

# Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>



deeplearning.ai

# NLP and Word Embeddings

---

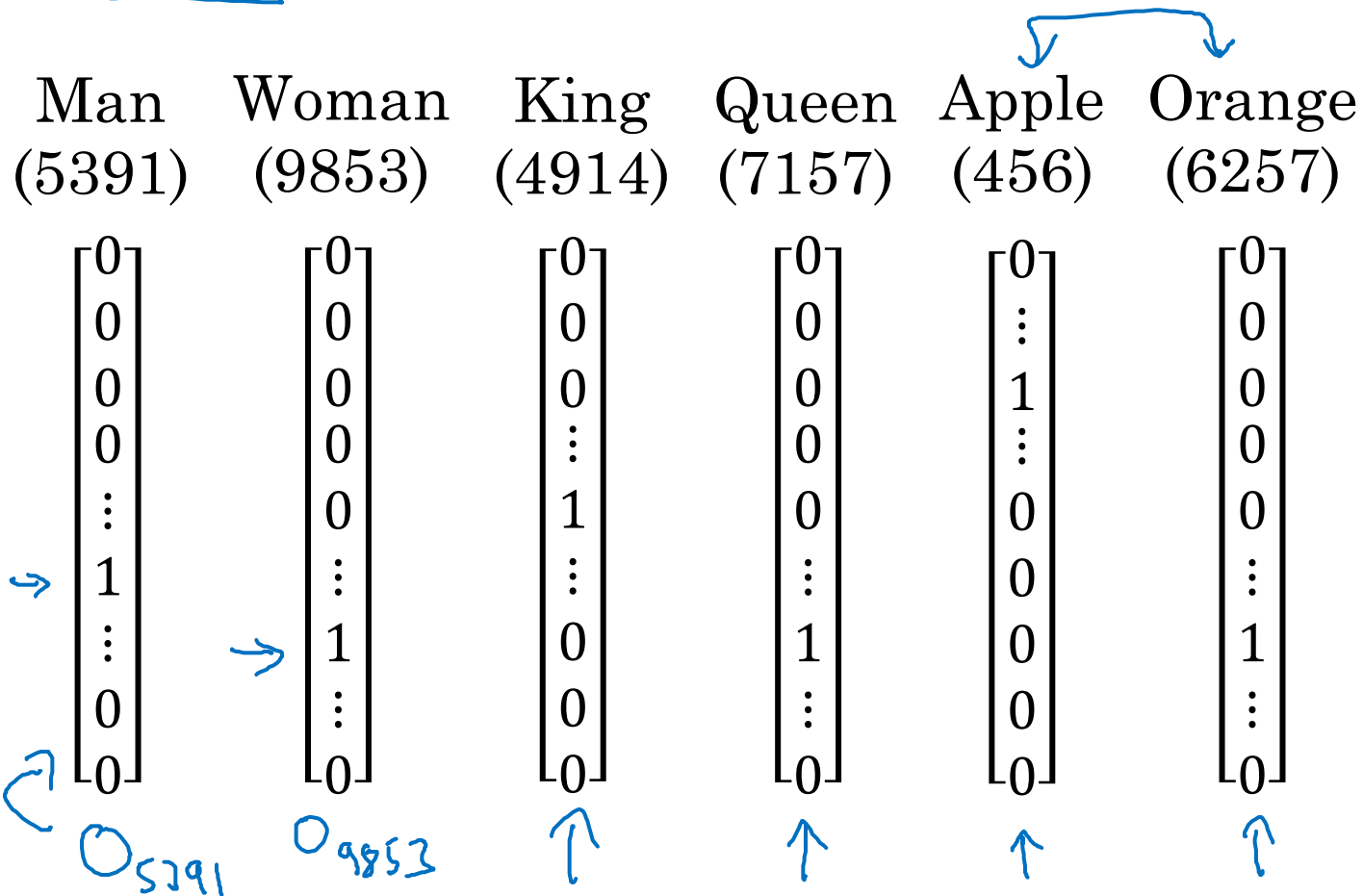
## Word representation

# Word representation

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

$|V| = 10,000$

1-hot representation not capture similarity between words



I want a glass of orange juice.

I want a glass of apple ?.

# Featurized representation: word embedding

word aspect	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						

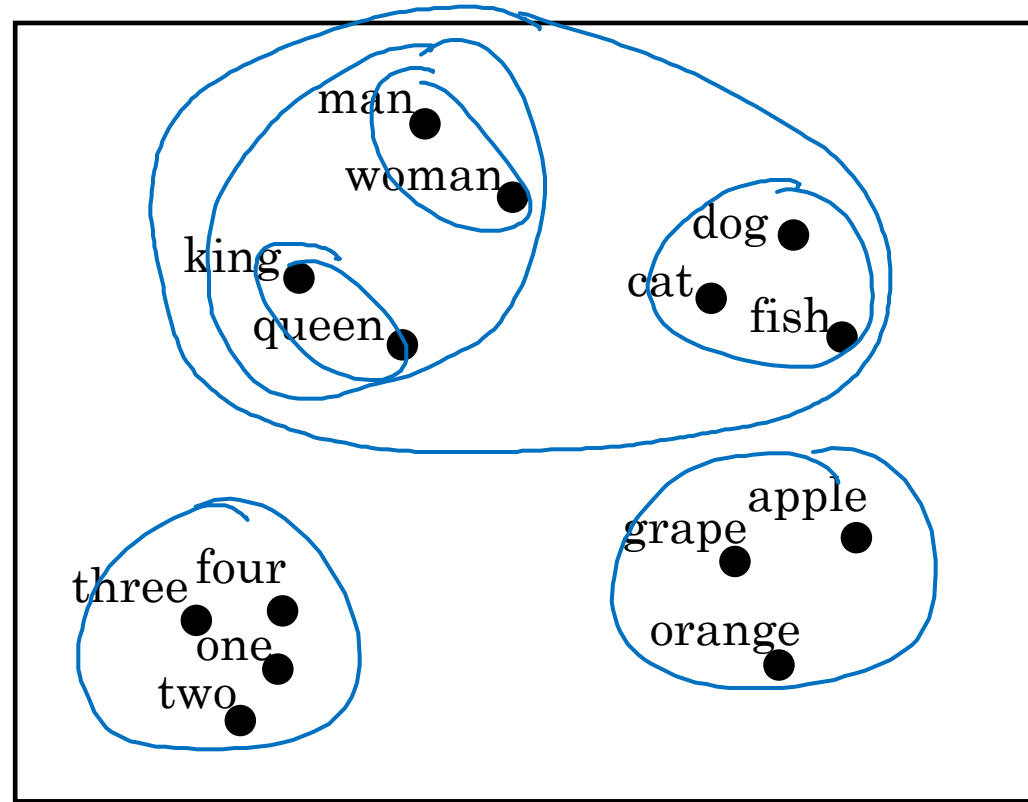
I want a glass of orange juice.

