

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

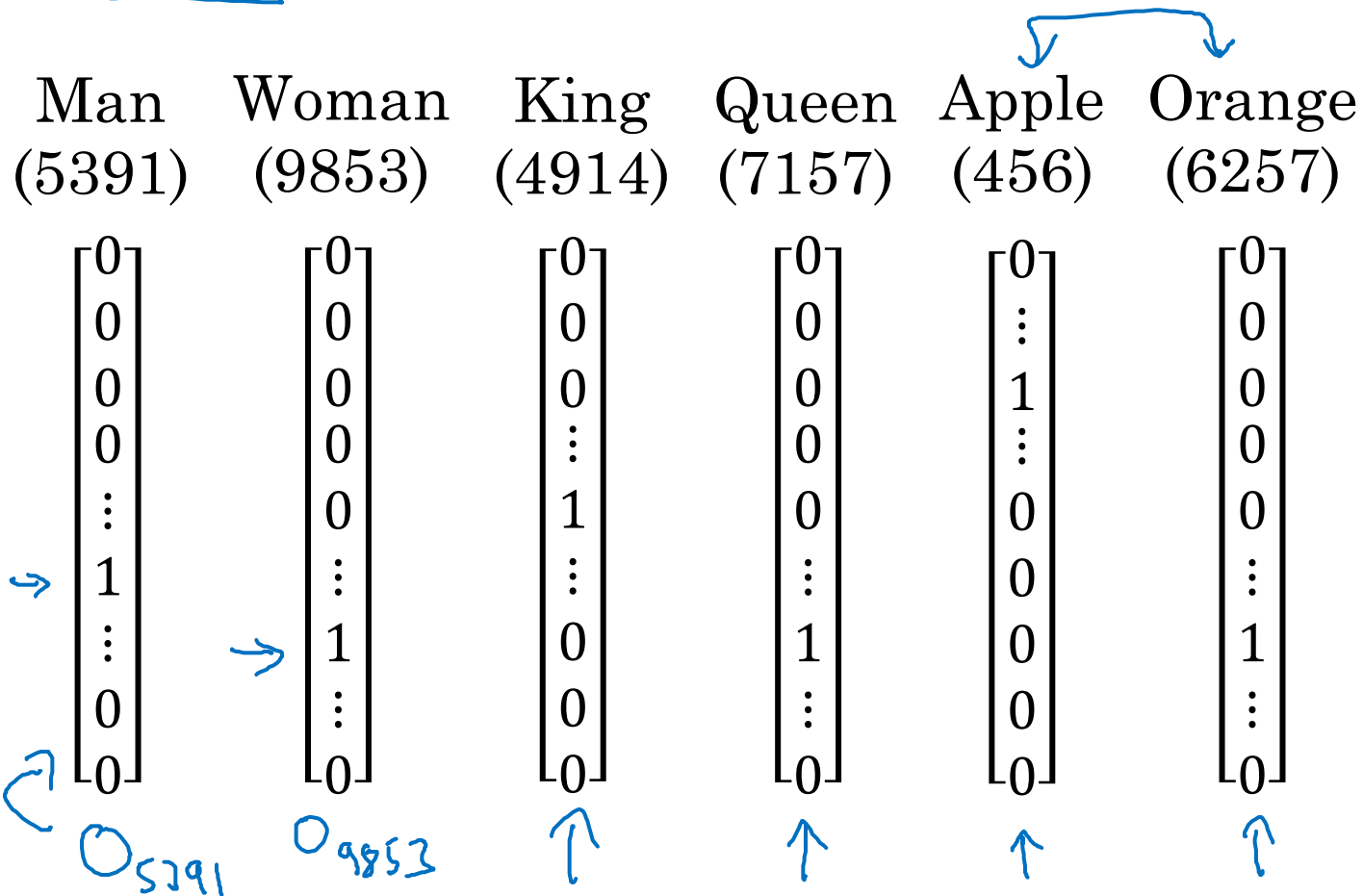
Word representation

Word representation

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

$|V| = 10,000$

1-hot representation



I want a glass of orange juice.

I want a glass of apple ?.

