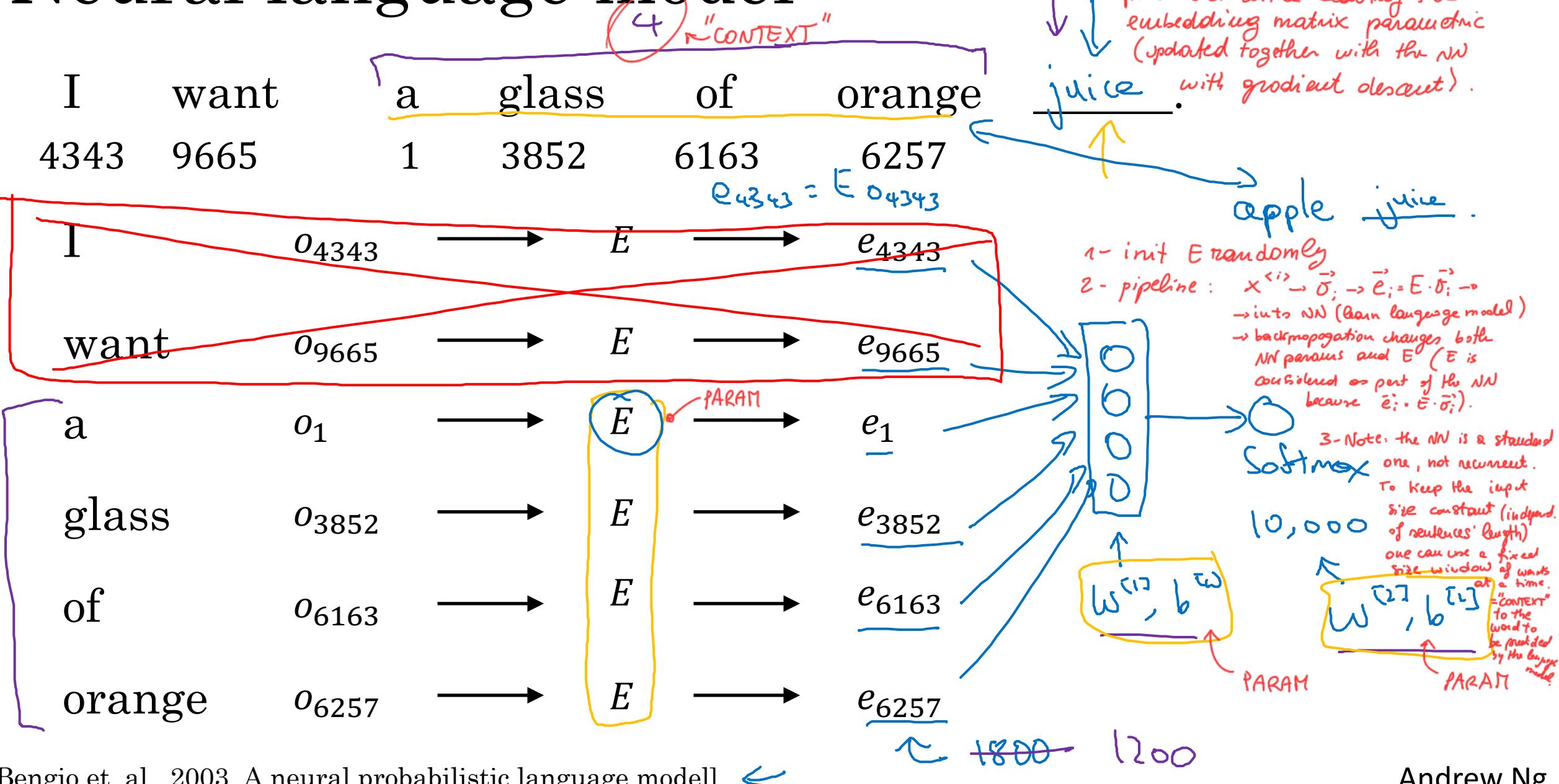
deeplearning.ai

NLP and Word Embeddings

Learning word embeddings

Neural language model

One way to learn an embedding is to use the embedding (matrix) in a language model learning procedure while leaving the embedding matrix parametric (updated together with the NN with gradient descent).



Other context/target pairs

With the aim of learning a word embedding
(and not necessarily a good language model)
other contexts can be used.

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

Context: Last 4 words.

→ 4 words on left & right

→ Last 1 word

→ Nearby 1 word (not necessarily neighboring).

skip gram

a glass of orange ? to go along with

orange ?

glass ?

An efficient word embedding, called WORD2VEC, is a word embedding built using a skip-gram model (another version is built with the continuous bag of words model. That is not discussed here).