I want a glass of apple juice.

Andrew Ng

# Visualizing word embeddings

similar words tend to be near in vector space

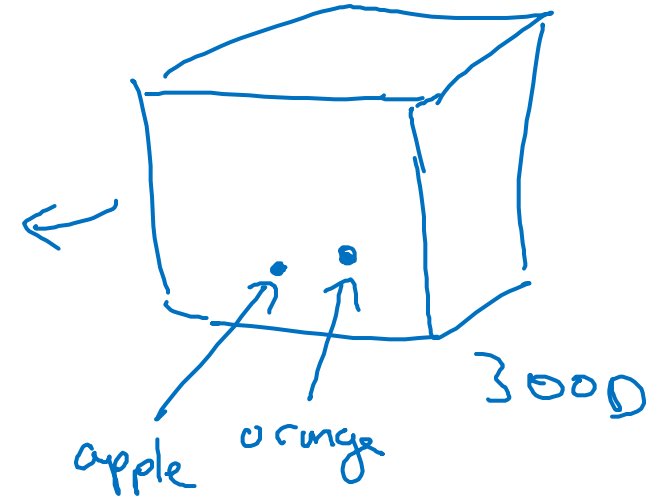


t-SNE

→ 300D



2D





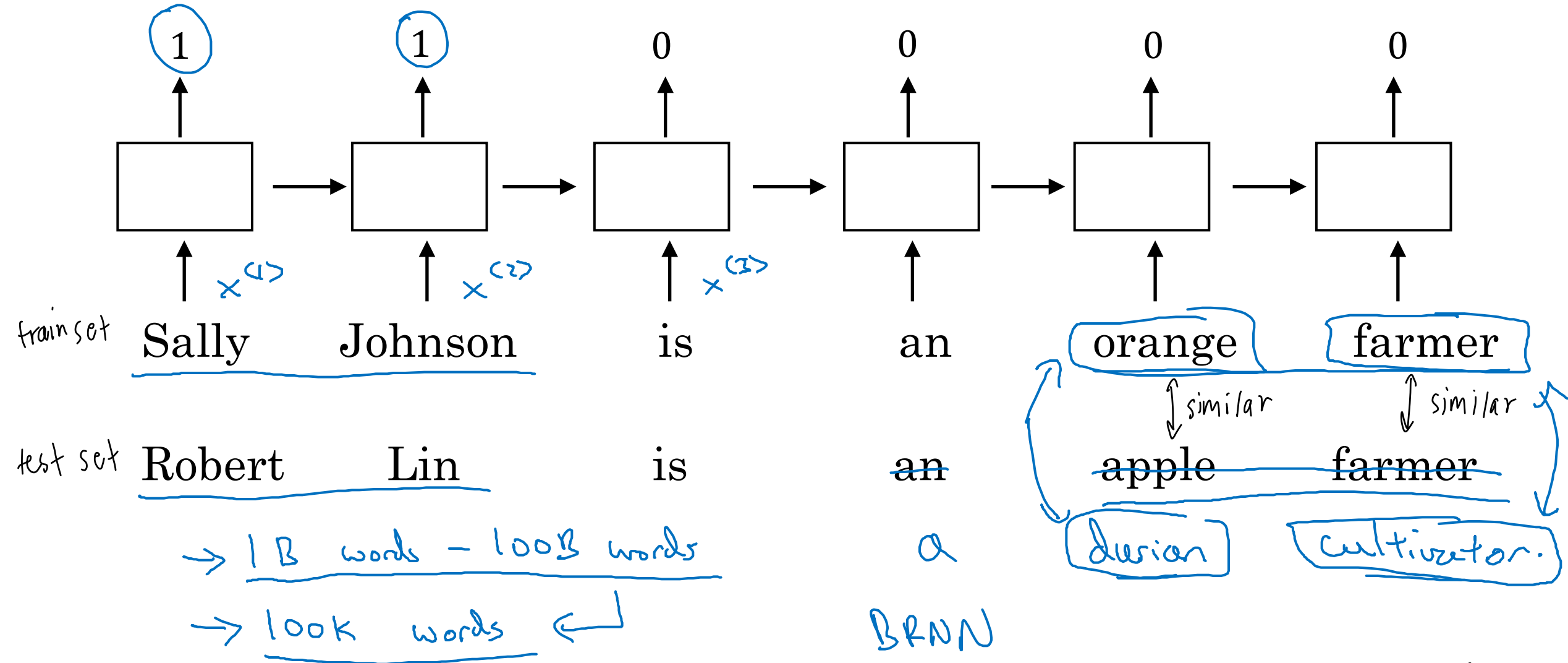
deeplearning.ai

# NLP and Word Embeddings

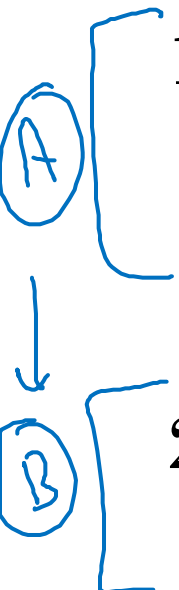
---

Using word  
embeddings

# Named entity recognition example



# Transfer learning and word embeddings

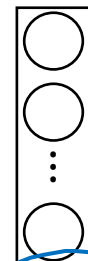
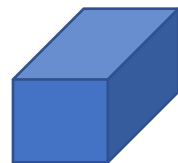
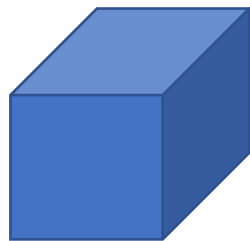
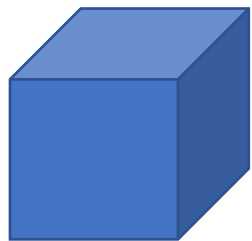
- 
1. Learn word embeddings from <sup>large train set</sup> large text corpus. (1-100B words)  
(Or download pre-trained embedding online.)
2. Transfer embedding to new task with smaller training set.  
(say, 100k words) → 10,000 → 300
3. Optional: Continue to finetune the word embeddings with new data.  
↳ need large amount of training data



