



deeplearning.ai

NLP and Word Embeddings

Word representation

Word representation

words with similar meanings have similar representations

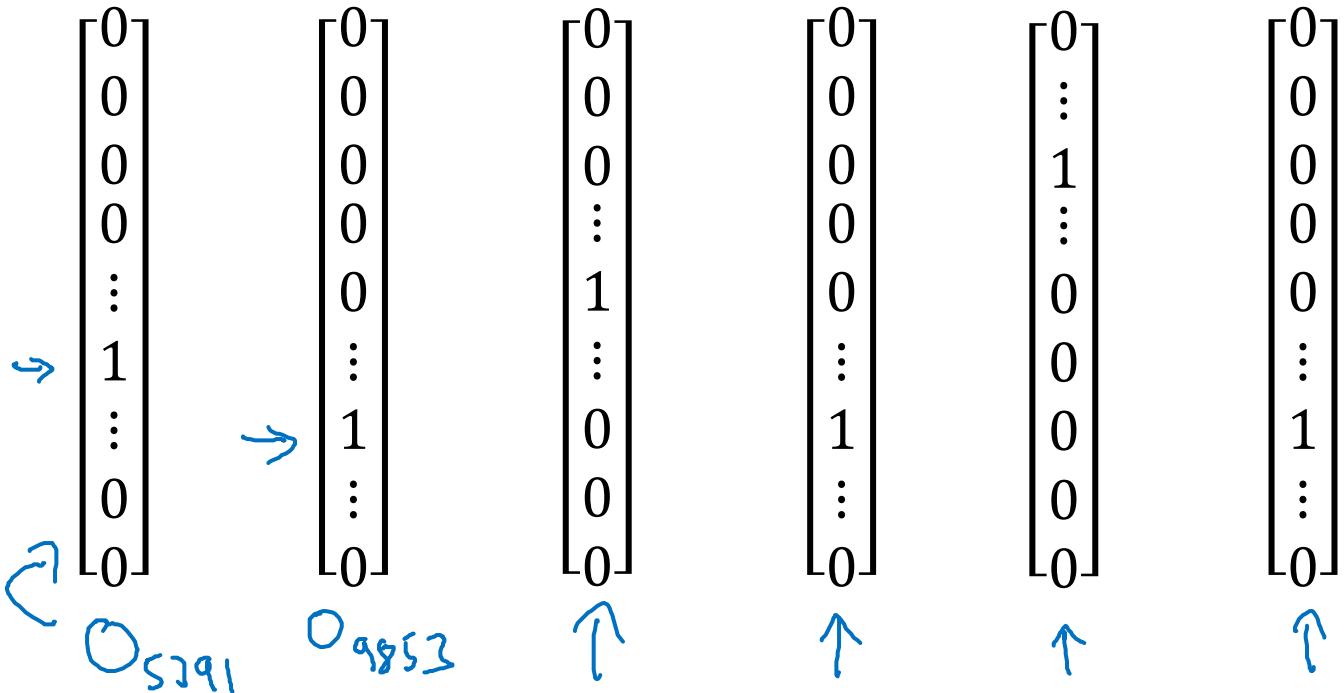
$$V = [a, \text{aaron}, \dots, \text{zulu}, \text{<UNK>}]$$

$$|V| = 10,000$$

1-hot representation

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

inner product shall
be zero, when words
are related



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

we aim to replace one-hot, which
is computationally expensive in
higher dimensions, with
something more agile & dense

Featurized representation: word embedding

| Feature list | 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 | | | | | | |

