

# 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

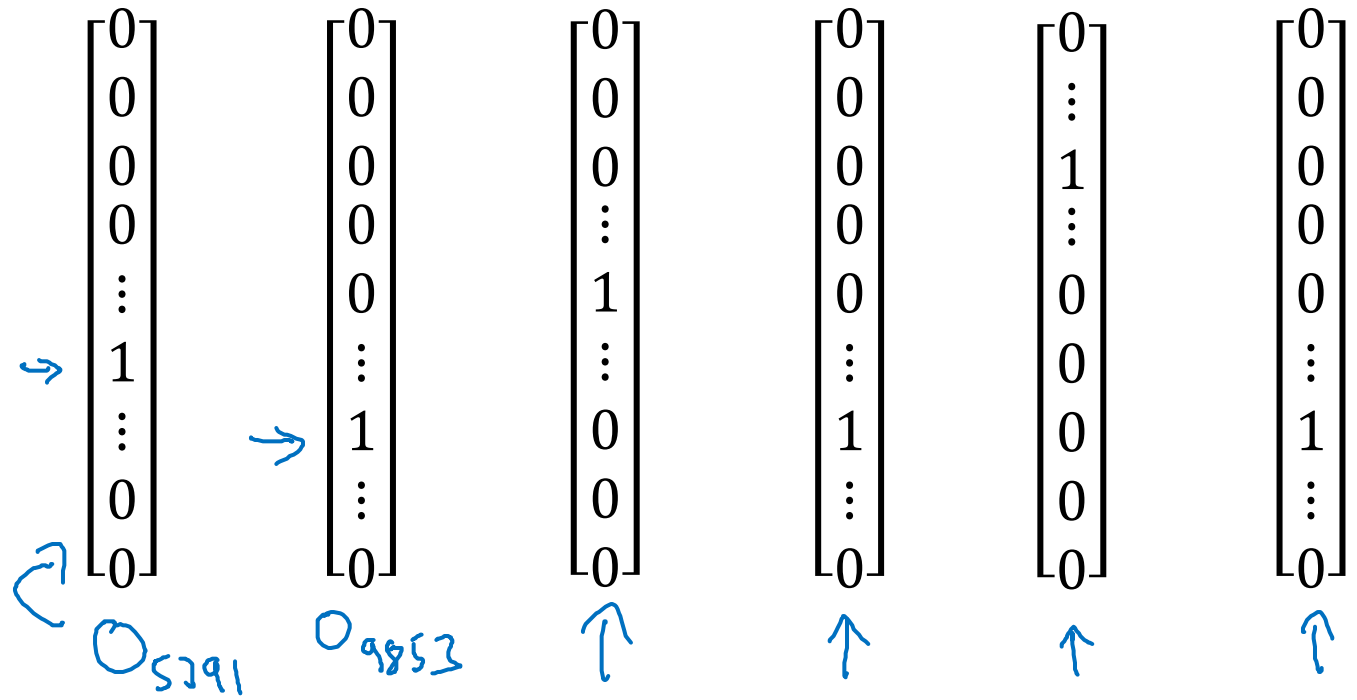
# Word representation

$V = [a, aaron, \dots, zulu, \langle \text{UNK} \rangle]$

$|V| = 10,000$

## 1-hot representation

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



I want a glass of orange juice.

I want a glass of apple \_\_\_\_\_?.

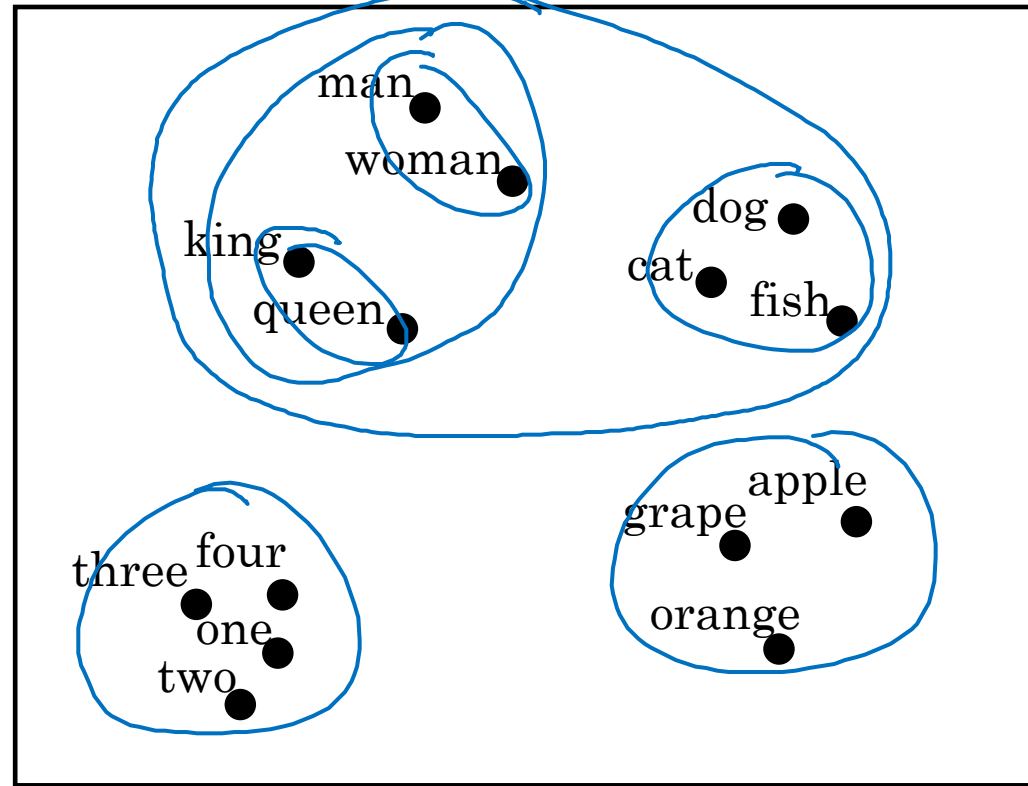
# Featurized representation: word embedding