# Relation to face encoding (embedding) 128D



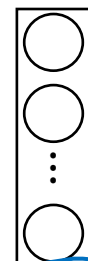
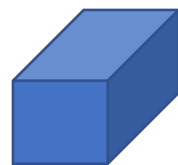
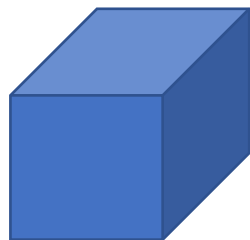
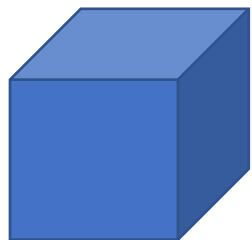
$x^{(i)}$



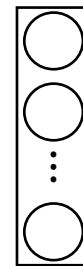
$f(x^{(i)})$



$x^{(j)}$



$f(x^{(j)})$



$\hat{y}$

$|V| = 10,000$

$e_1, \dots, e_{10,000}$



deeplearning.ai

# NLP and Word Embeddings

---

## 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

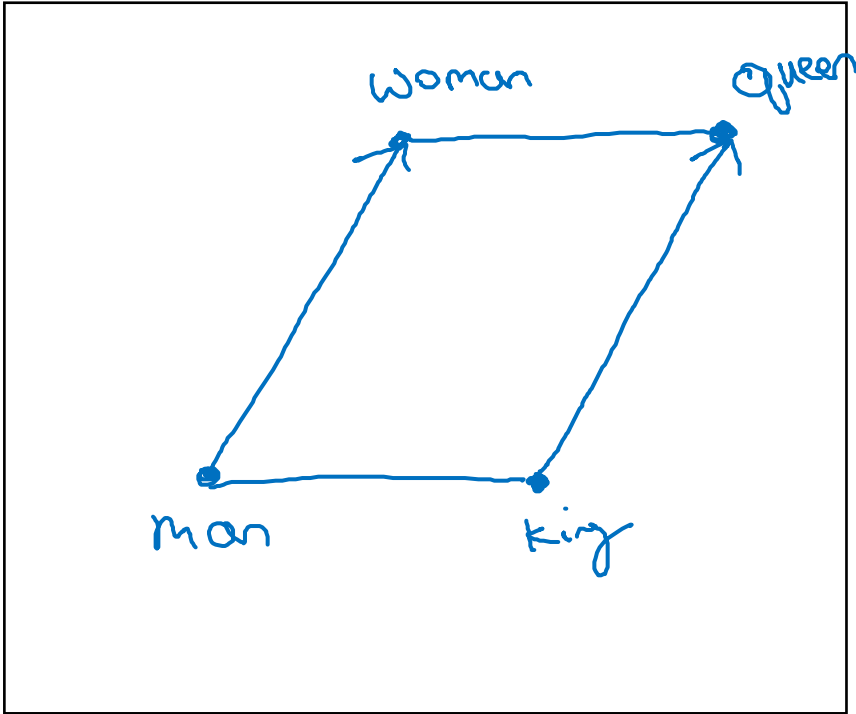
$\underbrace{e_{5391}}_{e_{\text{man}}} \rightarrow \underbrace{e_{9853}}_{e_{\text{woman}}} \quad \Leftrightarrow \quad \underbrace{e_{4914}}_{e_{\text{king}}} \rightarrow ? \quad \underbrace{e_{7157}}_{e_{\text{queen}}}$