To learn a word2vec embedding based
on the skip-gram
model:



deeplearning.ai

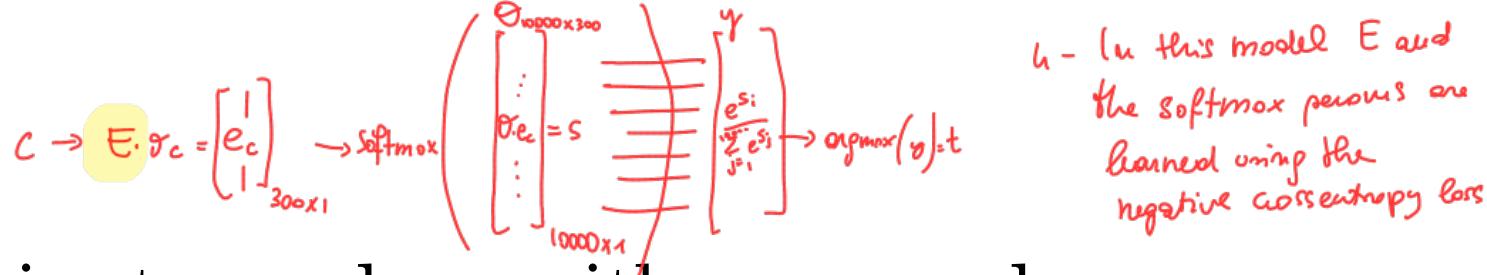
NLP and Word Embeddings

- 1- from the dataset sample pairs of words (CONTEXT, TARGET) such that TARGET is a word of the dataset and context is another word randomly sampled within n words before or after target. This sampling can use some heuristics to avoid sampling overly frequent and not informative words as context (such as articles).

Word2Vec
- 2- encode the pairs using the encoding E initialized at random (encoding matrix E initialized at random) whose entries will be updated by keeping E parametric in a skip-gram learning problem.
- 3- set a simple NN architecture to learn the skip-gram model context \rightarrow target. This architecture only consists of a softmax node preceded by E.
$$c \rightarrow \sigma_c \rightarrow e_c = E\sigma_c \rightarrow y = \text{SOFTMAX}(\Theta \cdot e_c)$$
 where $\Theta_{\text{vocab size} \times \text{encode size}}$ $\rightarrow t = \text{argmax}_f(y)$

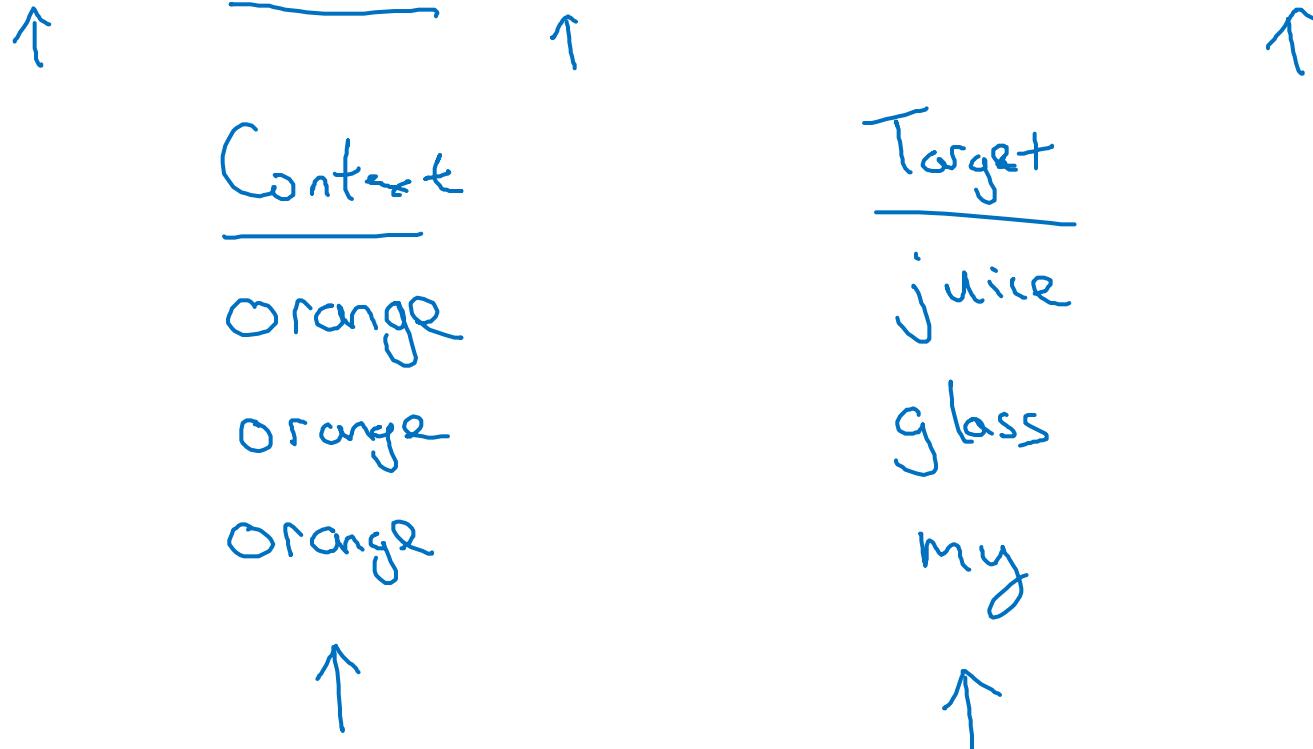
CONTINUE

Skip-grams



In this model E and the softmax neurons are learned using the negative cross-entropy loss.