	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						

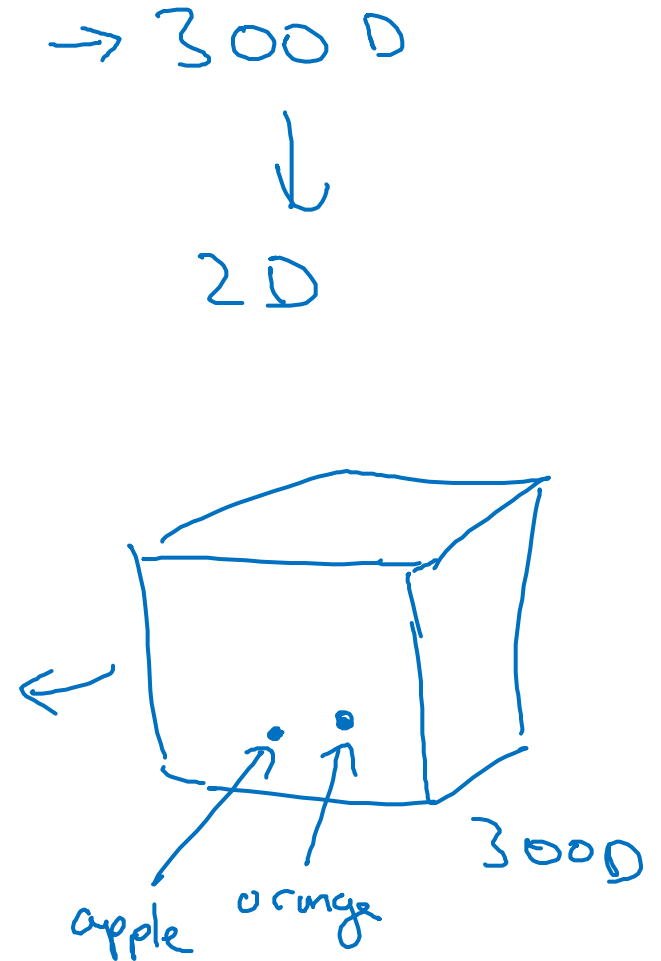
I want a glass of orange juice.  
 I want a glass of apple juice.

Andrew Ng

# Visualizing word embeddings



t-SNE





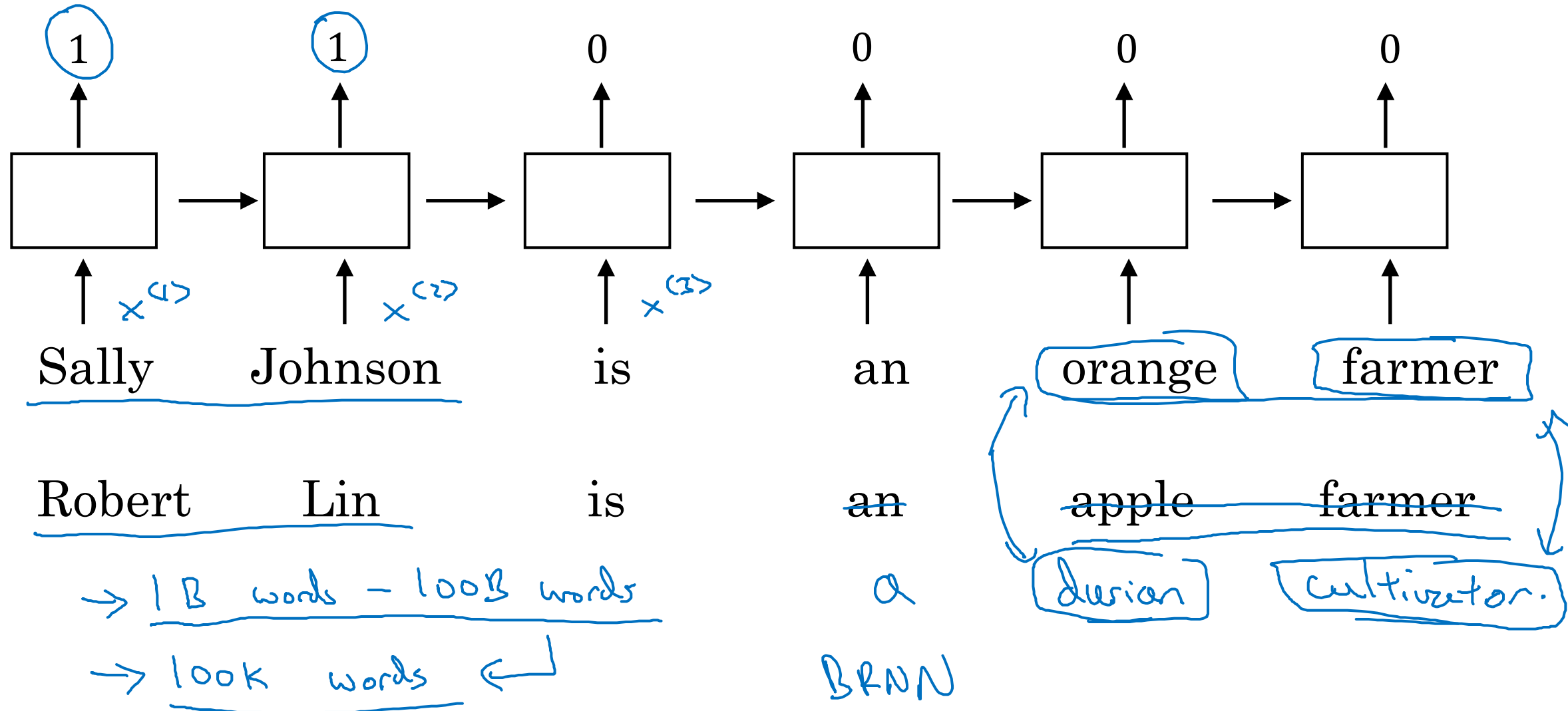
deeplearning.ai

# NLP and Word Embeddings


---

## Using word embeddings

# Named entity recognition example



# Transfer learning and word embeddings

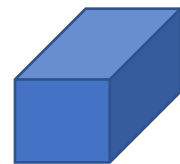
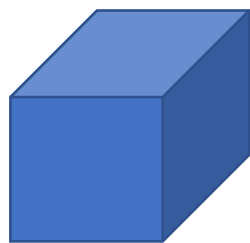
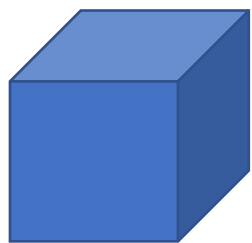
- 
1. Learn word embeddings from large text corpus. (1-100B words)  
(Or download pre-trained embedding online.)
  2. Transfer embedding to new task with smaller training set.  
(say, 100k words) → 10,000 → 300
  3. Optional: Continue to finetune the word embeddings with new data.