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

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

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

# Analogies using word vectors



300 D

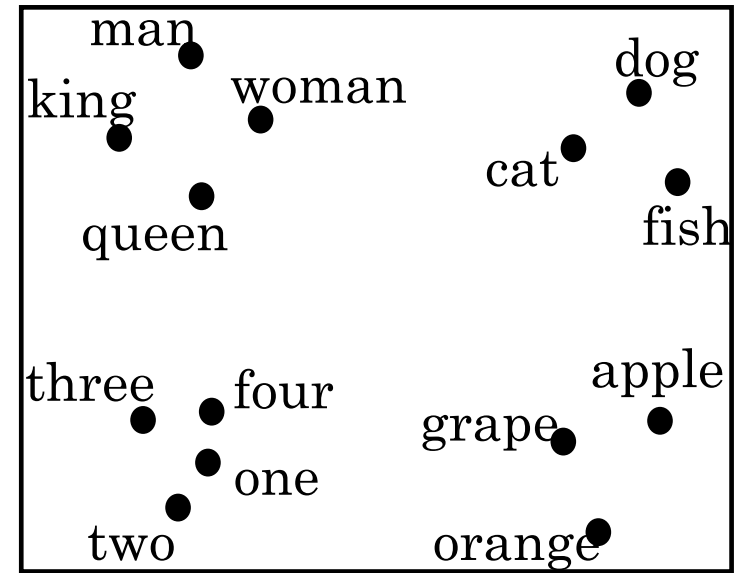
Find word  $w$  with  $\arg \max_w$

choose vector most similar to  $e_{king} - e_{man} + e_{woman}$

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

30 - 75%

3000  $\rightarrow$  20  
 $\uparrow$



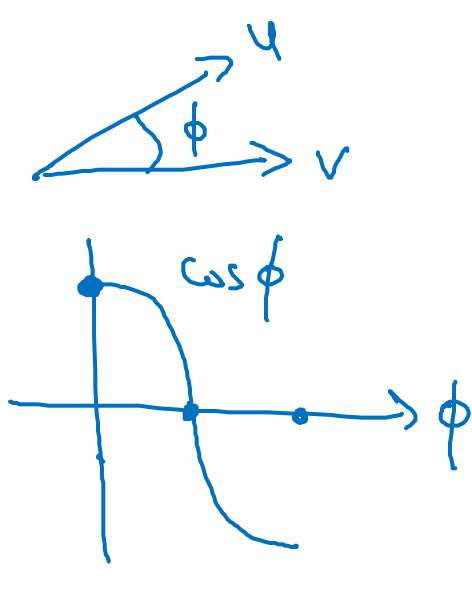
t-SNE

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

# Cosine similarity

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

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



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

Man:Woman as Boy:Girl

Ottawa:Canada as Nairobi:Kenya

Big:Bigger as Tall:Taller

Yen:Japan as Ruble:Russia



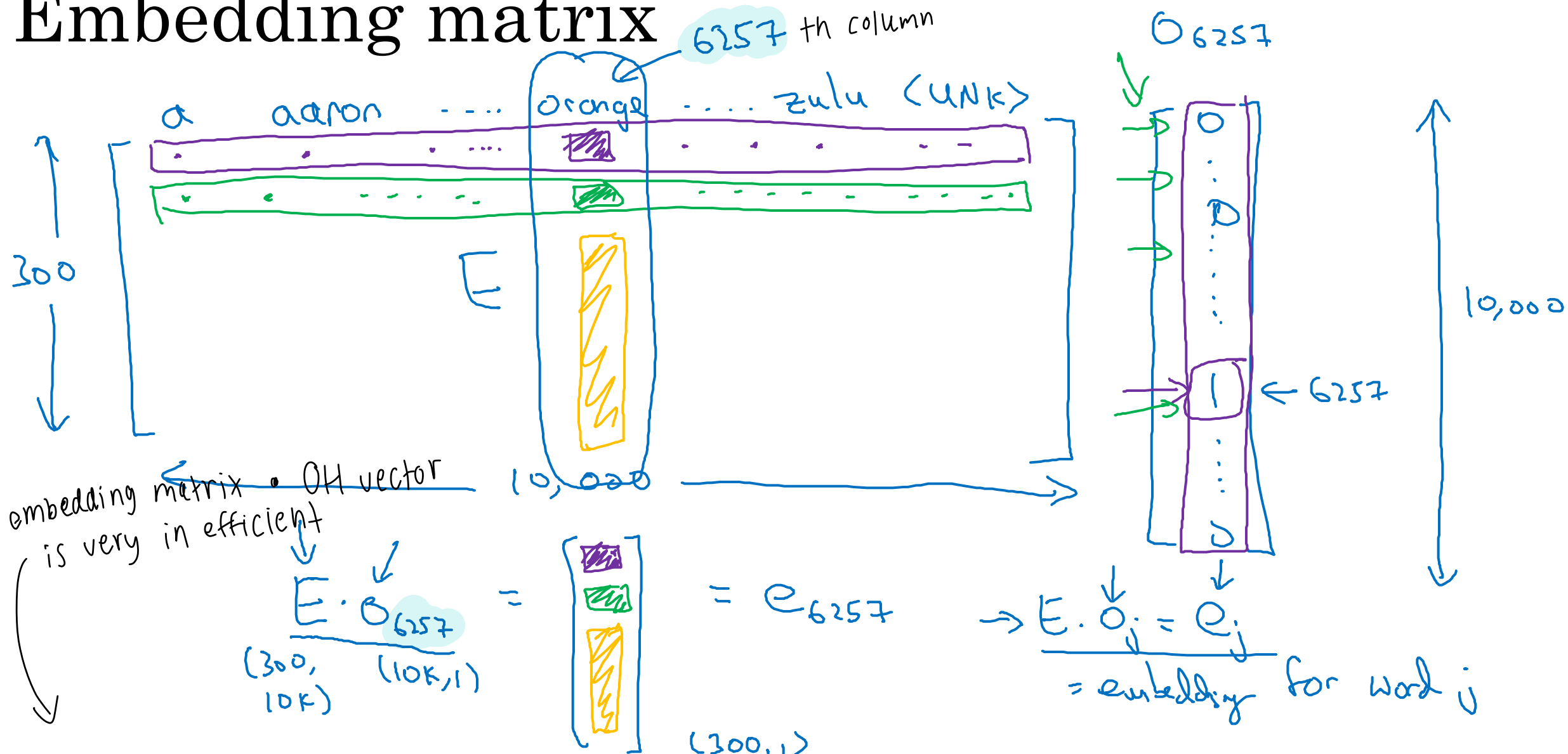
deeplearning.ai

# NLP and Word Embeddings

---

## Embedding matrix

# Embedding matrix



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

$\rightarrow \text{Embedding}$



deeplearning.ai

# NLP and Word Embeddings

---

## Learning word embeddings



# Neural language model

init embedding matrix with random values  
learn the matrix while training neural network

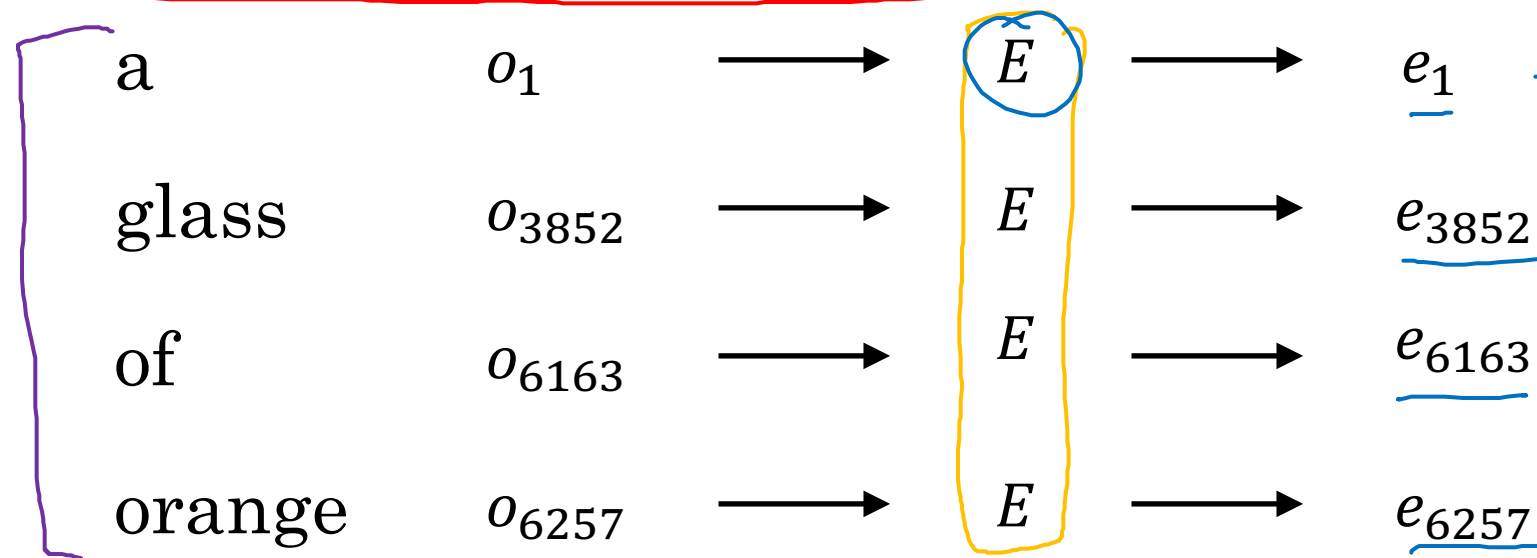
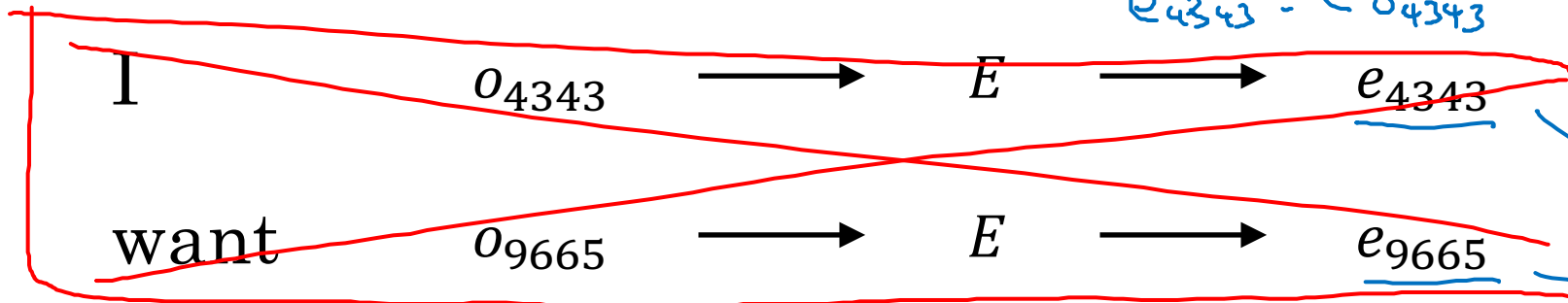
I want a glass of orange juice.