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

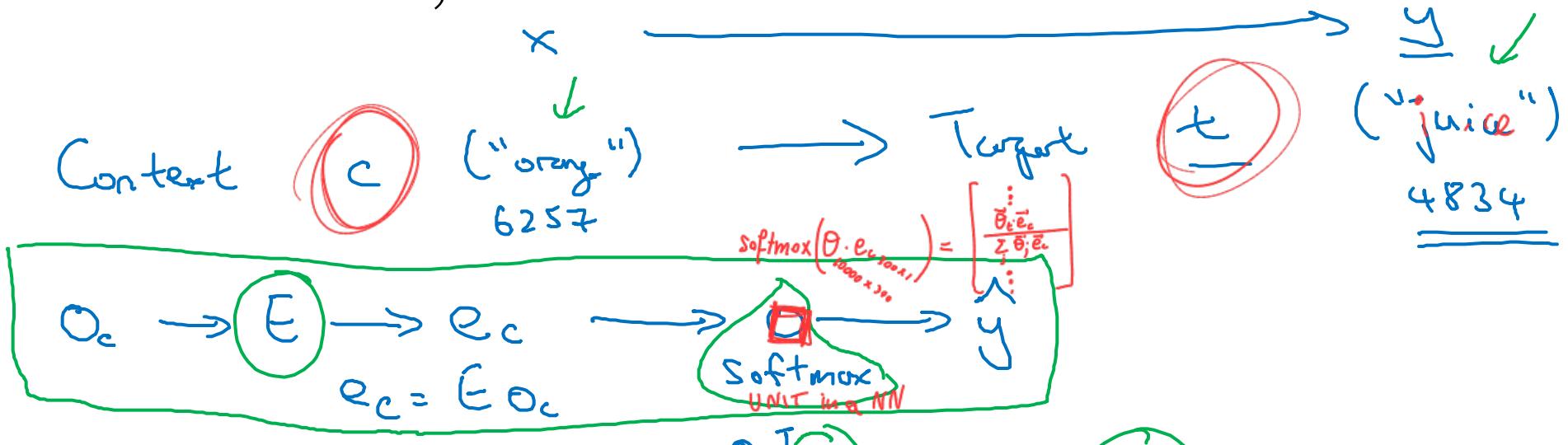


NOTE: Softmax classification has the problem of requiring vocabsize operations for each word (denom. of softmax) which doesn't scale well with the vocabsize.



Model

Vocab size = 10,000k



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

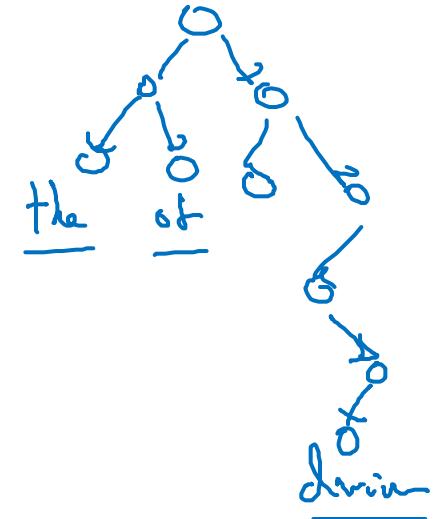
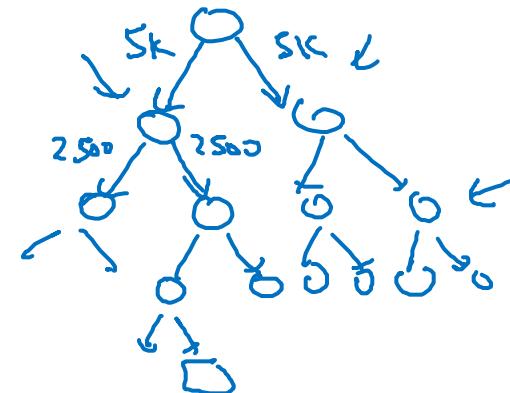
$$\rightarrow \ell(\hat{y}, y) = - \sum_{i=1}^{10,000} y_i \log \hat{y}_i$$

$$y = \begin{bmatrix} \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \leftarrow 4834$$

Problems with softmax classification

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

Hierarchical softmax.