↑ Gender ←
↑ Royal ←
↑ Age ←
↑ Food ←
↑ Size ←
↑ Cost ←
↑ Color ←
↓ Verb ↓

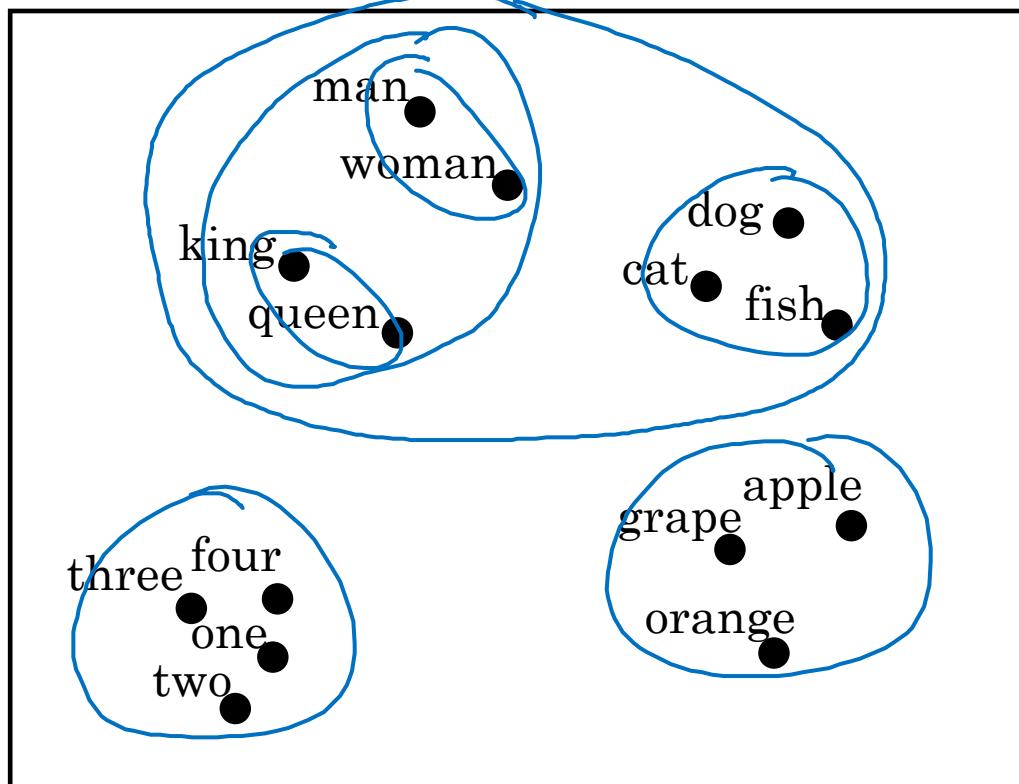
\mathbf{e}_{5391} \mathbf{e}_{9853}

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

Andrew Ng

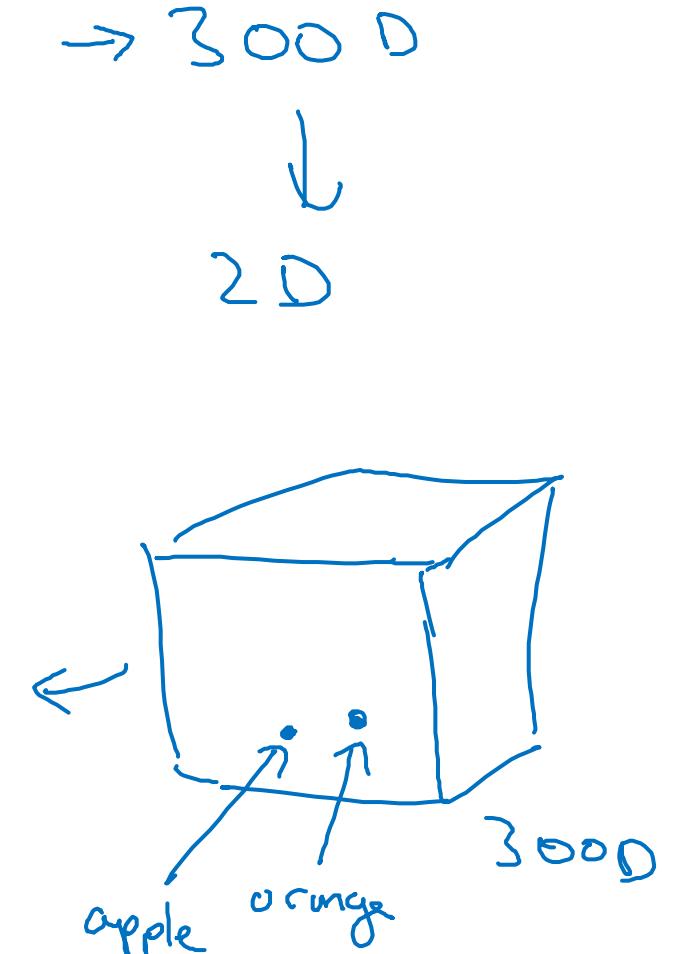
Visualizing word embeddings

Dimension of word vectors is usually smaller than the size of the vocabulary.
most common sizes for word vec ranges bet 50 & 400