Featurized representation: word embedding

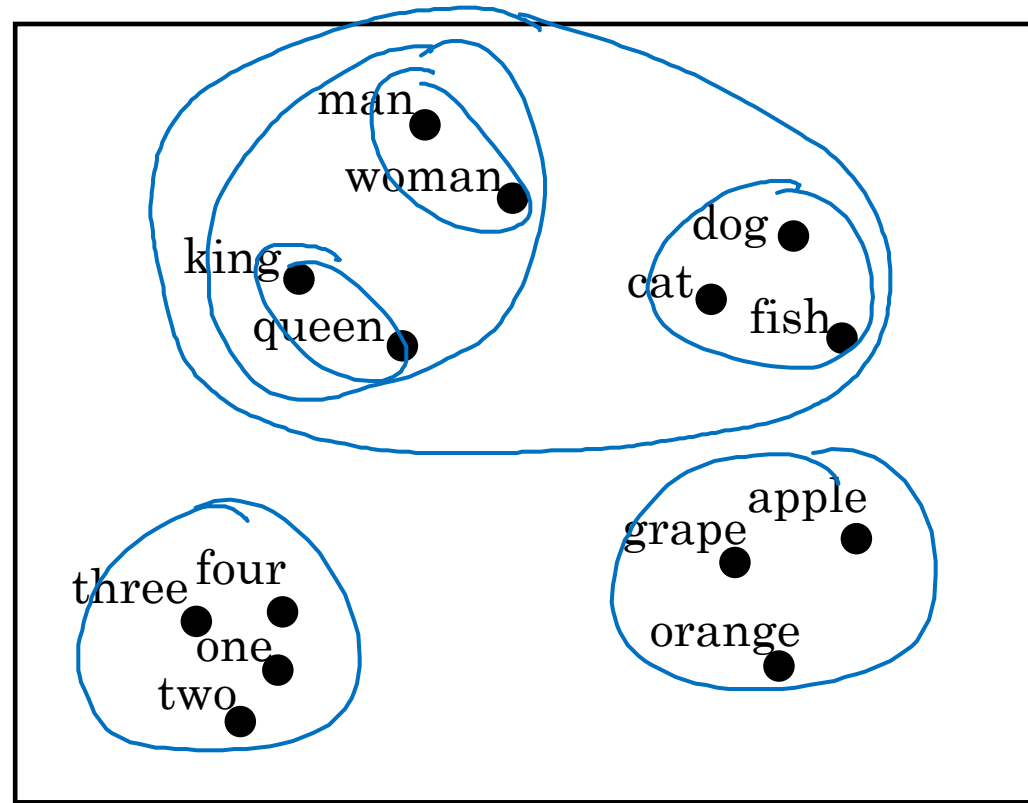
	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	<u>0.93</u>	<u>0.95</u>	-0.01	0.00
Age	0.03	0.02	0.7	0.69	0.03	-0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97
...				
size						
cost						
alive						
verb						

I want a glass of orange juice.

I want a glass of apple juice.