How to sample the context c ?

avoid too frequent words \rightarrow the, of, a, and, to, ...
 \rightarrow orange, apple, durian

P_{durian}

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

Negative sampling is an alternative way to learn the skip-gram model that is more efficient than softmax classification. In softmax classification one wants to learn a mapping context \rightarrow target and to do that one needs to update the weights for all the words in the vocabulary. In the negative sampling this is avoided by rephrasing the problem as a binary classification one.



deeplearning.ai

1 - for each word in the corpus
create 1 POSITIVE EXAMPLE

$\langle (\text{word}, \text{context}), 1 \rangle$

K NEGATIVE EXAMPLES

$\langle (\text{word}, \text{wrong context}), 0 \rangle$ obtained sampling random contexts from the entire corpus.

- K ≈ 5 for large corpore and ≈ 20 for small corpore.

- The sampling of random context words is made so to pick very frequent words with a lower probability.

2 - Instead of learning the usual skip-gram $P(\text{word} | \text{context})$ with softmax, we learn an approximated version of it $P(1 | \text{word}, \text{context})$ with logistic regression. This allows to learn the word2vec more efficiently because for each sample you only update the weight between word and context (as opposed to updating the weights between a context and all possible words). It's unclear to me if you can then use the

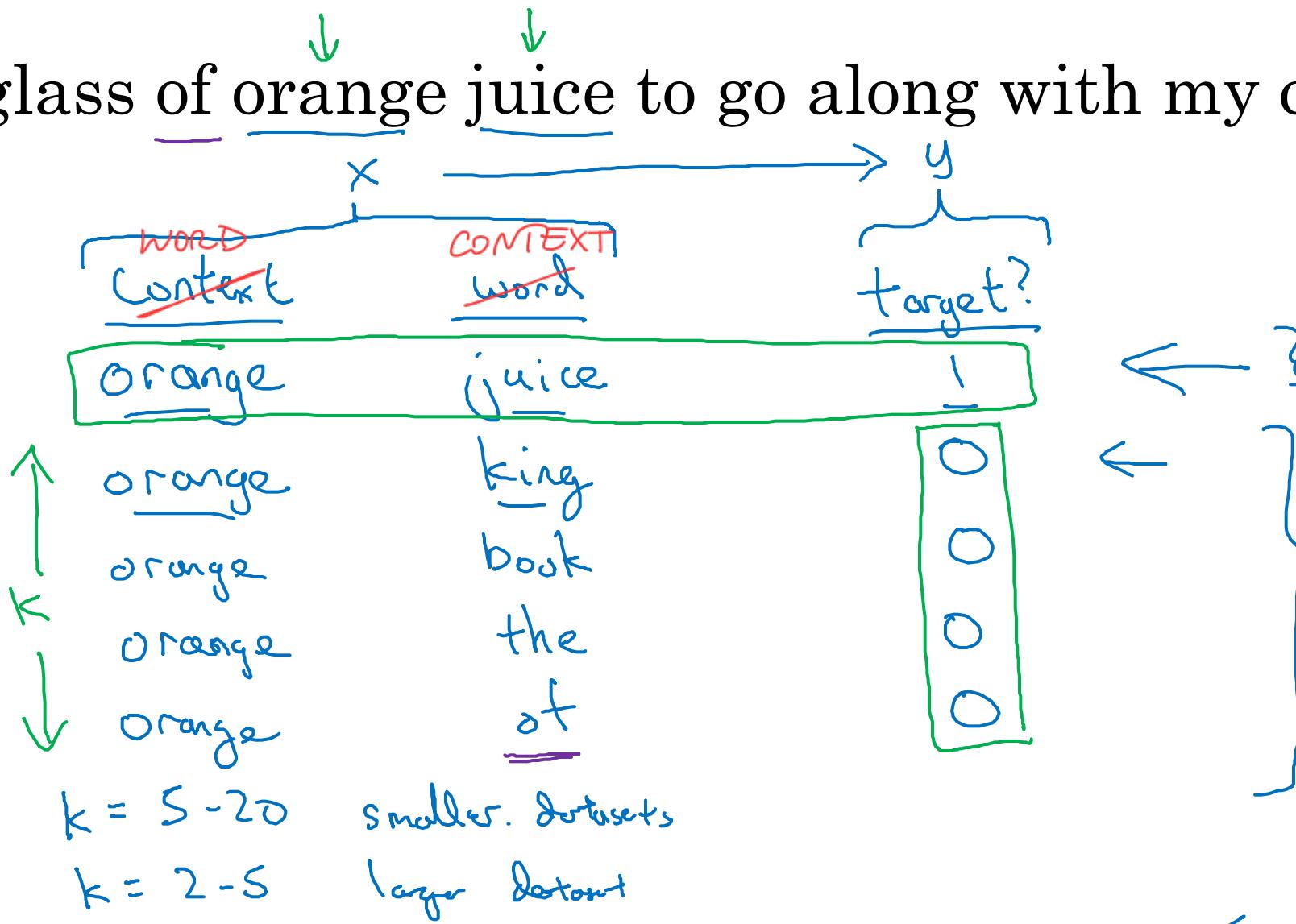
Negative sampling

logistic model from negative sampling as a skip-gram. Anyway people use it to learn the word2vec encoding mainly.