4343 9665 1 3852 6163 6257

$e_{4343} = e_{04343}$

juice.

apple juice.



Softmax

10,000

$w^{(1)}, b^{(1)}$

$w^{(2)}, b^{(2)}$

$\leftarrow 1800 \rightarrow 1200$

# Other context/target pairs

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

The diagram illustrates the context and target for the word 'juice'. A purple bracket under 'a glass of orange' is labeled 'context'. A blue bracket under 'juice' is labeled 'target'. A green arrow points from 'orange' to 'juice', and a blue arrow points from 'juice' to 'to go along with my cereal'.

Context: Last 4 words.

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

a glass of orange ? to go along with

orange ?

glass ?

skip gram



deeplearning.ai

# NLP and Word Embeddings

---

## Word2Vec

# Skip-grams

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



Context

orange

orange

orange



Target

juice

glass

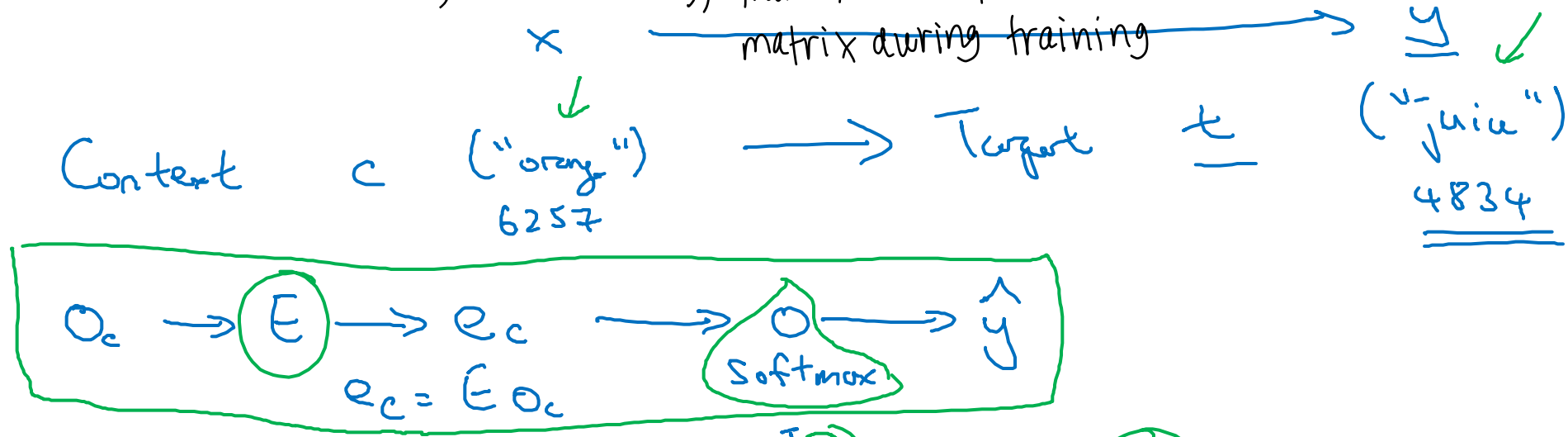
my



# Model

- context target pairs
- 1) randomly pick context words (treat common & rare words evenly)
  - 2) randomly pick target word for the context word (the closer to context word, the more chance to get picked as target)
  - 3) train NN to predict target given context and learn embedding

Vocab size = 10,000k



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

$\Theta_t$  = parameter associated with output  $t$

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

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

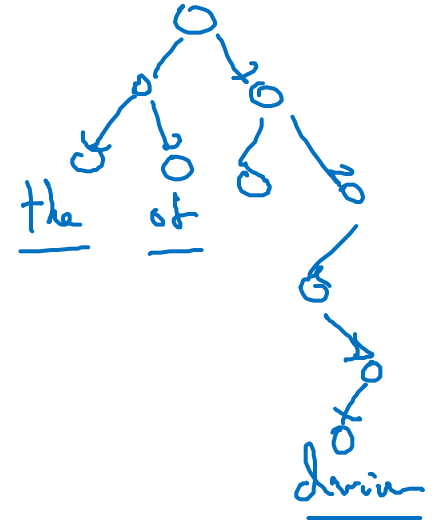
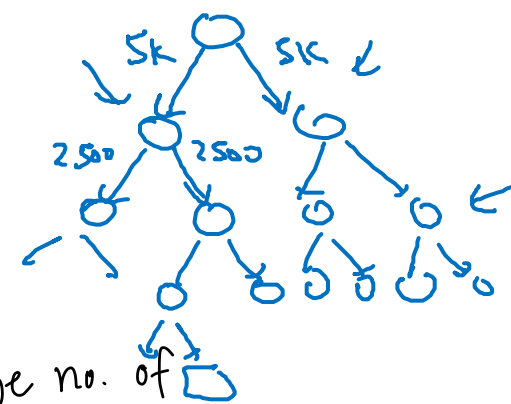
# Problems with softmax classification

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

might need to sum over large no. of words (inefficient)

Hierarchical softmax.

$\log |V|$



How to sample the context  $c$ ?

→ the, of, a, and, to, ...