t-SNE

non-linear dimensionality reduction technique



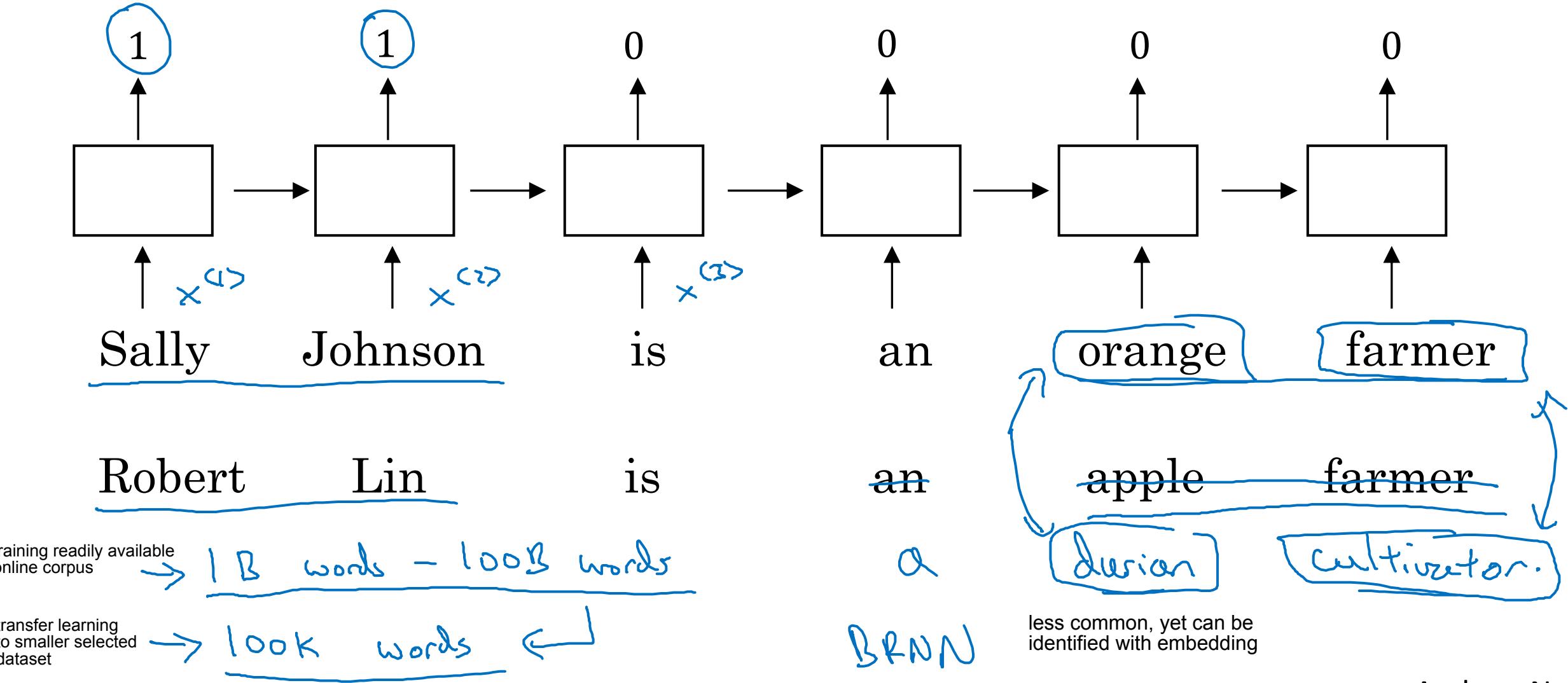


deeplearning.ai

NLP and Word Embeddings

Using word embeddings

Named entity recognition example



Transfer learning and word embeddings

makes huge difference when
trainingset is relatively small

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

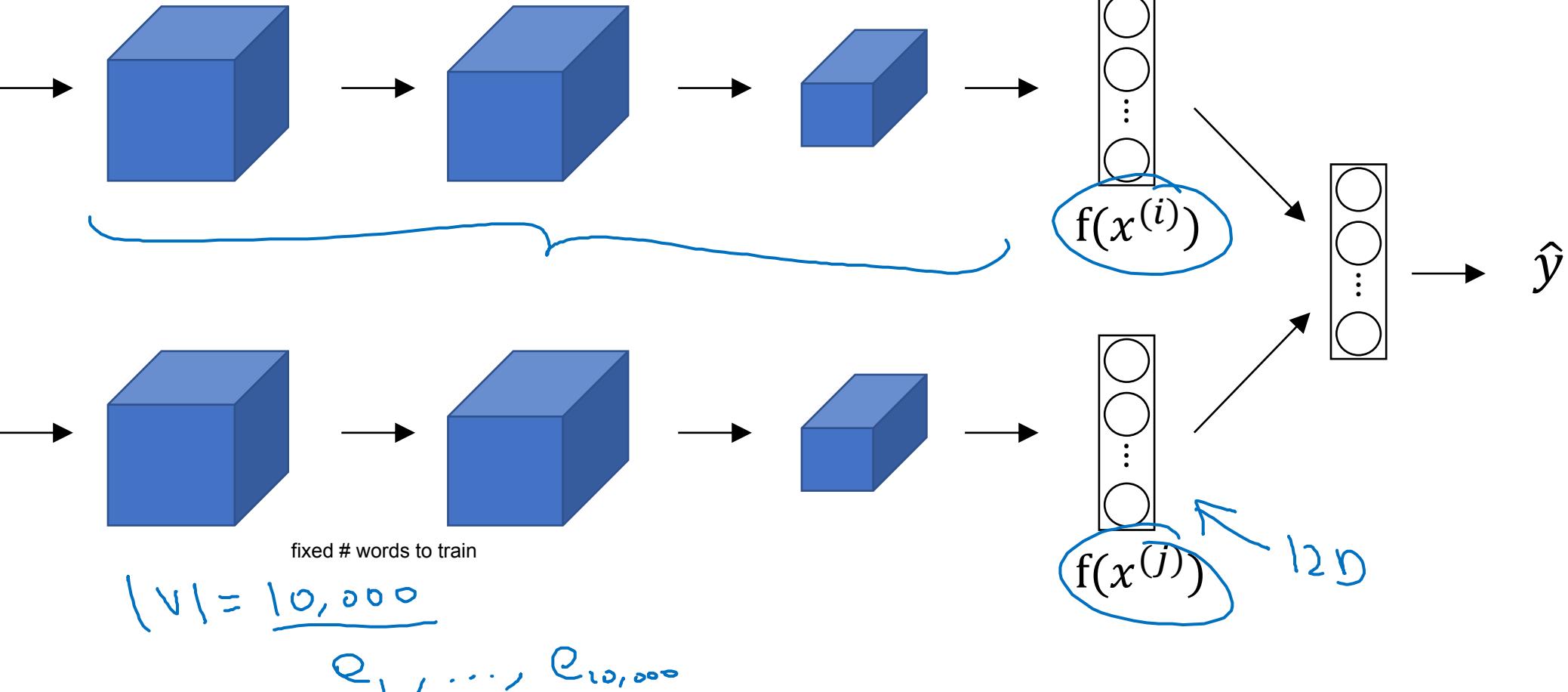
synonymous to each other

Relation to face encoding (embedding)

could be a totally new image



$x^{(i)}$





deeplearning.ai

NLP and Word Embeddings

Properties of word embeddings

Analogies

A good word embedding:
 $e_{\text{boy}} - e_{\text{girl}} = e_{\text{brother}} - e_{\text{sister}}$
 $e_{\text{boy}} - e_{\text{brother}} = e_{\text{girl}} - e_{\text{sister}}$

| assuming 4-dim 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 | 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 |

assuming..

$\underline{\text{Man} \rightarrow \text{Woman}}$ vs $\underline{\text{King} \rightarrow ? \text{ Queen}}$

required to compute $\underline{e_{\text{man}} - e_{\text{woman}}}$ $\approx \underline{e_{\text{king}} - e_{? \text{ Queen}}}$ vec difference

e_{5391}
 e_{man}

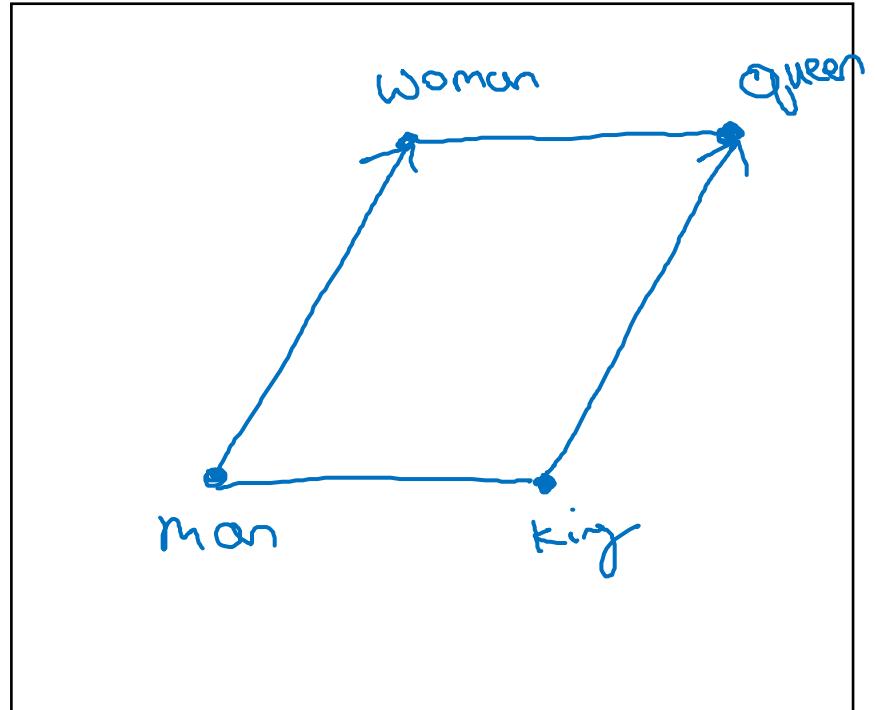
e_{woman}

e_{king}

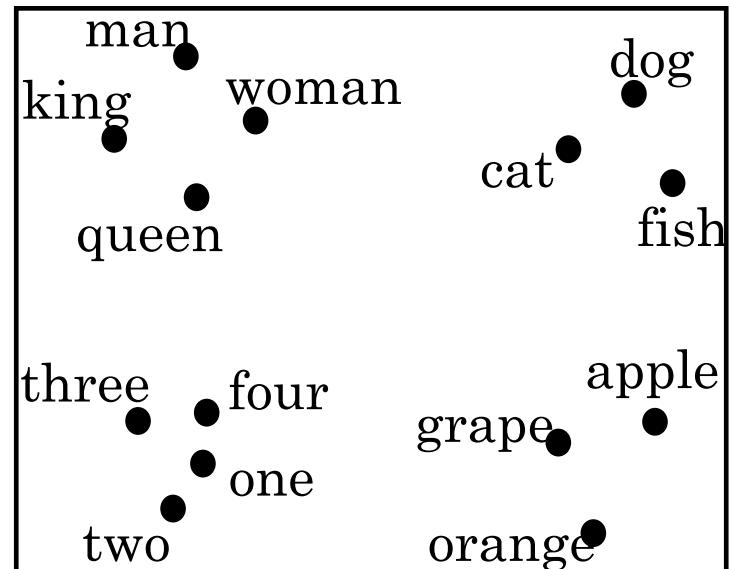
$e_{? \text{ Queen}}$

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

Analogies using word vectors



$300D \rightarrow 2D$



t-SNE

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

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

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

similarity function

30 - 75%

research papers acc

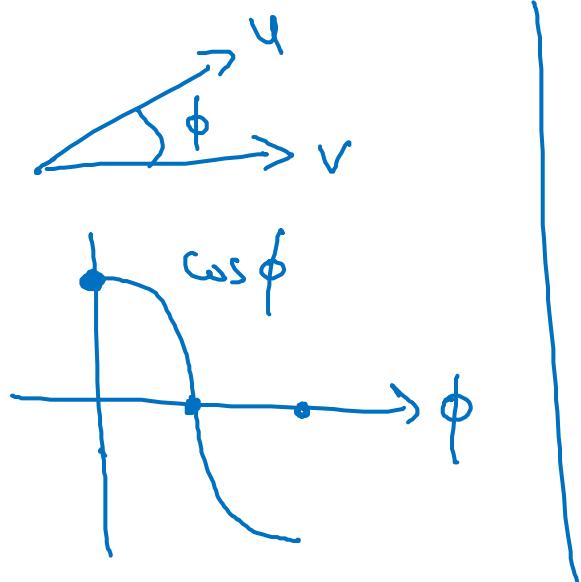
Andrew Ng

Cosine similarity

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

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

it's a judgement of orientation, not magnitude.
i.e 2 vec with same orientation have a cosine similarity of 1