NLP and Word Embeddings

Defining a new learning problem

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



Model

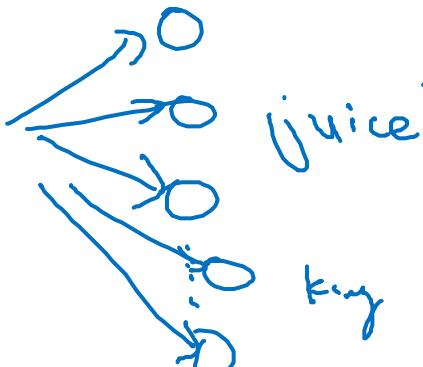
Softmax:

LOGISTIC REGRESSION PROBLEM

$$P(y=1 | c, t) = \sigma(\theta_t^T e_c)$$

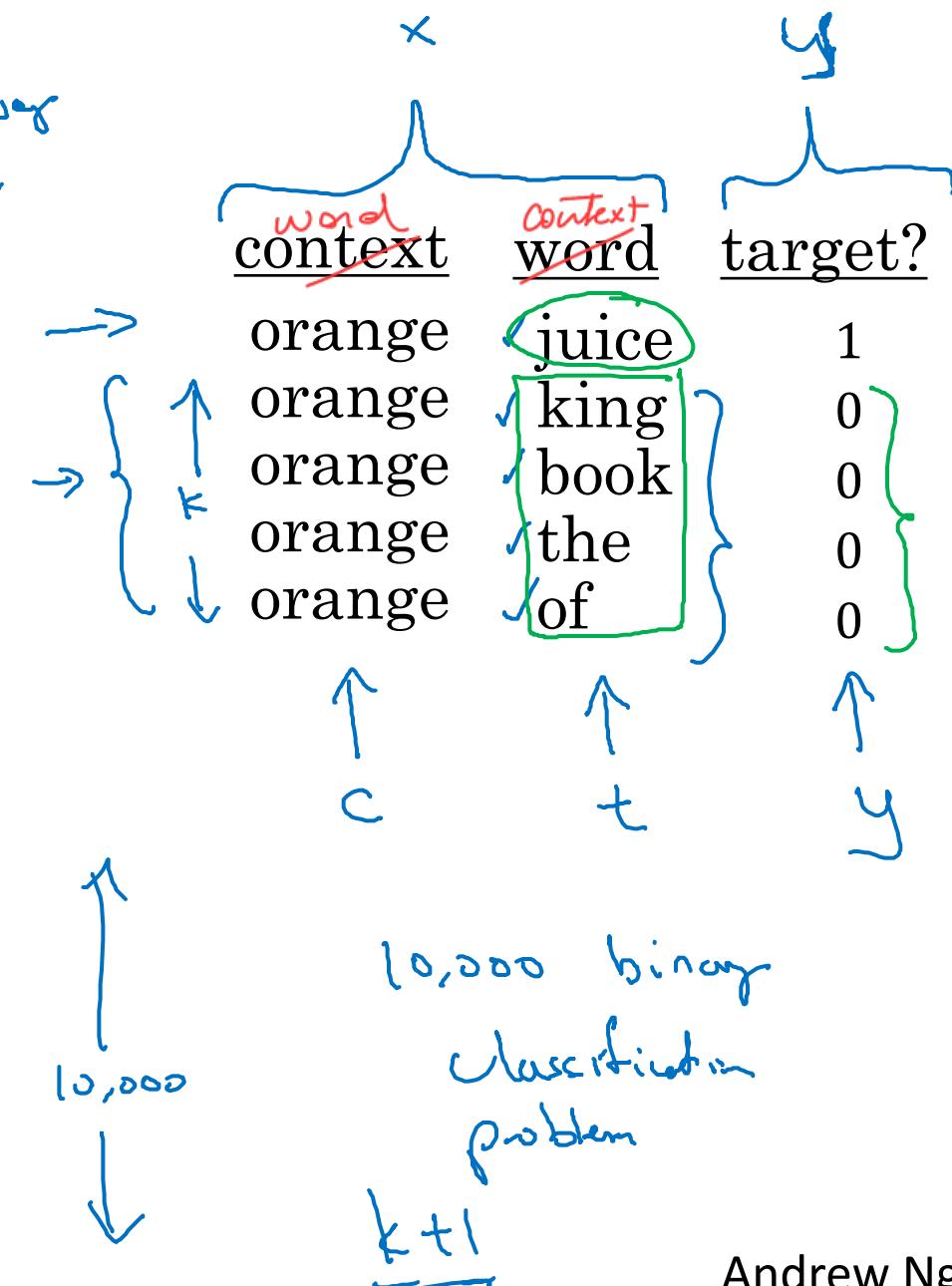
Orange
6257

$$\theta_{6257} \rightarrow E \rightarrow e_{6257}$$



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

10,000-way softmax



Selecting negative examples

<u>context</u>	<u>word</u>	<u>target?</u>
orange	juice	1
orange	king	0
orange	book	0
orange	the	0
orange	of	0

the , of, and, ...

t

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

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

Empirically it's best to sample negative examples
slightly inversely proportional to their frequency in the vocabulary.