→ orange, apple, durian

$P_{\text{durian}}$

$P(c)$

$t$   
 $c \rightarrow t$



deeplearning.ai

# NLP and Word Embeddings

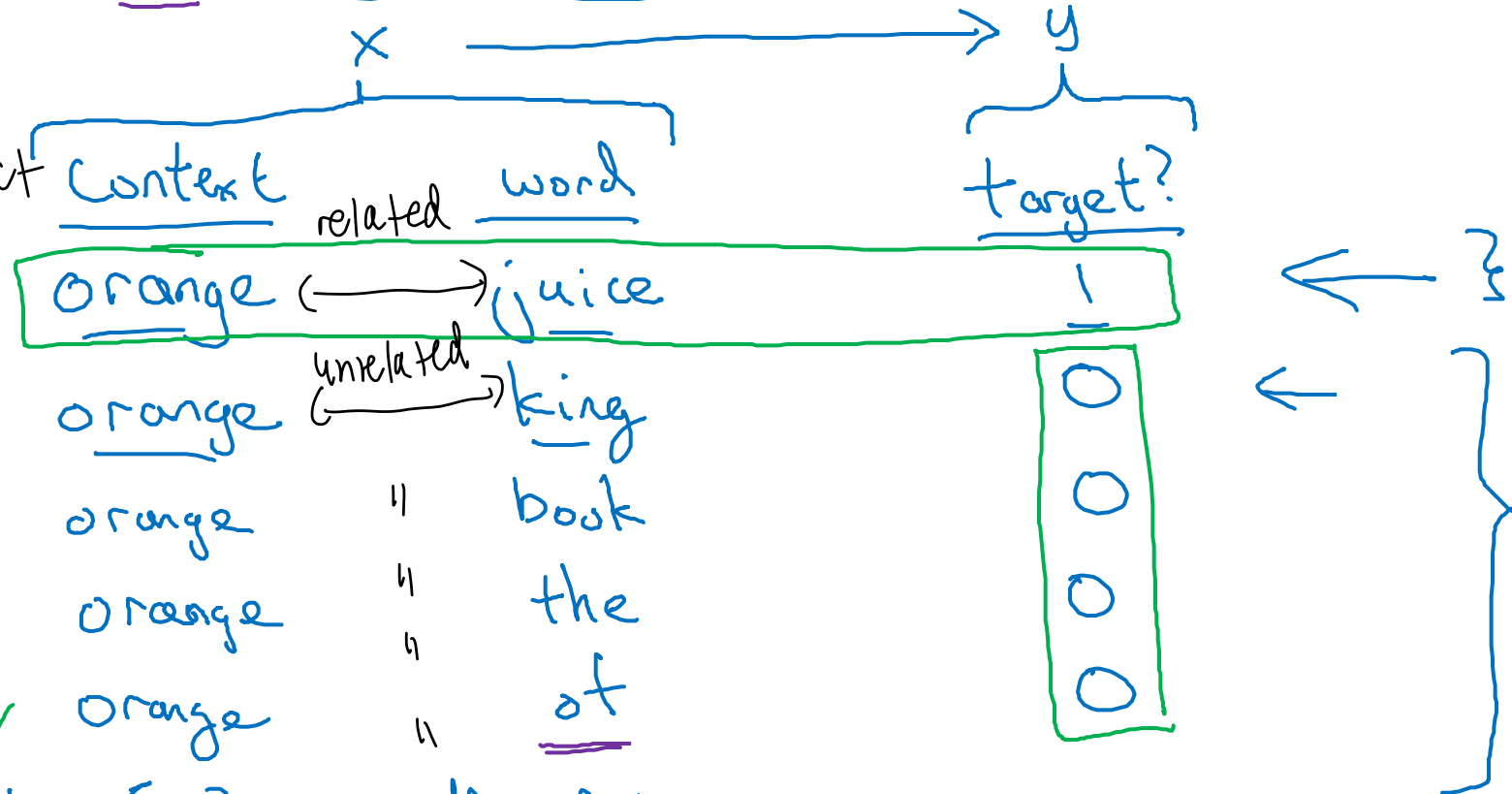
---

## Negative sampling

# Defining a new learning problem

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

- 1) pick context target pair
- 2) add fake examples
- 3) train NN to tell correct context target pair and learn embedding matrix



no. of negative (fake) samples

$$k = 5 - 20$$

$$k = 2 - 5$$

smaller datasets

larger dataset



# Model

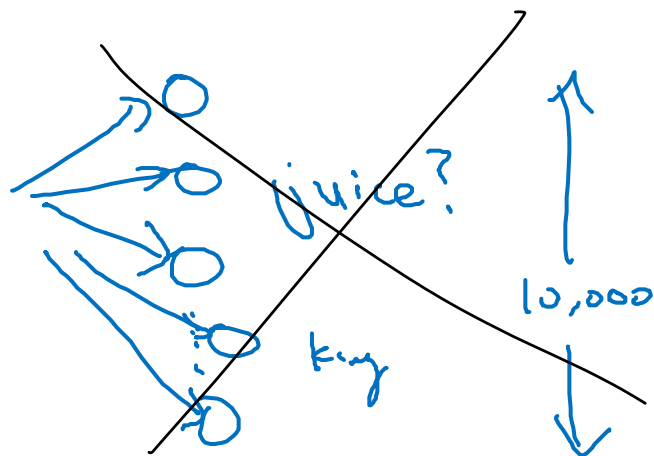
Softmax:

$$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}} \quad \left. \vphantom{\sum_{j=1}^{10,000}} \right\} \begin{array}{l} \text{10,000-way} \\ \text{softmax} \end{array}$$

$$P(y=1 | c, t) = \underbrace{\sigma}_{\text{sigmoid}}(\theta_t^T e_c) \leftarrow$$

Orange  
6257

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



$x$		$y$
context	word	target?
orange	juice	1
orange	king	0
orange	book	0
orange	the	0
orange	of	0
$\uparrow$ $c$	$\uparrow$ $t$	$\uparrow$ $y$

10,000 binary  
classification  
problem  
 $\underline{k+1}$

# 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

the, of, and, ...

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

frequency of word i

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

↑



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.