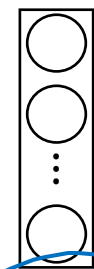
# Relation to face encoding (embedding) 128D



$x^{(i)}$



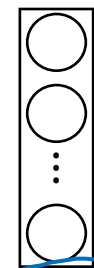
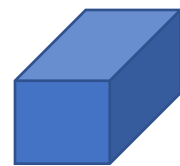
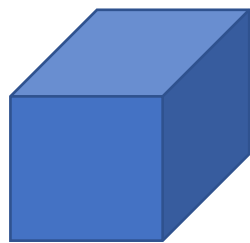
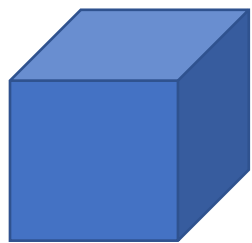
128D



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



$x^{(j)}$

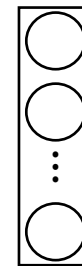


$f(x^{(j)})$

128D

$|V| = 10,000$

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



$\hat{y}$



deeplearning.ai

# NLP and Word Embeddings

---

## Properties of word embeddings

# 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

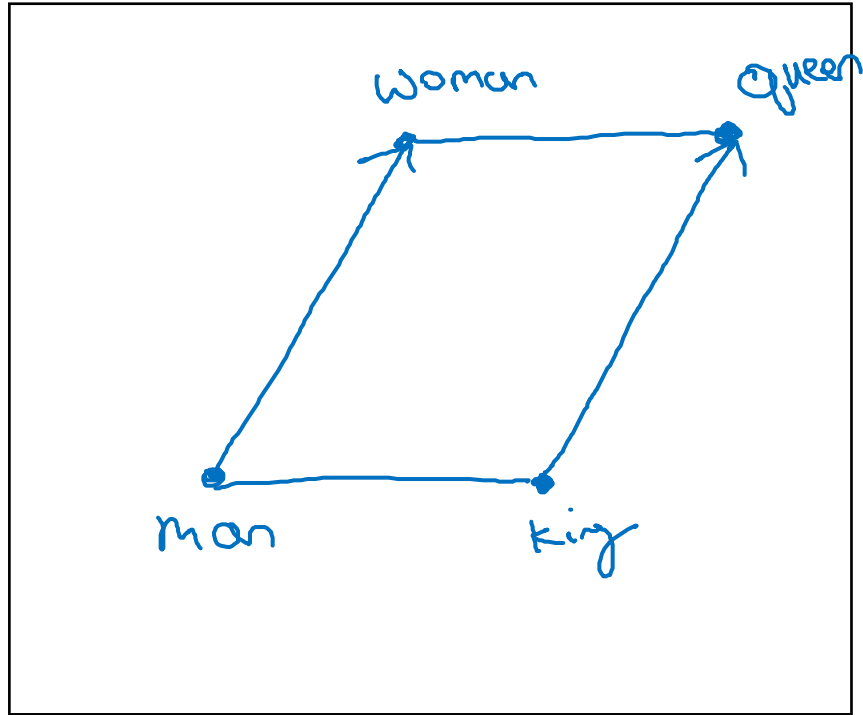
$$\underbrace{e_{5391}}_{e_{\text{man}}} - \underbrace{e_{9853}}_{e_{\text{woman}}} \approx \underbrace{e_{4914}}_{e_{\text{king}}} - \underbrace{e_{7157}}_{e_{\text{queen}}}$$

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

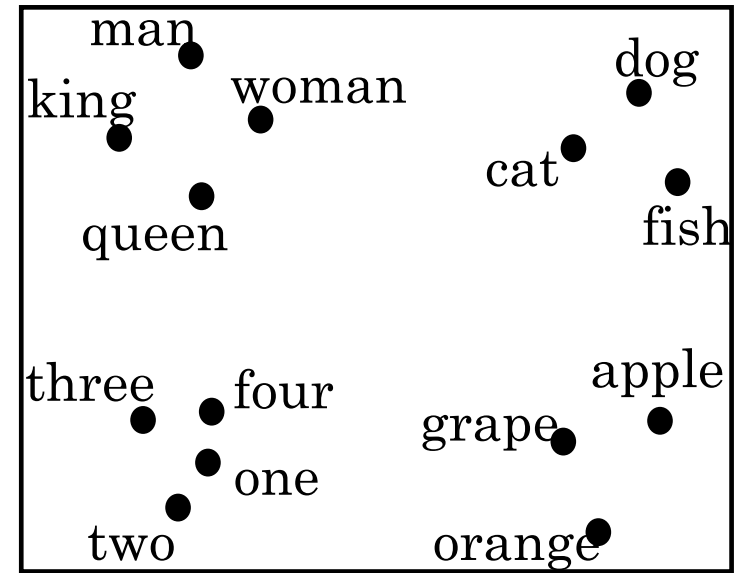
Man  $\rightarrow$  Woman  $\approx$  King  $\rightarrow$  ? Queen



# Analogies using word vectors



300D → 20  
↑



t-SNE

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

300 D

Find word  $w$ :  $\arg \max_w$

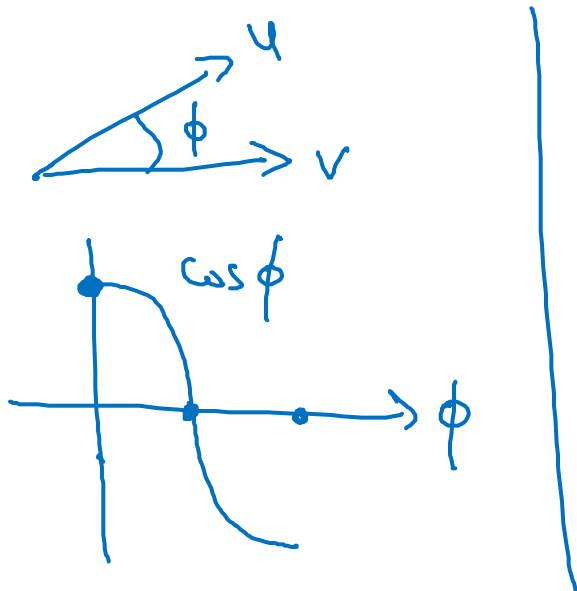
$$\text{Sim}(e_w, e_{king} - e_{man} + e_{woman})$$