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

another function: squared difference, measures the dissimilarities though

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

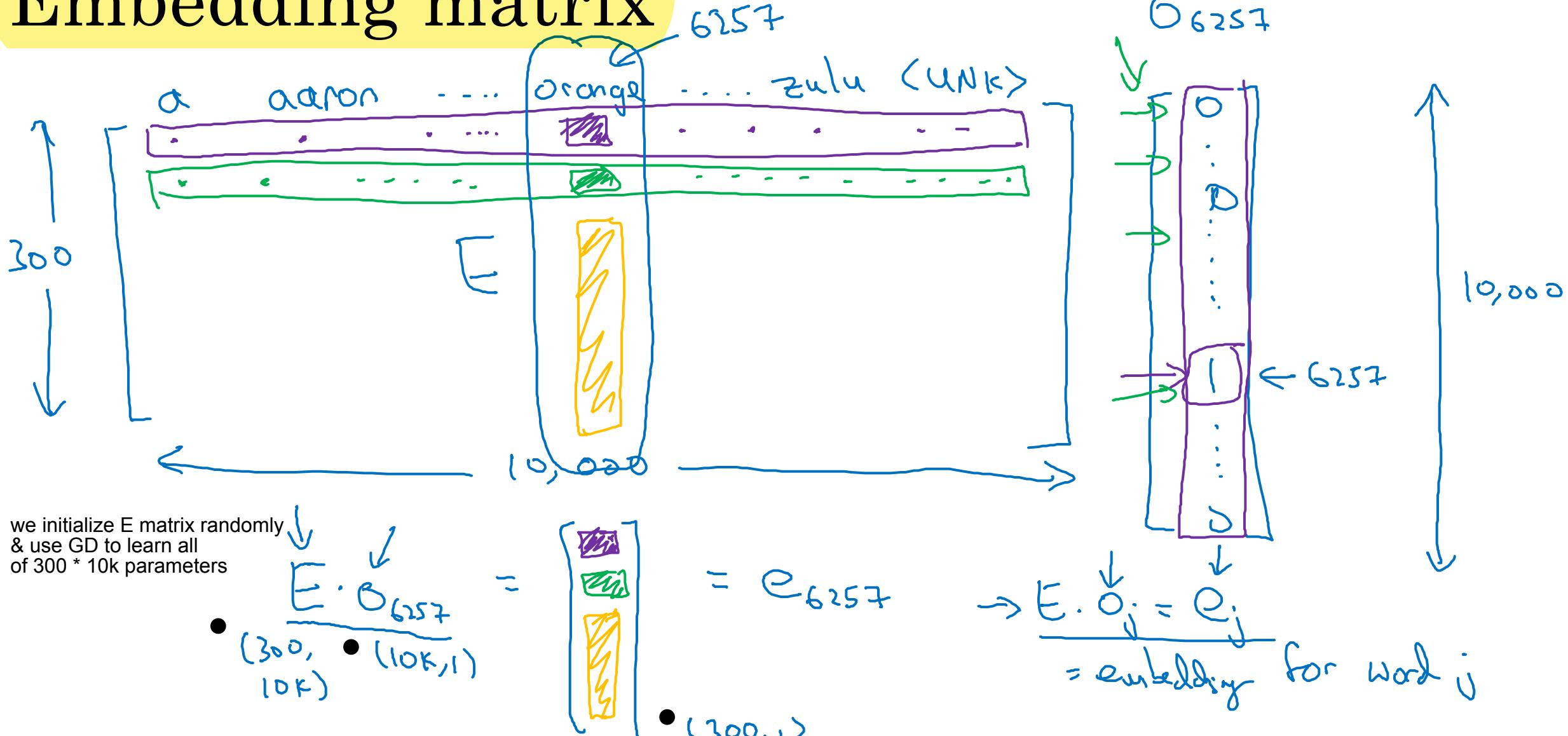


deeplearning.ai

NLP and Word Embeddings

Embedding matrix

Embedding matrix



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

keras \rightarrow Embedding

more efficient to be used than doing mat mul.
because of such high dim in one-hot