Unclear to me... (esp. what are b_i and b_j ?)



deeplearning.ai

NLP and Word Embeddings

GloVe word vectors

- Computes the embedding of a word based on its global co-occurrence with all the other words in the corpus (within a limited window).
- loss s.t. similarity between 2 embeddings proportional to co-occurrence of those 2 words.
- Difference with word2vec: to my understanding is the procedure. While in word2vec we learn embedding through the task of learning context \rightarrow target, here we explicitly count cooccurrences...

GloVe (global vectors for word representation)

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

c, t

$x_{i,j} = \# \text{ times } i \text{ appears in context of } j$.

$x_{i,j}$ i
↑ ↑ ↑
c t

j
↑
c

$$x_{ij} = x_{ji} \leftarrow$$



Model

Minimize

$$\sum_{i=1}^{10,000} \sum_{j=1}^{100,000} f(x_{ij}) (\theta_i^T e_j + b_i + b_j' - \log \frac{x_{ij}}{c})$$

↑ ↓ ↑ ↓ ↗

$\theta_i^T e_j$

" $\theta_t^T e_c$ "

weight_{ij} term

$f(x_{ij}) = 0$ or $x_{ij} = 0$. "0 log 0" = 0

this, is, of, a, ...

derian

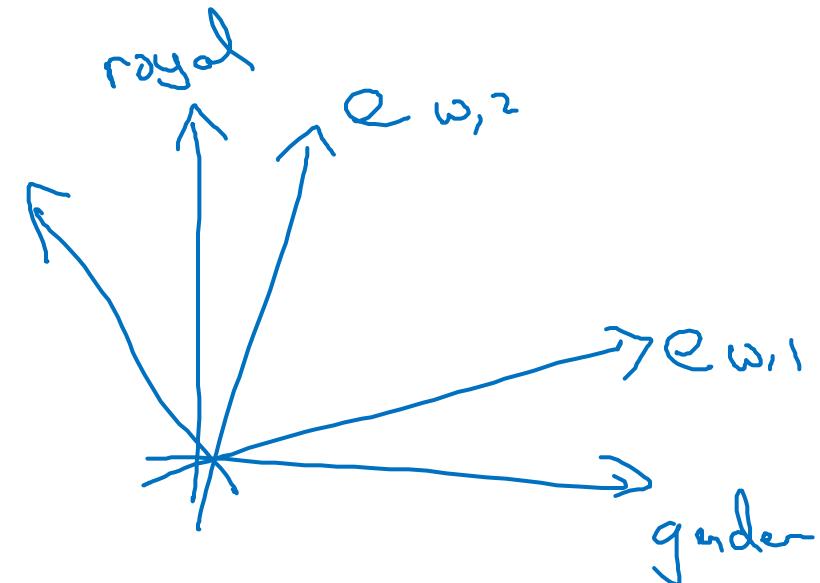
θ_i, e_j are symmetric

$\theta_w^{(\text{final})} = \frac{\theta_w + \theta_w}{2}$

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



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

$$\langle A\theta_i \rangle^T (A^T e_j) = \cancel{\theta_i^T A^T A} \cancel{e_j} = \theta_i^T e_j$$

In general, use RNN for sentiment classification (input = embedded words)



deeplearning.ai

NLP and Word Embeddings

Sentiment classification

Sentiment classification problem



The dessert is excellent.



Service was quite slow.



Good for a quick meal, but nothing special.



Completely lacking in good taste, good service, and good ambience.



10,000 → 100,000 words

Simple sentiment classification model

The dessert is excellent

8928 2468 4694 3180



The o_{8928} $\rightarrow E \rightarrow e_{8928}$

One possibility is "averaging the meaning" of the words in the sentence. This works well in practice but can lead to artifacts related to uncommon uses of a word.