30 - 75%

# 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



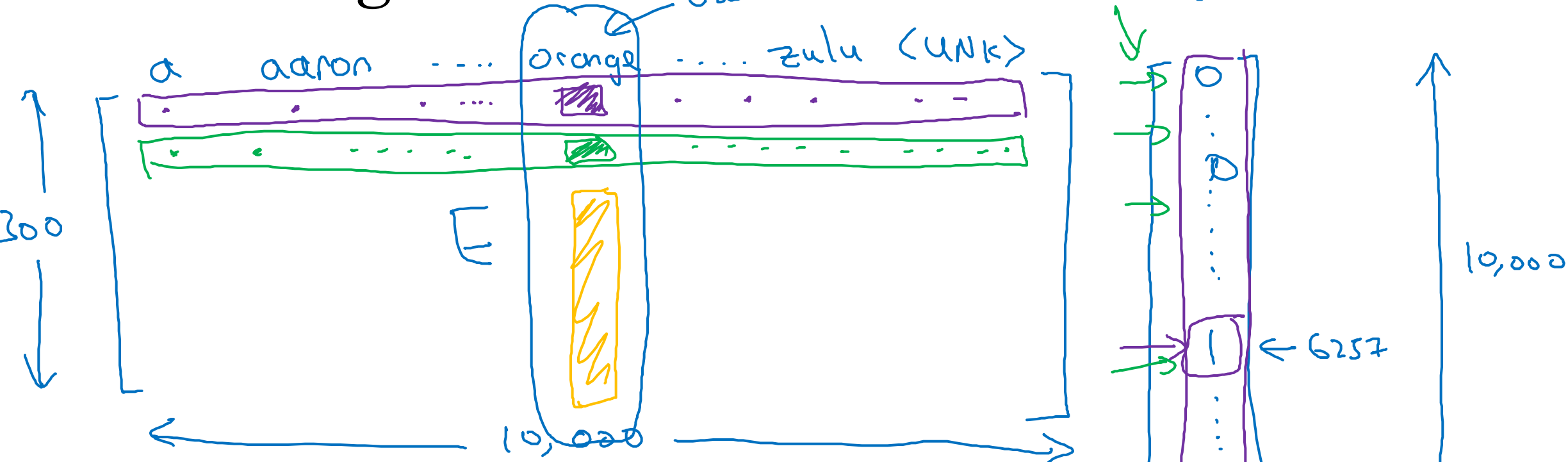
deeplearning.ai

# NLP and Word Embeddings

---

## Embedding matrix

# Embedding matrix



$$\begin{matrix} \downarrow \\ \mathbf{E} \cdot \mathbf{O}_{6257} \\ \begin{matrix} (300, \\ 10K) \end{matrix} \end{matrix} = \begin{matrix} \begin{bmatrix} \text{purple box} \\ \text{green box} \\ \text{orange box} \end{bmatrix} \\ (300, 1) \end{matrix} = \mathbf{e}_{6257} \rightarrow \mathbf{E} \cdot \mathbf{O}_j = \mathbf{e}_j = \text{embedding for word } j$$

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



deeplearning.ai

# NLP and Word Embeddings

---

## Learning word embeddings



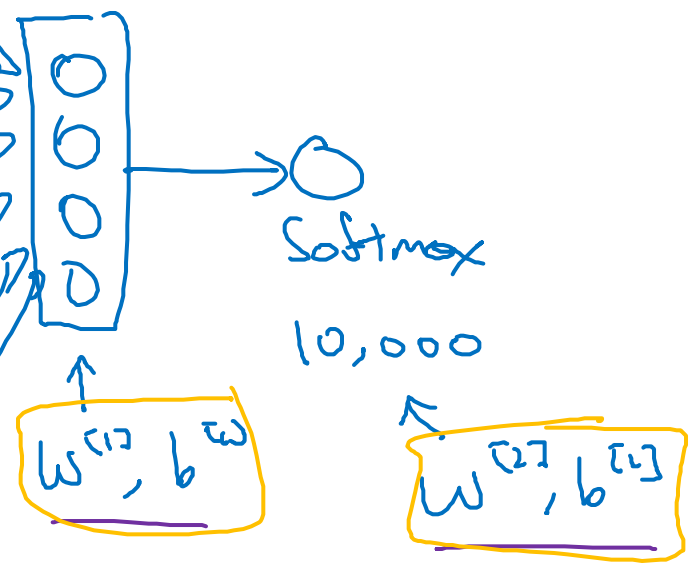
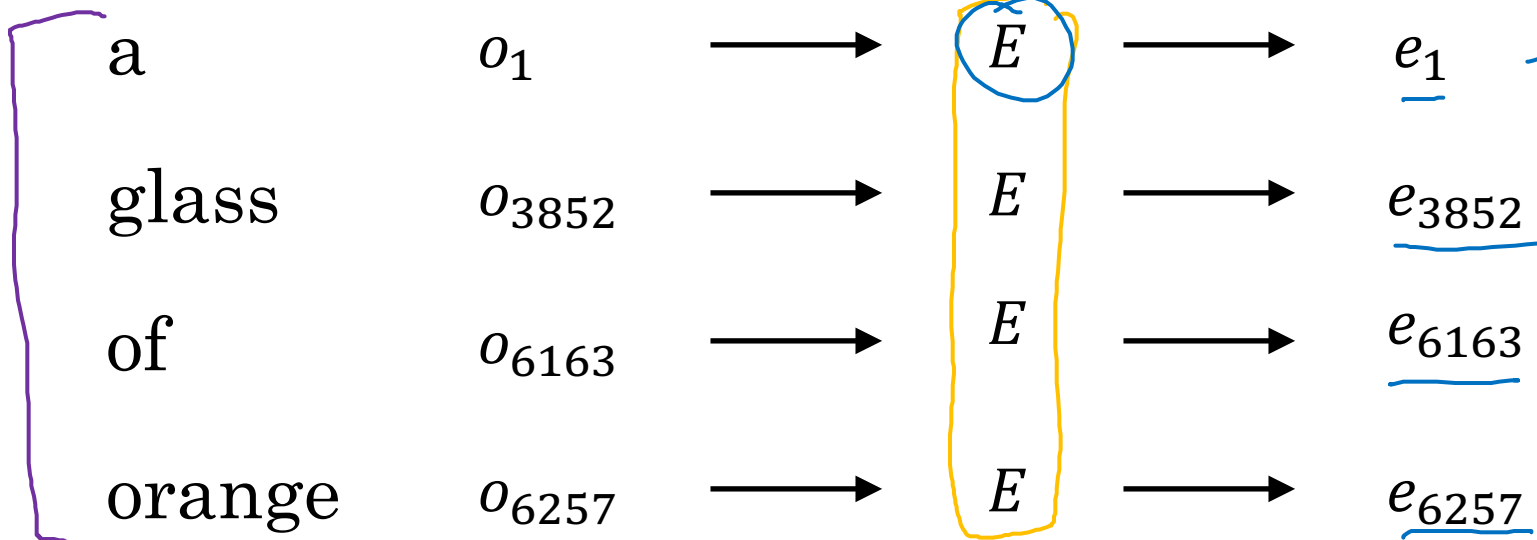
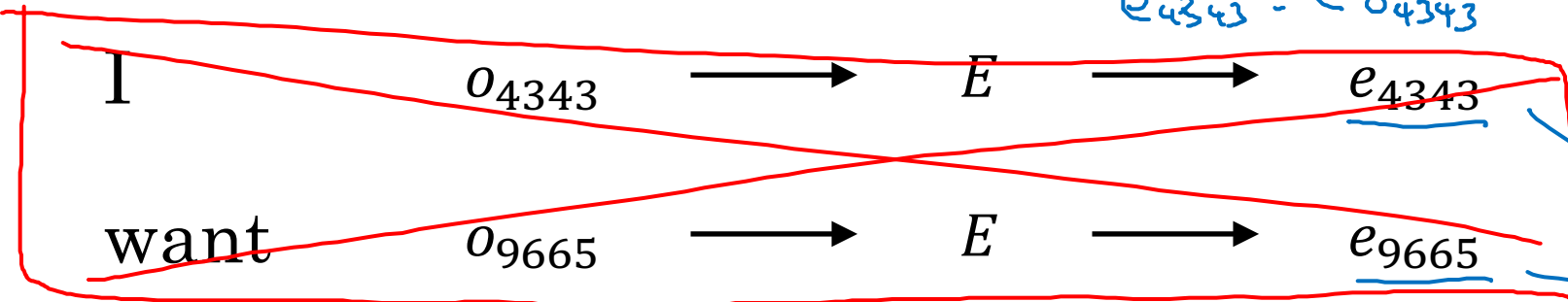
# Neural language model