Andrew Ng



deeplearning.ai

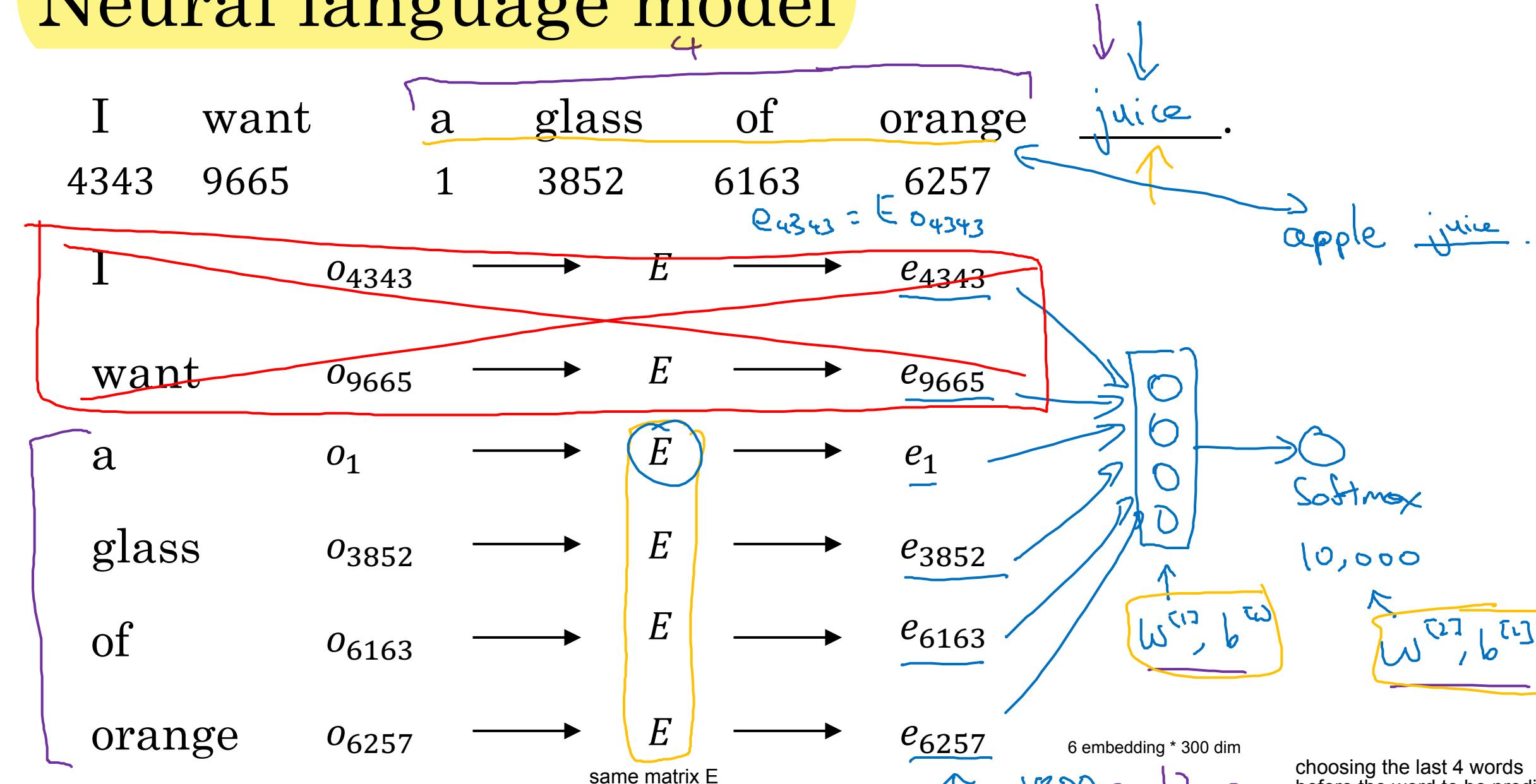
NLP and Word Embeddings

Learning word embeddings

Neural language model

earlier complex models

this is supervised learning



Other context/target pairs

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

Context

target

- Context: Last 4 words.

multiple options available

- 4 words on left & right
- Last 1 word
- Nearby 1 word

a glass of orange ? to go along with

orange ?

glass ?

skip gram

* building language model, use last few words
 * learning word embedding, use other contexts

- Important Note: It is okay if we do poorly on $P(t \text{ given } c)$, the more important by-product of this task is that we learn a useful set of word embeddings.



deeplearning.ai

NLP and Word Embeddings

Word2Vec

Skip-grams

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



Context

Context

orange

orange

orange



randomly choose
one word to be

Target

juice

glass

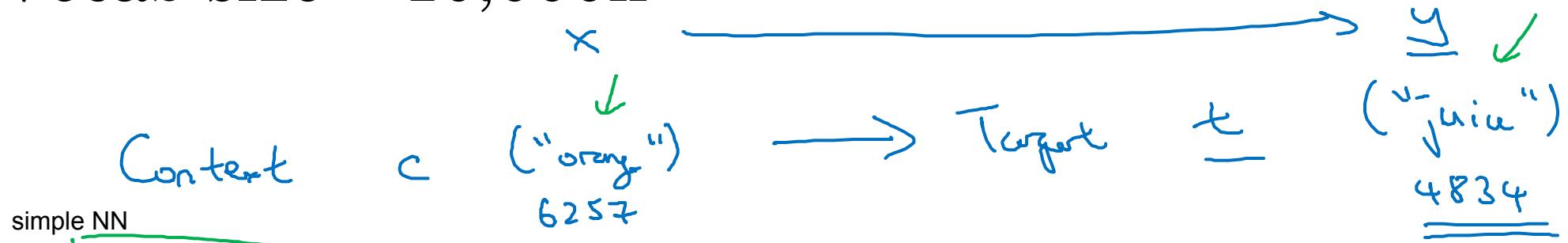
my



Model

Skip Gram Model

Vocab size = 10,000k



whose circled in green: e_c & softmax, contain learnable parameters

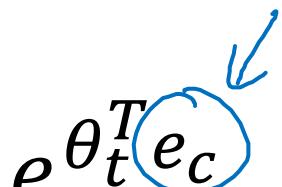
Softmax: $p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$

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

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

Problems with softmax classification

computational speed
of softmax is one
main problem

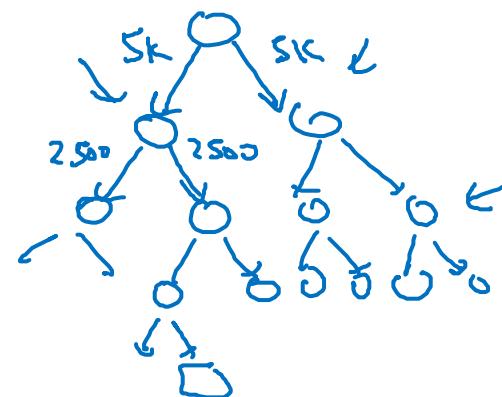
$$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$$


summing over all 10k,
is computationally expensive

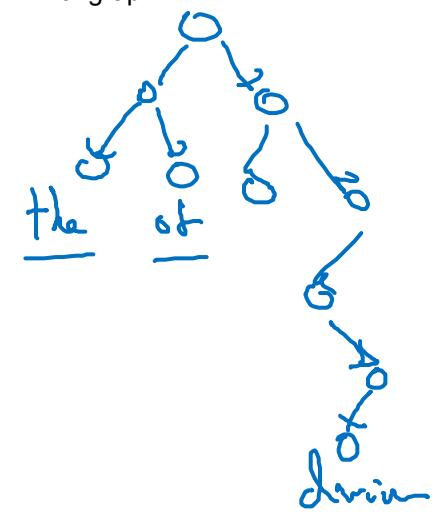
one solution:

Hierarchical softmax.

$\log |V|$



in practice, often-words
appear in earlier layers
of graph



How to sample the context c ?

frequently occurring
words

\rightarrow the, of, a, and, to, ...

\rightarrow orange, apple, durian

not entirely random,
some heuristics are used

c & t shall be chosen to
be nearby words

t

$c \rightarrow t$

$P(c)$

P_{durian}

- θ_t & e_c will have the same # dimensional vector as the word embedding
- they are both trained with optimization



deeplearning.ai

NLP and Word Embeddings

Negative sampling

Defining a new learning problem

this is an update of Word2Vec model;
to address the problem of having a huge
NN, requiring lots of training data