$c, t$

$X_{ij} = \# \text{ times } \overset{j}{\underset{\substack{\uparrow \\ t}}{x}} \text{ appears in context of } \overset{i}{\underset{\substack{\uparrow \\ c}}{x}}.$

$X_{ij} = X_{ji} \leftarrow$  if consider both sides of target as context  
in glove model, words are related if they appear together frequently

# Model

minimize

$$\sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(x_{ij}) (\underbrace{\Theta_i^T e_j}_{\substack{t \quad c \\ \text{"}\Theta_t^T e_c\text{"}}} + b_i + b_j' - \log x_{ij})^2$$

wrong (see next page)

0?

weighting  
term

$$f(x_{ij}) = 0 \text{ at } x_{ij} = 0.$$

$$0 \log 0 = 0$$

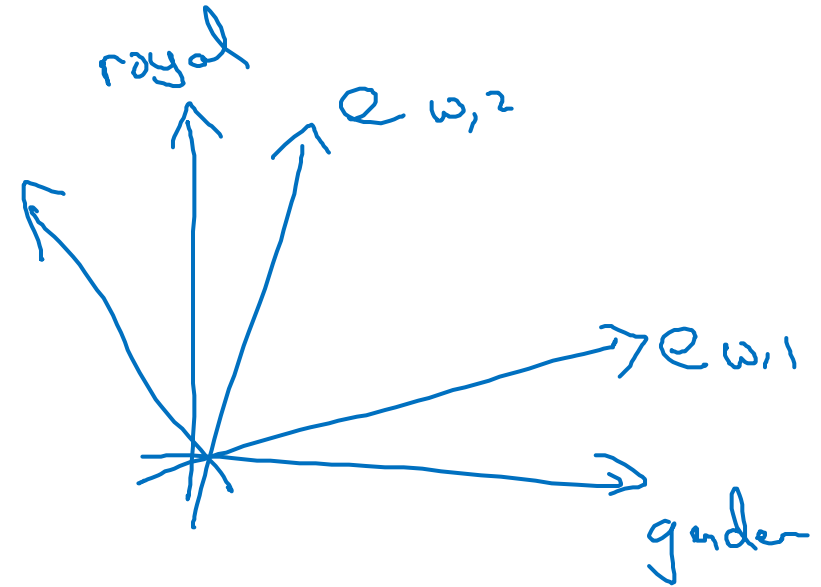
this, is, of, a, ...  
derivation

$\Theta_i, e_j$  are symmetric

$$e_w^{(final)} = \frac{e_w + \Theta_w}{2}$$

# 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	←



objective function

no. of unique words in dictionary

correct formula

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

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



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  $\rightarrow$  100,000 words