I want a glass of orange juice.

4343 9665 1 3852 6163 6257

$e_{4343} = E_{04343}$

juice.  
apple juice.



~~1800~~ 1200

# Other context/target pairs

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

*Context* (purple bracket under "a glass of orange")  
*target* (blue arrow pointing to "juice")

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

Target

orange

juice

orange

glass

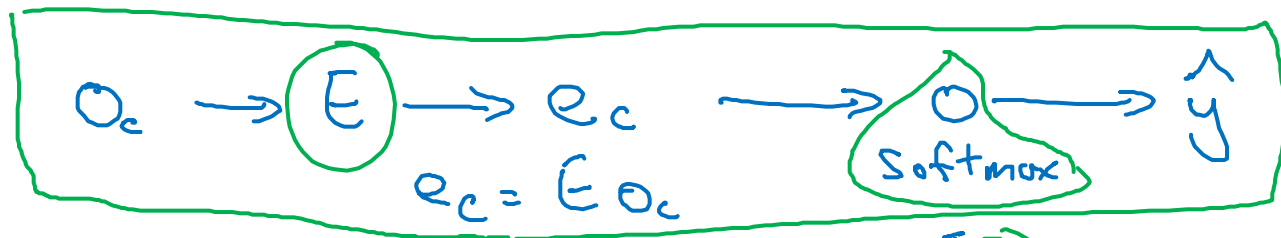
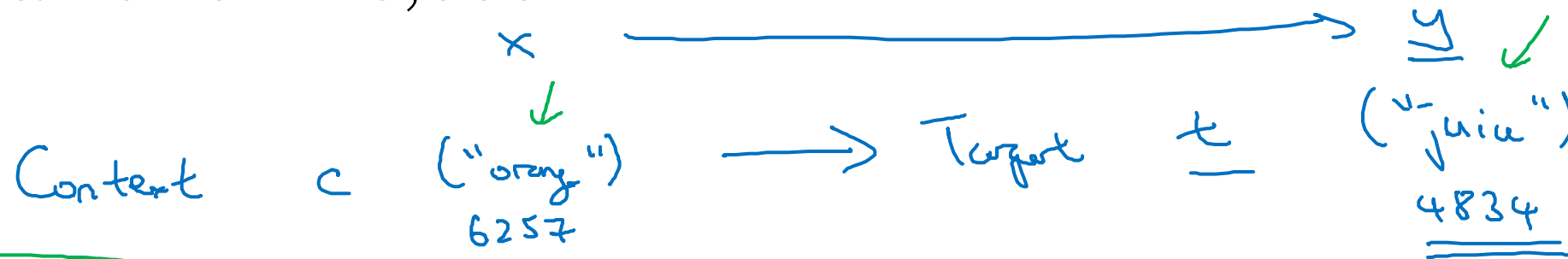
orange

my



# Model

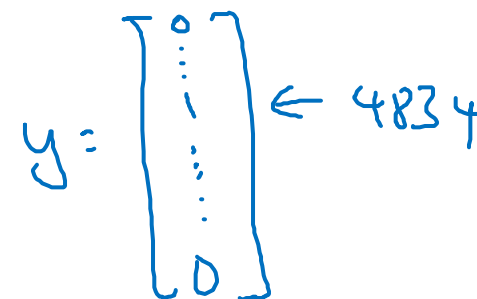
Vocab size = 10,000k



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

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

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

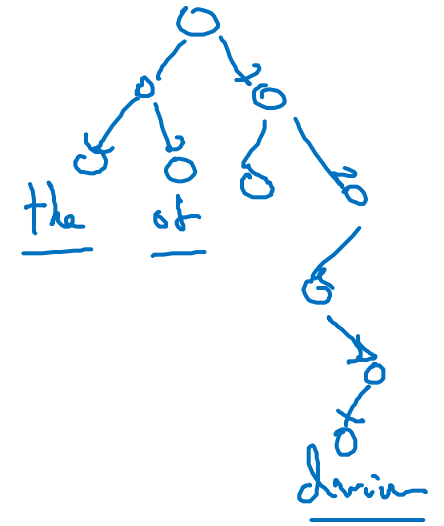
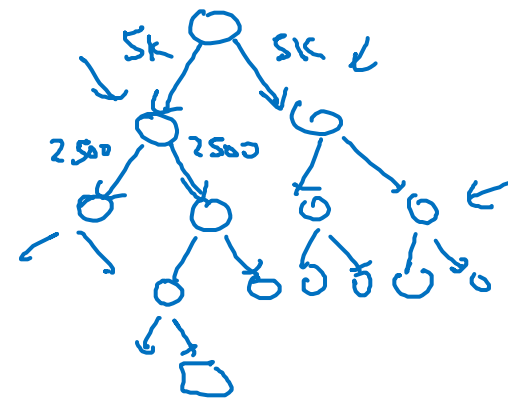


# Problems with softmax classification

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

Hierarchical softmax.