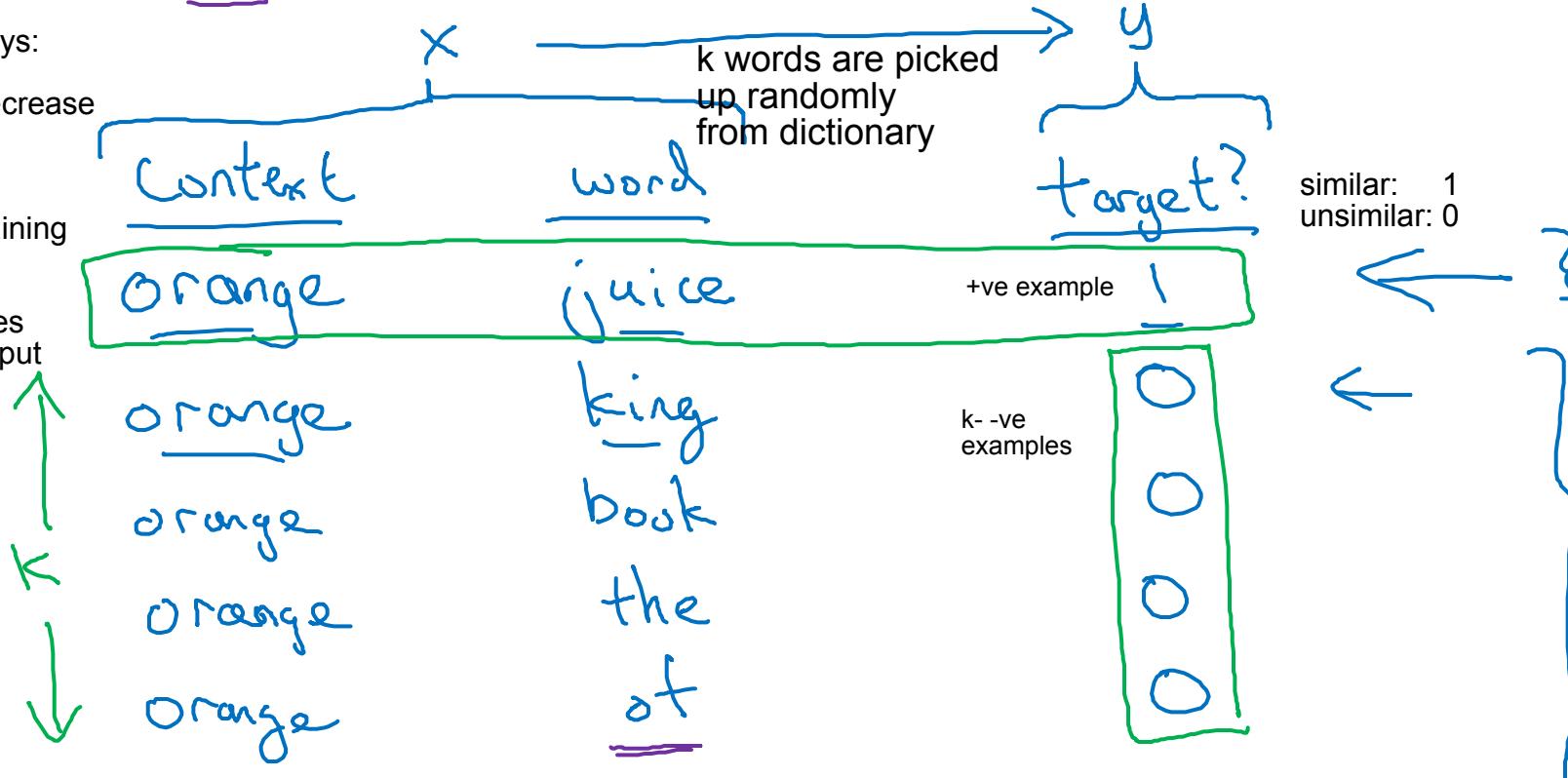
I want a glass of orange juice to go along with my cereal.

Authors tackled this problem in 2 ways:

1. subsampling frequent words to decrease
training examples.

2. modifying optimization through
-ve sampling, which causes each training
sample to update only small % of
the huge model weights.

10,000 words in corpus * 300 features
a total of 3M in each of hidden & output
layers



Check:
Word2Vec_McCormick notes

choose: $k = 5-20$ for smaller datasets

$k = 2-5$ larger dataset

Model

Negative Sampling

we deliberately generate -ve examples

Softmax:

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

↓ ↓

possible target word possible embedding word

instead of
10,000-way softmax

logistic regression model

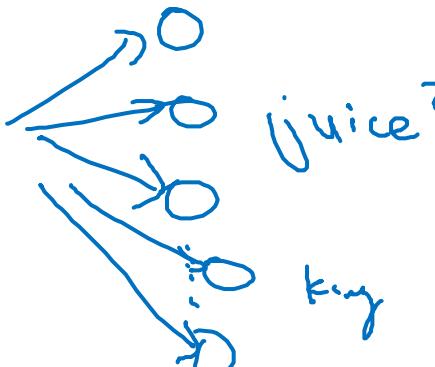
$$P(y=1 | c, t) = \sigma(\theta_t^T e_c)$$

input:

Orange
6257

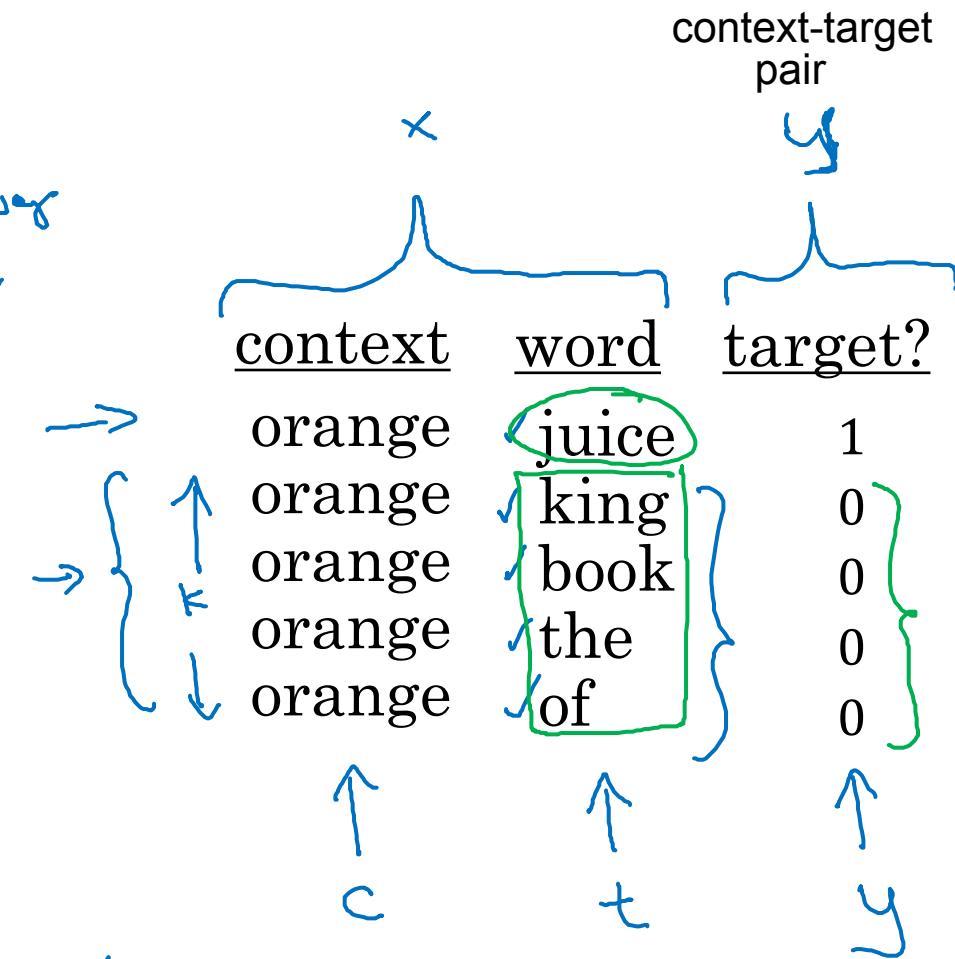
$o_{6257} \rightarrow E \rightarrow e_{6257}$
one-hot embedding matrix

embedding



● we update

10,000 binary classification problem
-ve examples +ve example
 $k+1$



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 |

high rep. of
such words

the , of, and, ...