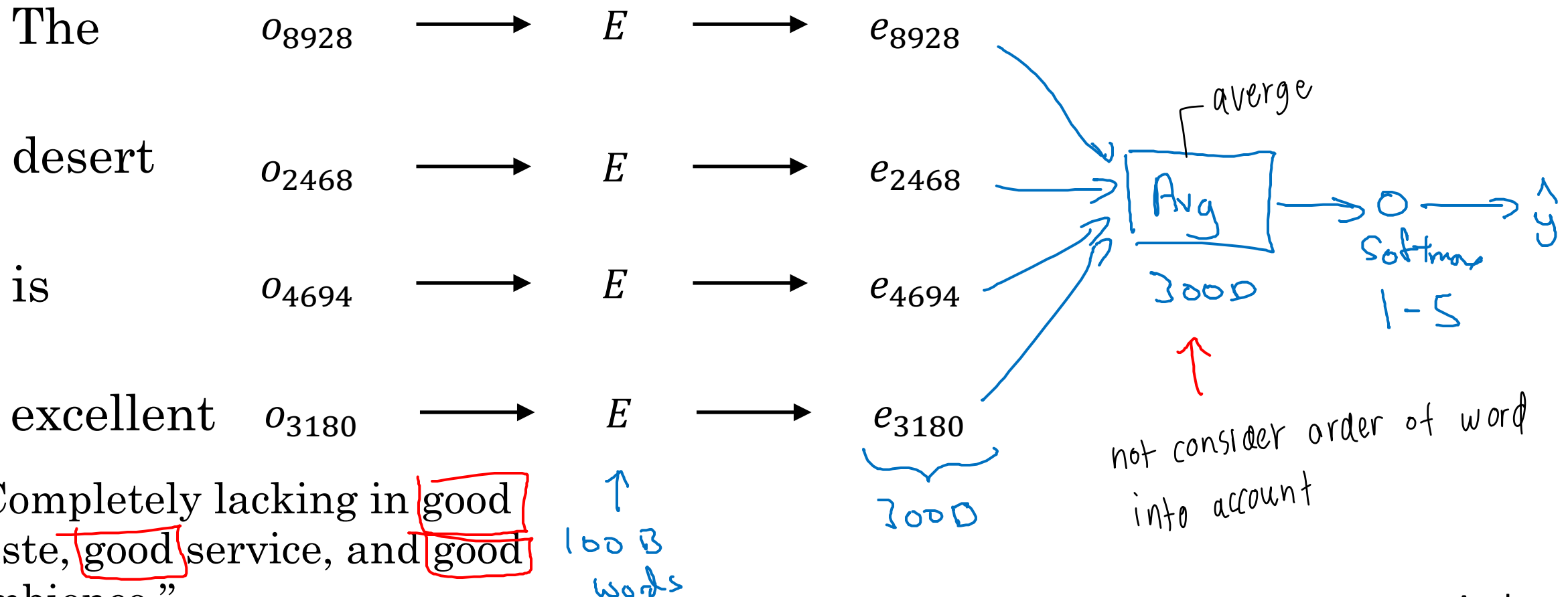


# Simple sentiment classification model

The dessert is excellent

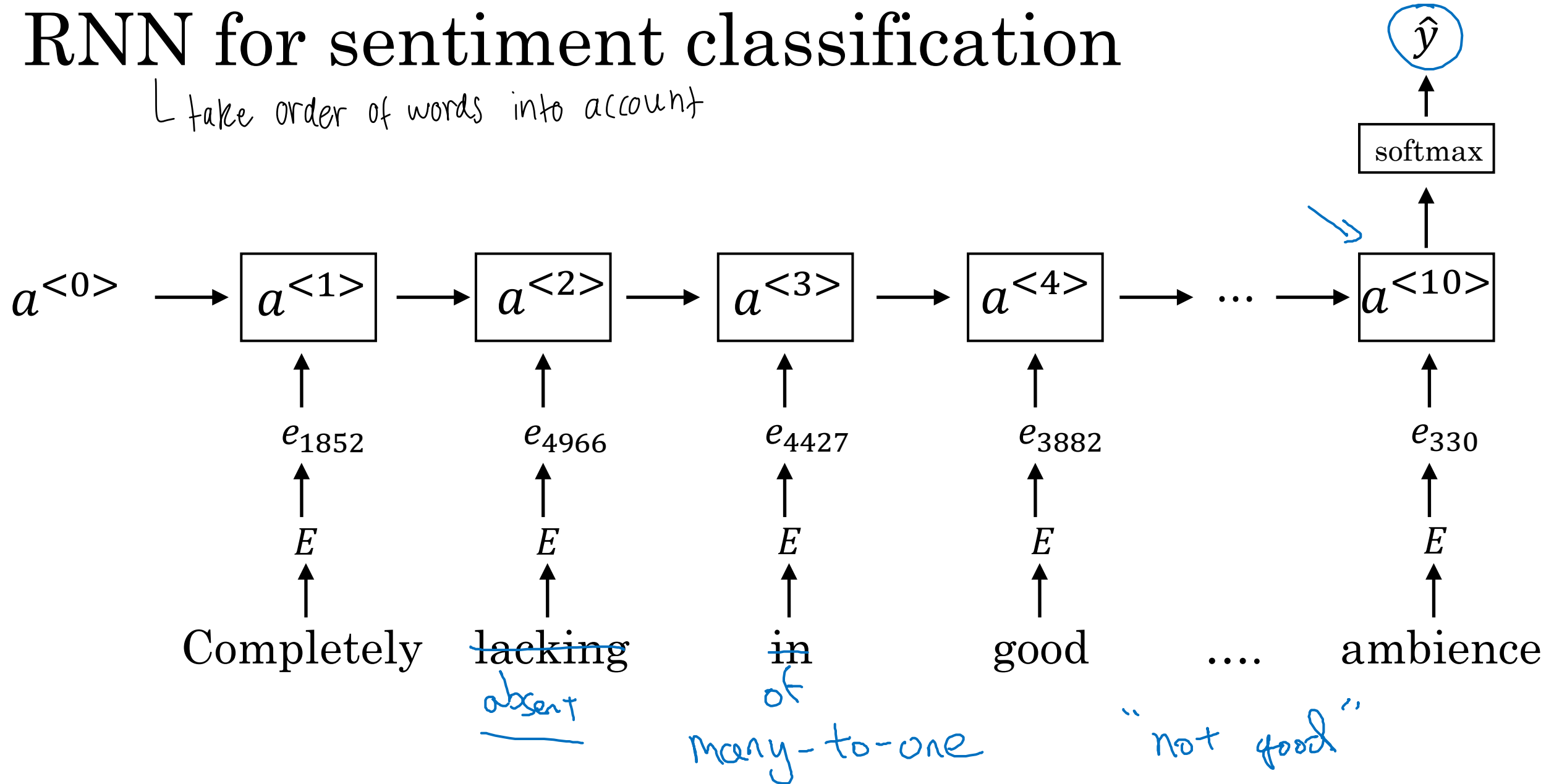


8928 2468 4694 3180



# RNN for sentiment classification

↳ take order of words into account





deeplearning.ai

# NLP and Word Embeddings

---

## Debiasing word embeddings

# The problem of bias in word embeddings

Man:Woman as King:Queen

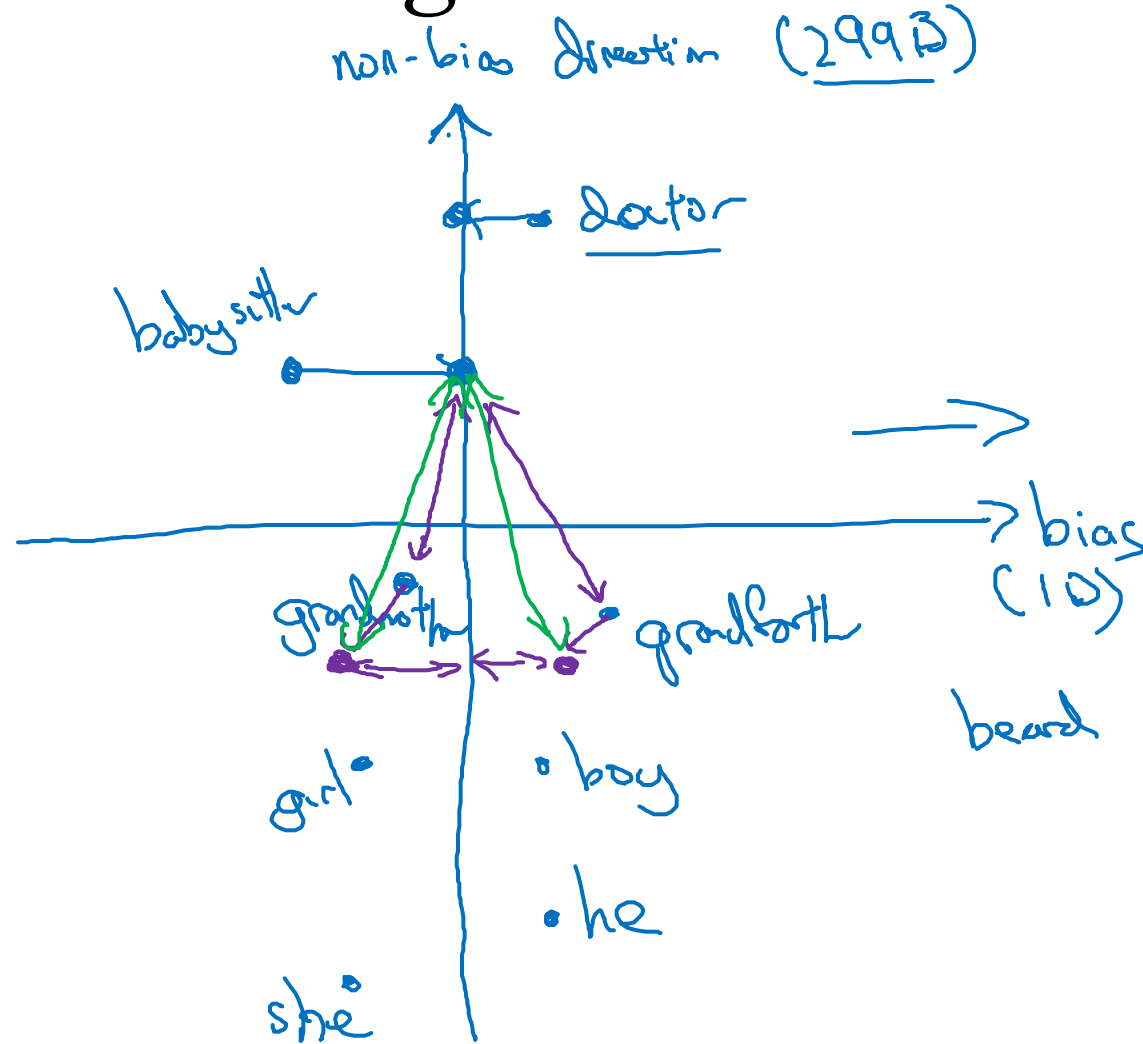
Man:Computer\_Programmer as Woman:Homemaker X

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.

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