$\log |V|$



How to sample the context  $c$ ?

→ the, of, a, and, to, ...

→ orange, apple, divin

$P_{divin}$

$P(c)$

$t$   
 $c \rightarrow t$



deeplearning.ai

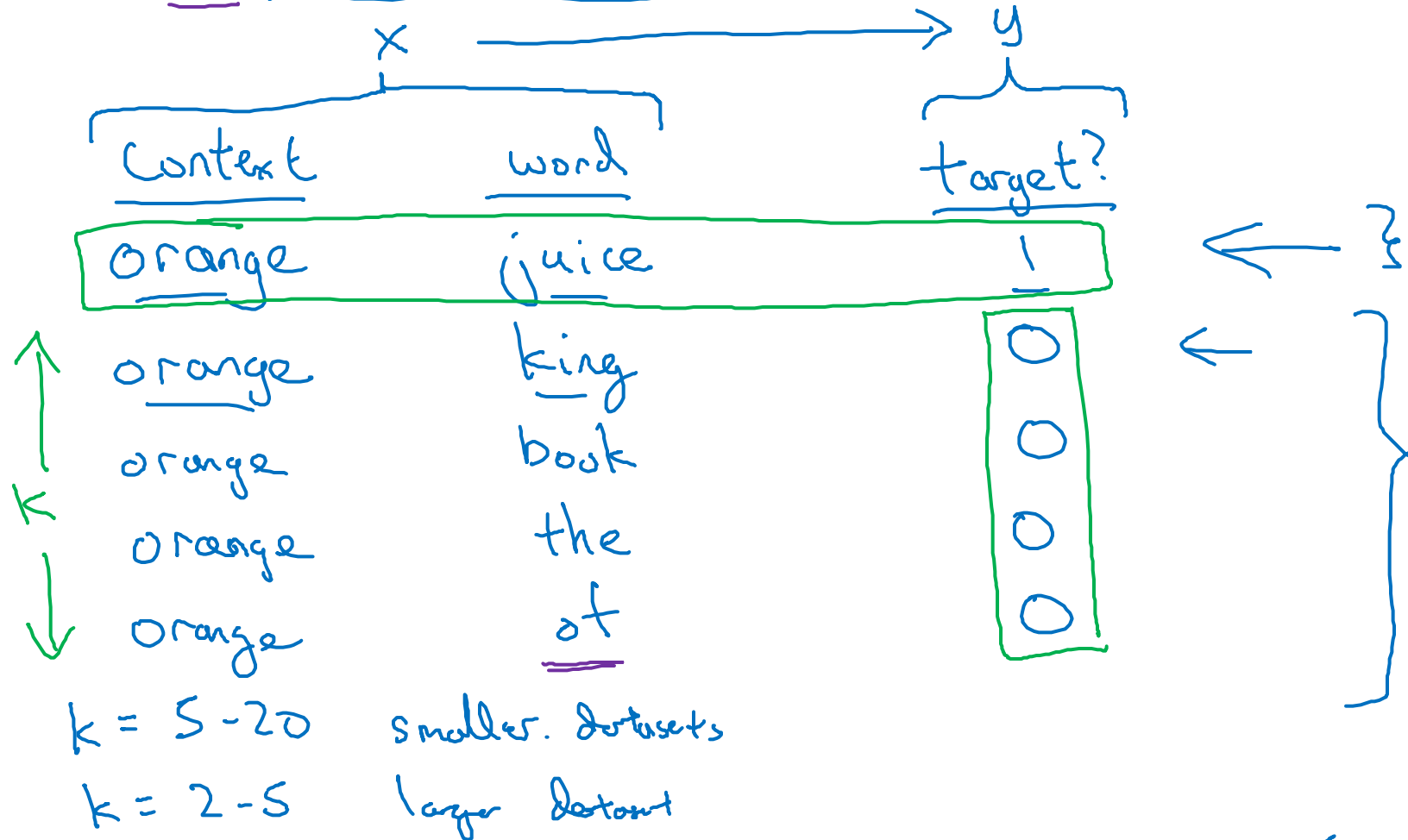
# NLP and Word Embeddings

---

## Negative sampling

# Defining a new learning problem

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





# Model

Softmax:

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

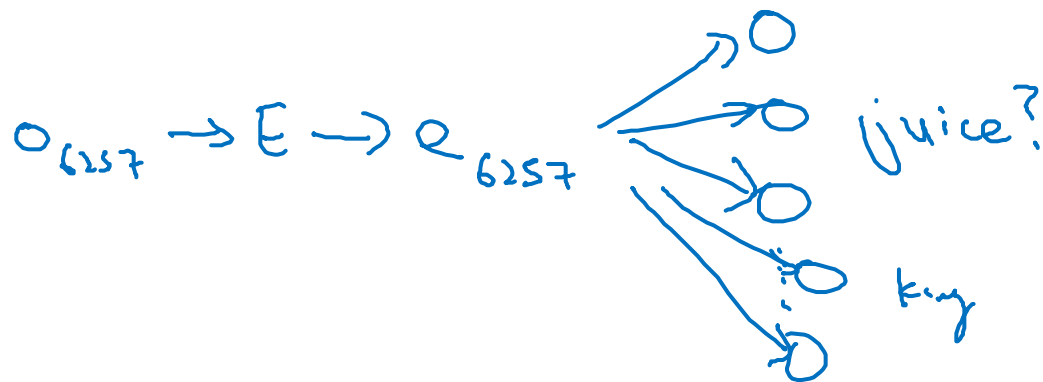
*10,000-way softmax*

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

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

↑ c
↑ t
↑ y

Orange  
6257



10,000  
10,000 binary classification problem  
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

↑  
t

the, of, and, ...

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

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

↑



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.

$c, t$

$X_{ij}$  = # times  $i$  appears in context of  $j$ .