one option: sampling
according to empirical
frequency of words
in corpus

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

another option, yet
not very representative

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

vocab size



deeplearning.ai

NLP and Word Embeddings

GloVe word vectors

GloVe (global vectors for word representation)

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

previously
we used:

c, t

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

\uparrow i \uparrow \uparrow
 c t j

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

It's a count capturing how many times i & j appear close to each other i.e +/- 10 words



Model

optimizing
by GD

minimize

- θ_i & e_j shall be initialized randomly at the beginning of training
- X_{ij} : # times word j appears in the context of word i

$$\sum_{i=1}^{10,000} \sum_{j=1}^{100,000} f(X_{ij}) (\theta_i^T e_j + b_i - b_j) - \log X_{ij}^2$$

Annotations:

- $f(X_{ij})$ is the weighting term, which is zero if $X_{ij} = 0$.
- The term $\theta_i^T e_c$ represents the "relatedness" between words i and c.
- $\theta_i^T e_j$ is the context vector for word j.
- b_i and b_j are bias terms.
- $\log X_{ij}^2$ is a measure of how often words i and j occur together.
- A question mark above the log term indicates it might be a typo or a placeholder.
- A note on the right says: "how related i & j, as measured by how often they occur with each other".

heuristics can be used to define weighting term; giving weight to both.
more frequent weights shall be given more weight though

weighting fn $f(\cdot)$ must satisfy $f(0)=0$

it helps prevent learning only from extremely common word pairs.
It is not necessary that it satisfies this function

stop words
this, is, of, a, ...

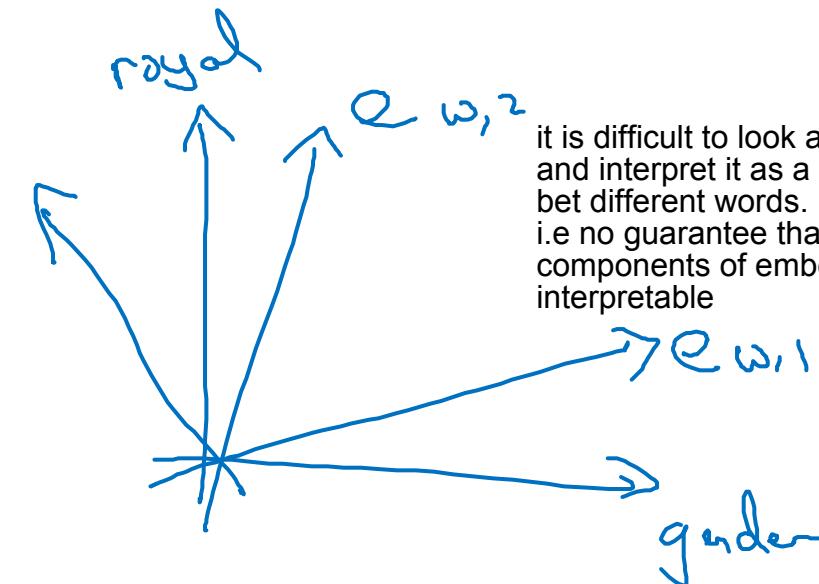
more frequent → duration ← less frequent

• θ_i , e_j are symmetric

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

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 |



it is difficult to look at a single row and interpret it as a direct relation bet different words.
i.e no guarantee that individual components of embeddings are interpretable

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

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



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.



- one challenge:
no huge datasets

10,000 \rightarrow 100,000 words

Simple sentiment classification model

The dessert is excellent

8928 2468 4694 3180



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

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

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

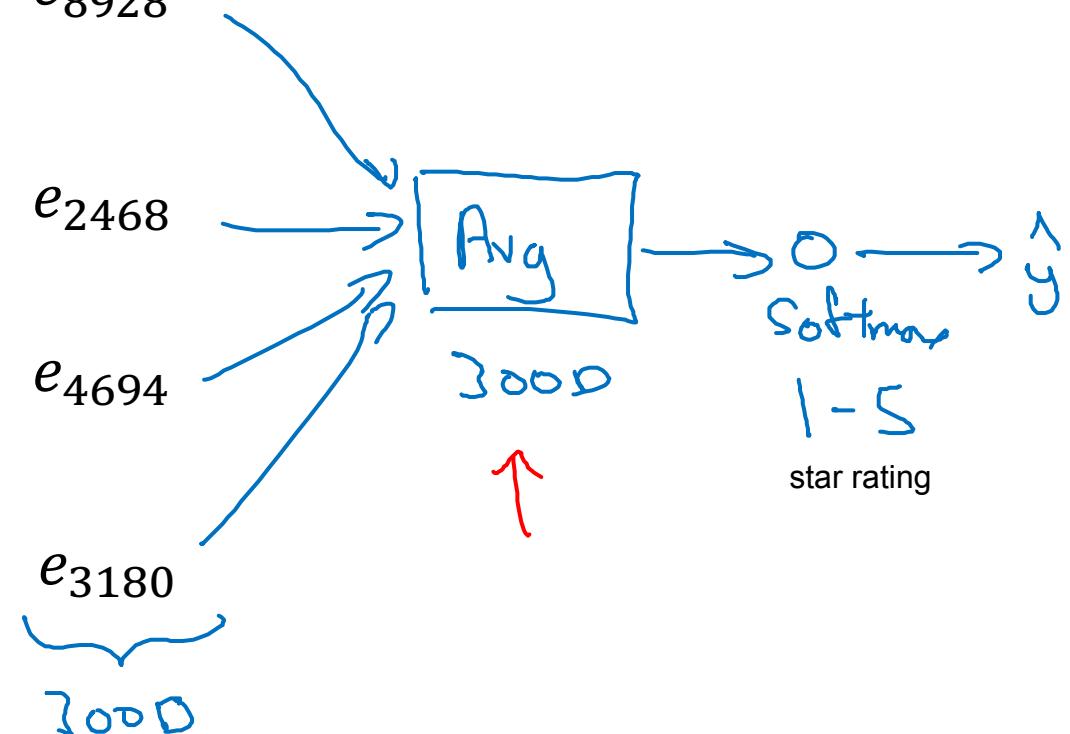
excellent $o_{3180} \rightarrow E \rightarrow e_{3180}$