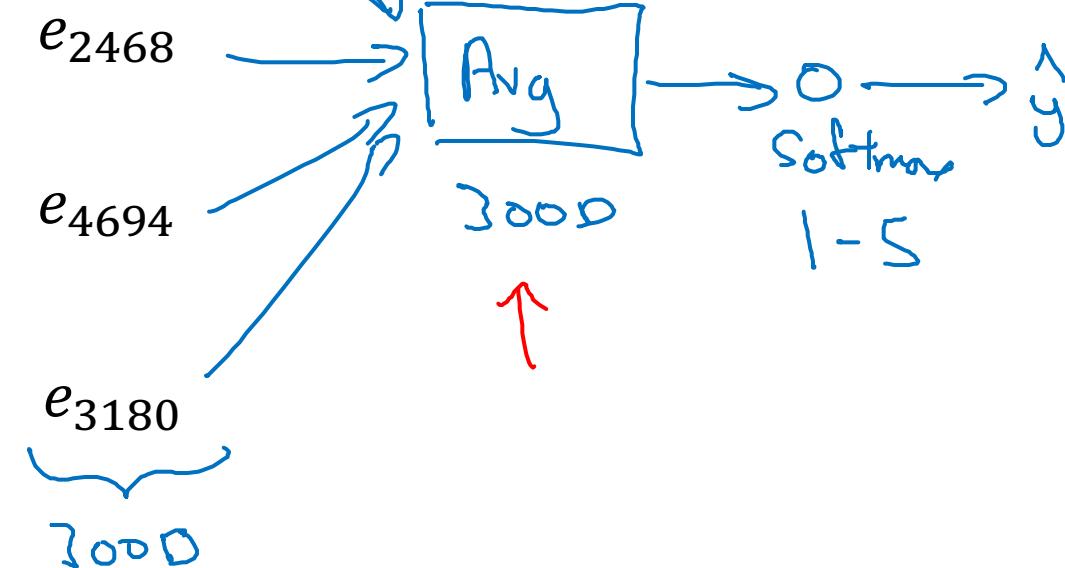
desert $o_{2468} \rightarrow E \rightarrow e_{2468}$

is $o_{4694} \rightarrow E \rightarrow e_{4694}$

excellent $o_{3180} \rightarrow E \rightarrow e_{3180}$

"Completely lacking in good taste, good service, and good ambience."

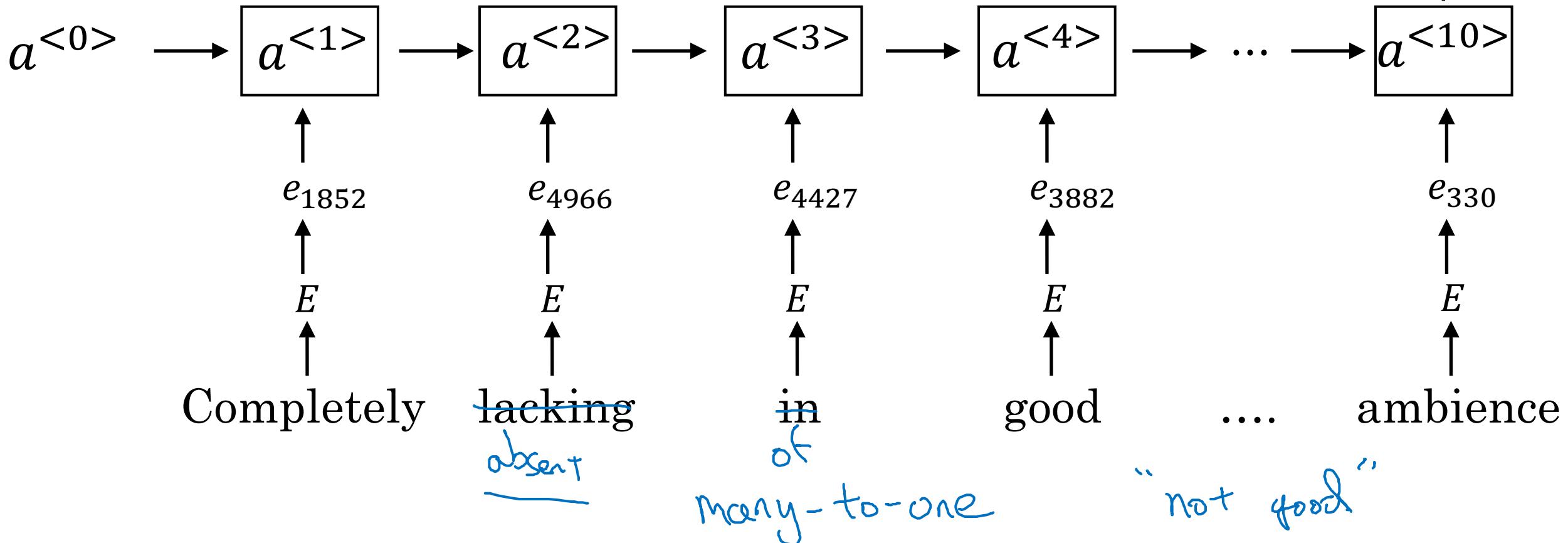
↑
100 B
words



RNN for sentiment classification

\hat{y}

softmax



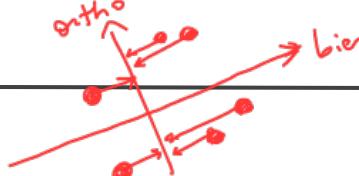
To remove bias from word encoding one pipeline could be :

- 1 - identify the bias direction in the embedding space by averaging the direction of pairs that are known to have that bias. e.g. average directions of boy-girl, man-woman, ...



to find the direction of "gender" bias.