$\begin{matrix} \uparrow & \uparrow & & \uparrow \\ c & t & & c \end{matrix}$

$X_{ij} = X_{ji}$  ←

# Model

Minimize

$$\sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(x_{ij}) \left( \underbrace{\Theta_i^T e_j}_{\substack{t \quad c \\ \text{"}\Theta_t^T e_c\text{"}}} + b_i + b_j - \log x_{ij} \right)^2$$

←

weighting term

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

" $0 \log 0$ " = 0

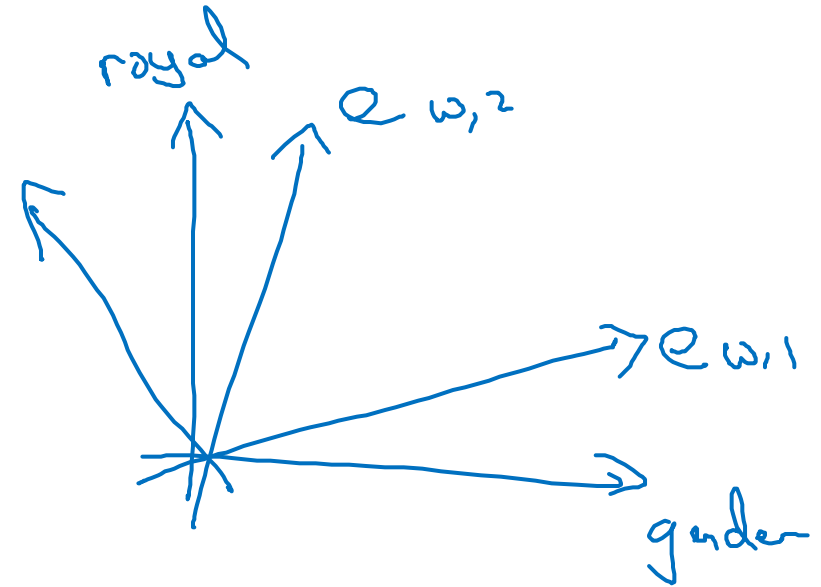
→ this, is, of, a, ...  
 → derivation

$\Theta_i, e_j$  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	←



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

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



deeplearning.ai

# NLP and Word Embeddings

---

## Sentiment classification

# Sentiment classification problem

$x$



$y$

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