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

↑
100 B words

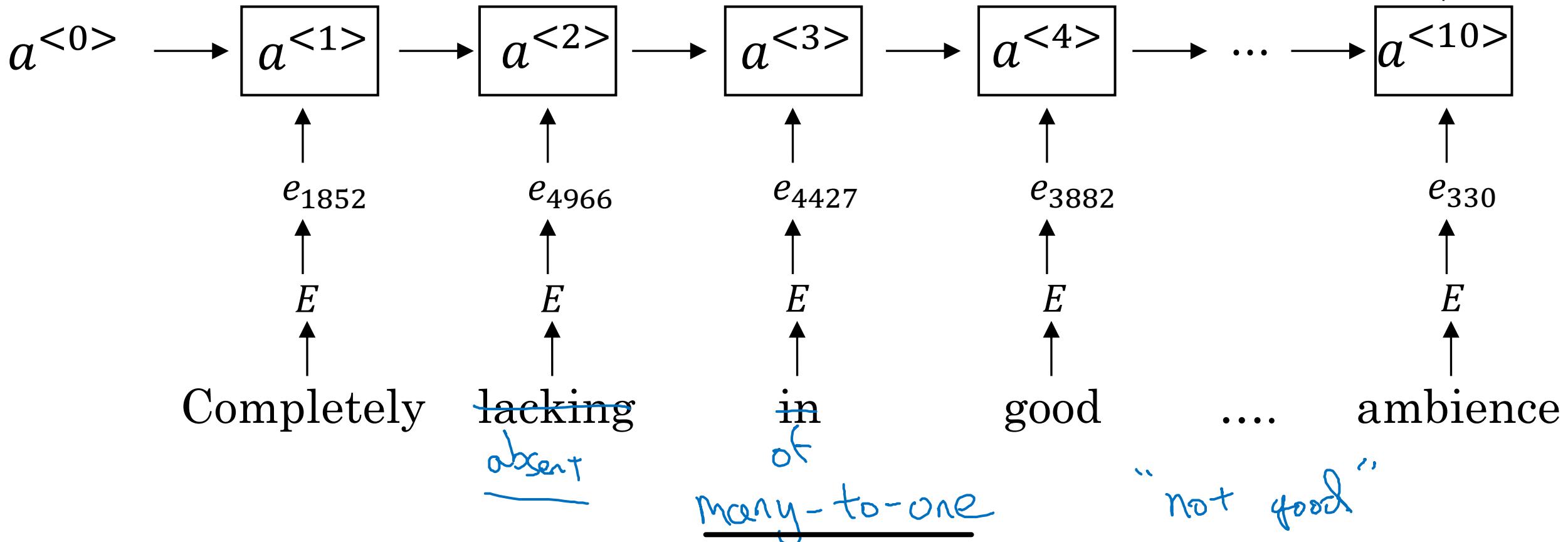
It's a process of summing all words in the example.

one con: is ignoring word order



RNN for sentiment classification

avoiding the prev con





deeplearning.ai

NLP and Word Embeddings

Debiasing word embeddings

The problem of bias in word embeddings

Man:Woman as King:Queen

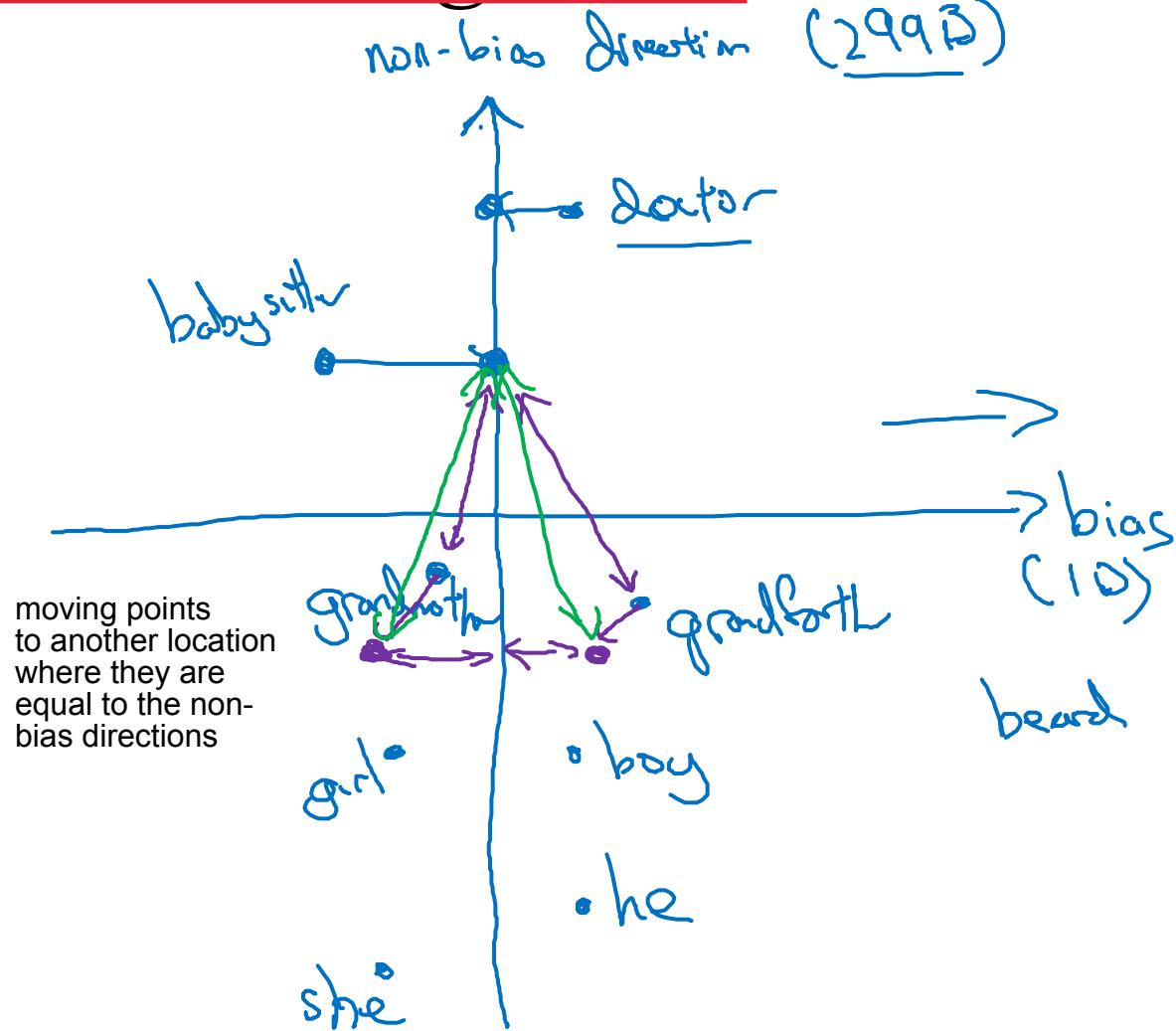
Man:Computer_Programmer as Woman:Homemaker 

Father:Doctor as Mother:Nurse 

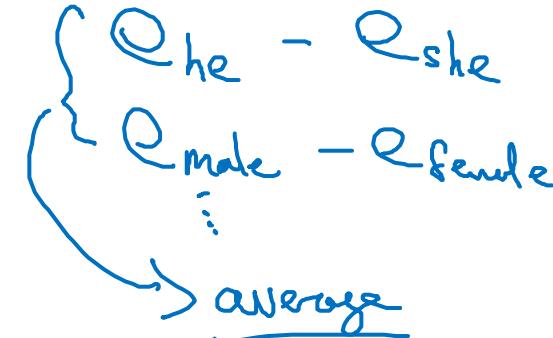
- ✓ Word embeddings can reflect gender, ethnicity, age, sexual orientation, and other biases of the text used to train the model.



Addressing bias in word embeddings



1. Identify bias direction. e.g gender



words like grandpa & grandma are by construction gender included, unlike doctor or babysitter

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

such equalized pairs (grandma, grandpa) shall have the same similarity (distance) with gender neutral words (doctor)

3. Equalize pairs.

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