- 2 - compress the embedding along the bias direction (often you can just eliminate it because the bias could be definitional for some words, i.e. employed in the definition of some words).



just eliminate it because the bias could be definitional for some words, i.e. employed in the definition of some words).

- 3 - for words for which the bias is definitional (e.g. man, woman, boy...)

make them

equally distant from bias-neutral words by operating a transformation of the space.

NLP and Word Embeddings

Debiasing word embeddings

The problem of bias in word embeddings

Man:Woman as King:Queen

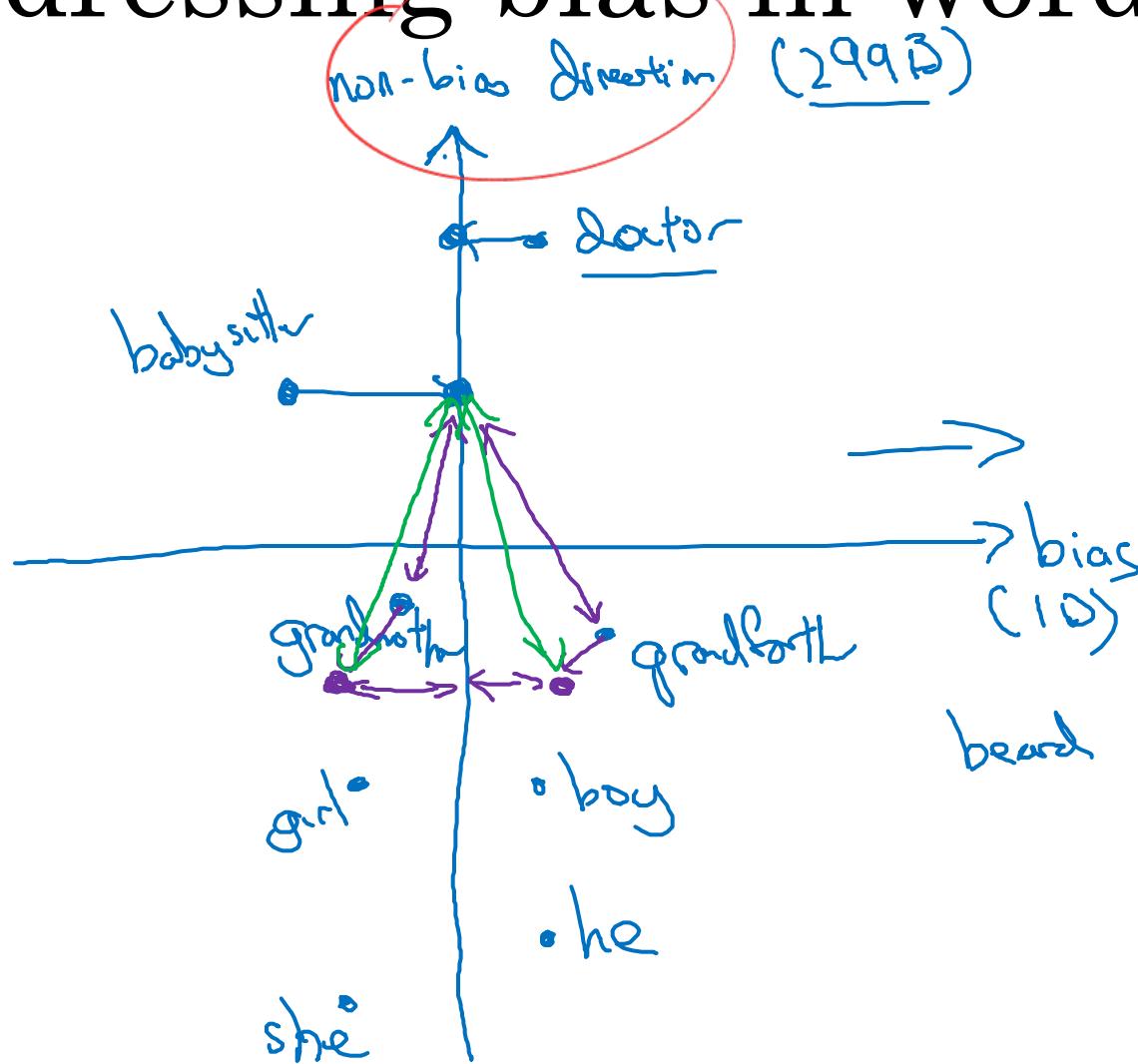
Man:Computer_Programmer as Woman:Homemaker 

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.



Addressing bias in word embeddings



1. Identify bias direction.

$$\begin{cases} \mathbf{e}_{\text{he}} - \mathbf{e}_{\text{she}} \\ \mathbf{e}_{\text{male}} - \mathbf{e}_{\text{female}} \\ \vdots \\ \text{average direction} \end{cases}$$

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

bias is part of the definition

3. Equalize pairs. *using linear algebra*

$$\rightarrow \text{grandmother} - \text{grandfather} = \text{girl} - \text{boy}$$