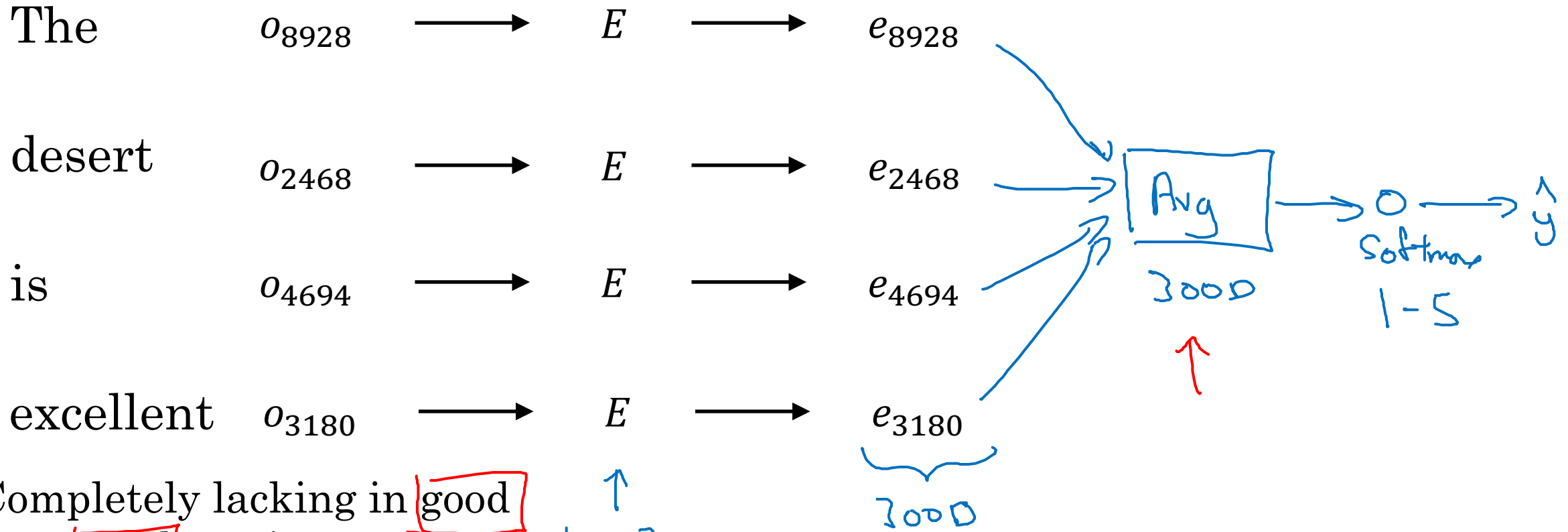


# Simple sentiment classification model

The dessert is excellent



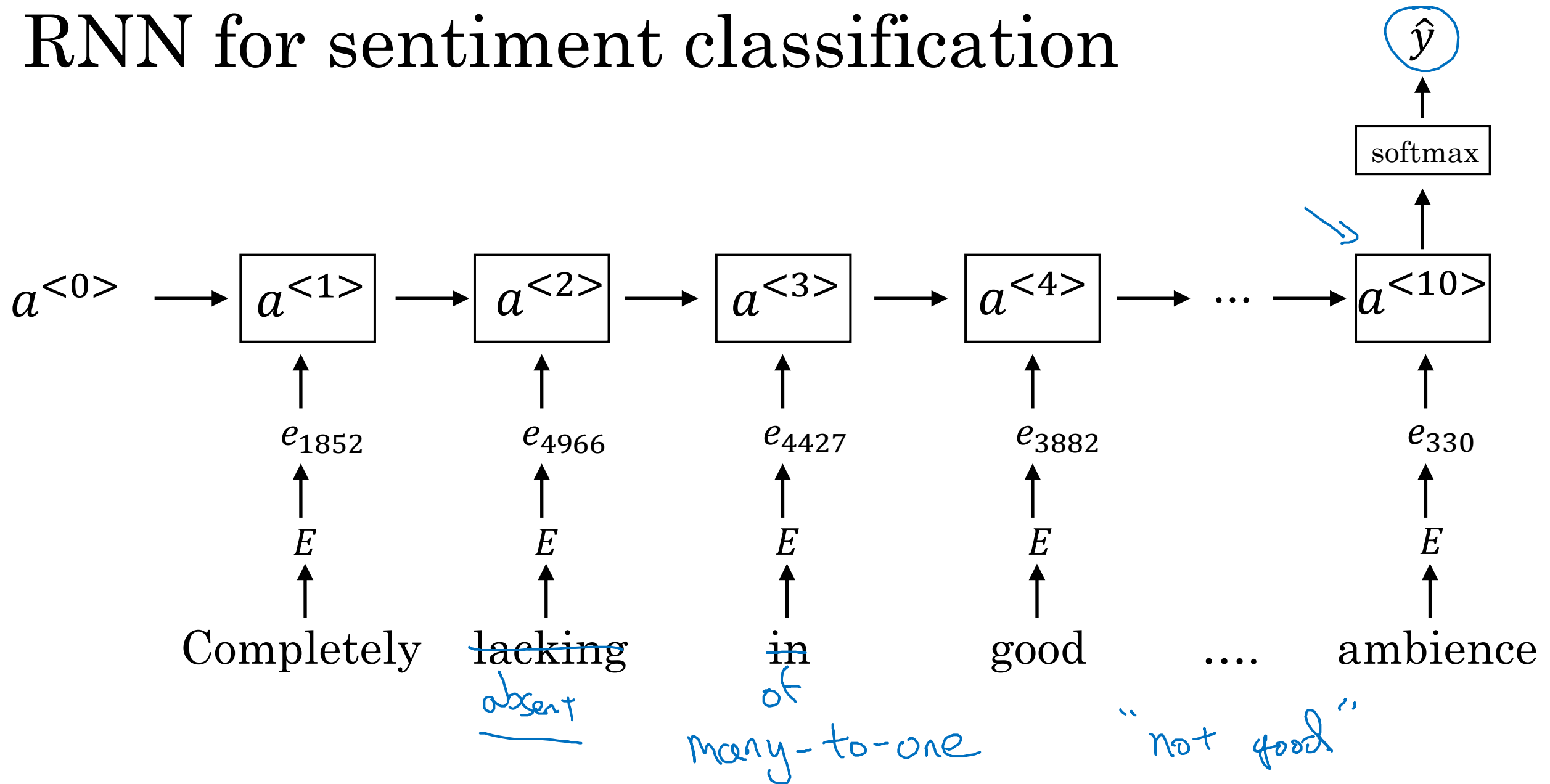
8928 2468 4694 3180



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

↑  
100 B  
words

# RNN for sentiment classification





deeplearning.ai

# NLP and Word Embeddings

---

## Debiasing word embeddings

# The problem of bias in word embeddings

Man:Woman as King:Queen

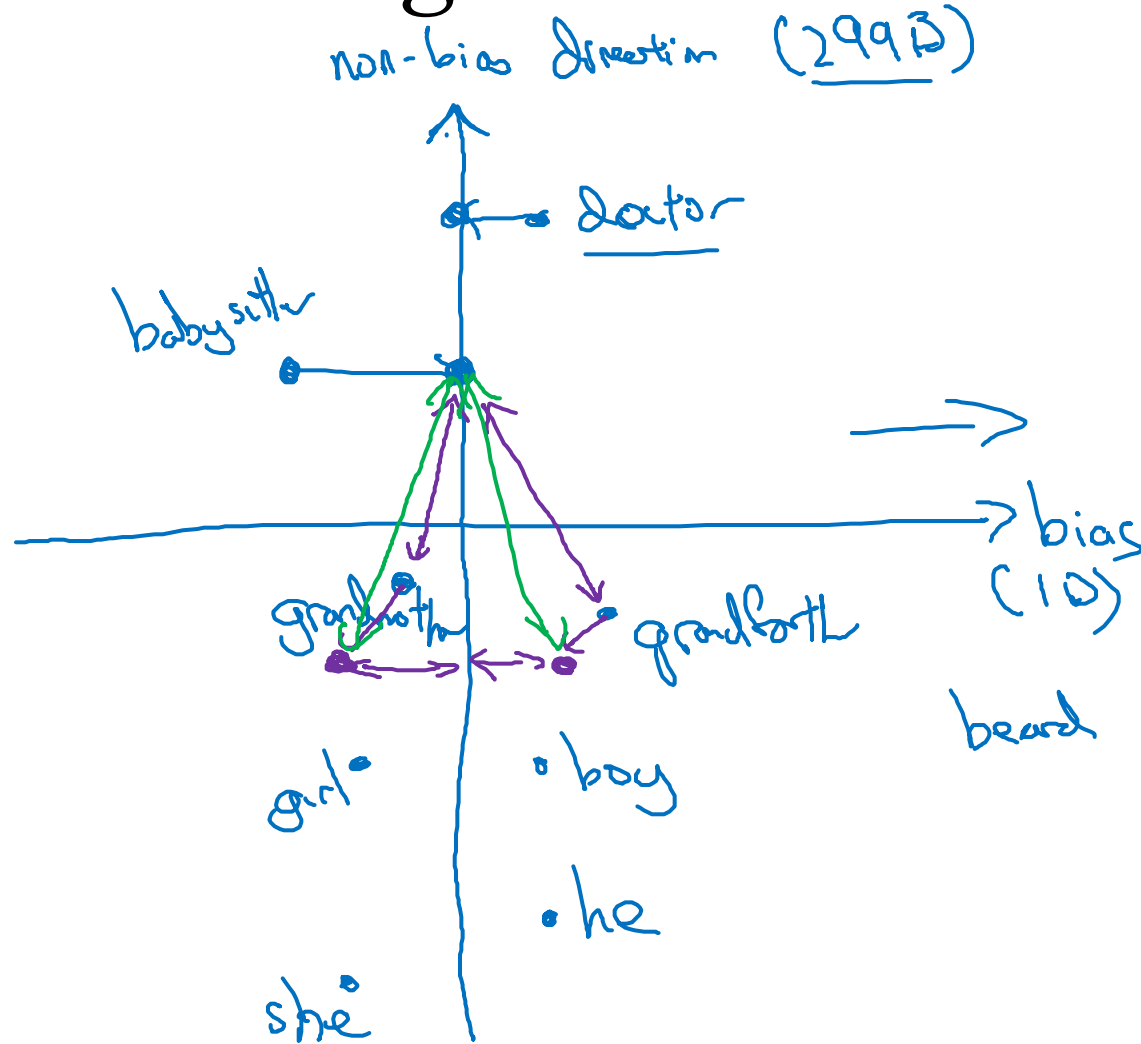
Man:Computer\_Programmer as Woman:Homemaker ✗

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.

$\{ \begin{array}{l} e_{he} - e_{she} \\ e_{male} - e_{female} \\ \vdots \end{array} \}$   
→ average

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

3. Equalize pairs.

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