Andrew Ng

Visualizing word embeddings



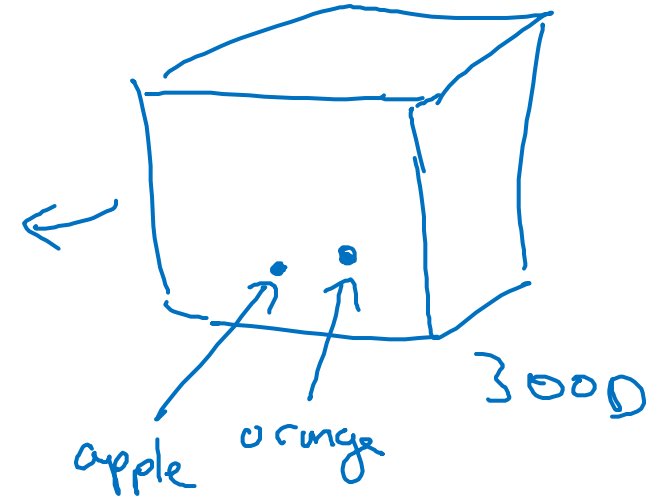
t-SNE

A non-linear
dimensionality
reduction technique

→ 300D



2D



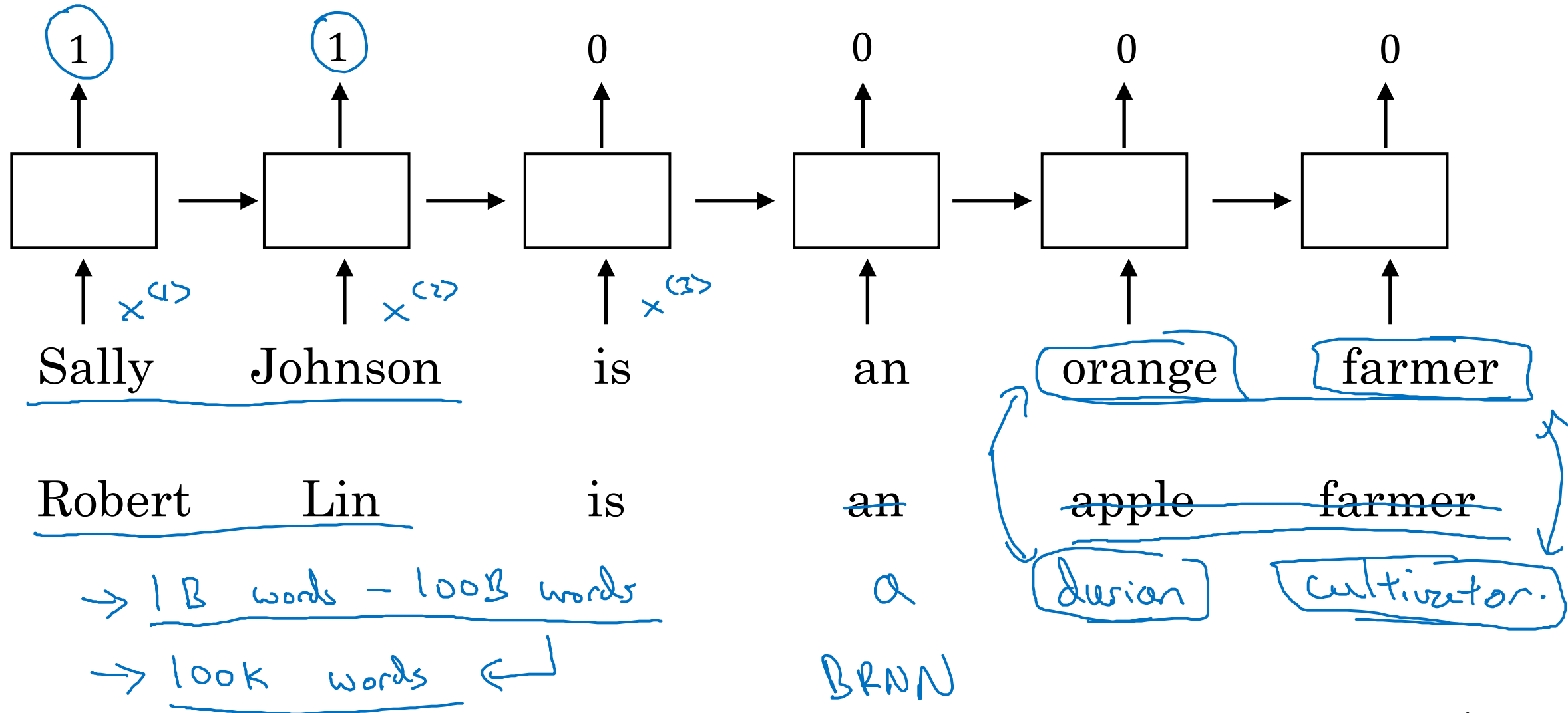


deeplearning.ai


NLP and Word Embeddings

Using word
embeddings

Named entity recognition example



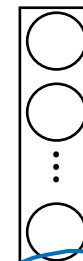
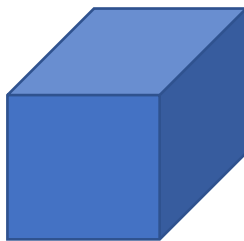
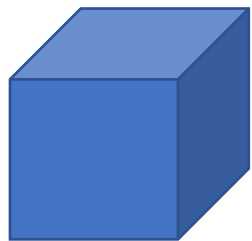
Transfer learning and word embeddings

- 
1. Learn word embeddings from 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



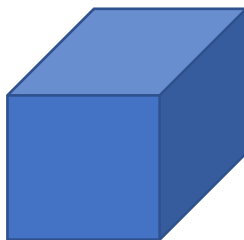
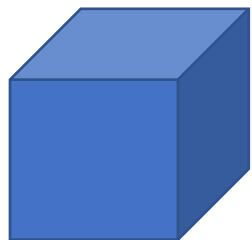
$x^{(i)}$



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



$x^{(j)}$



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



\hat{y}

$|V| = 10,000$

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



deeplearning.ai

NLP and Word Embeddings

Properties of word embeddings

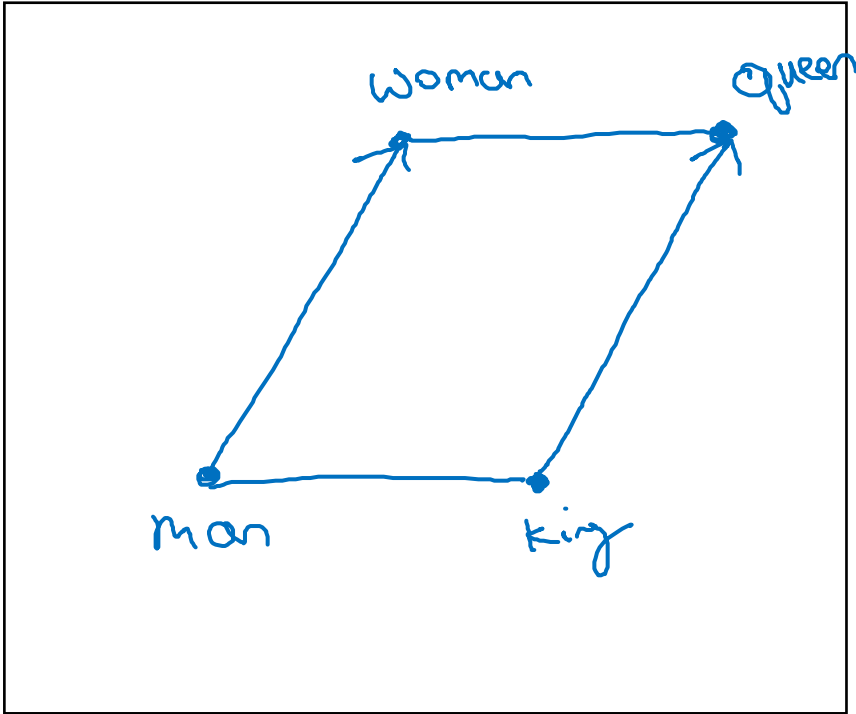
Analogy

	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

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

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

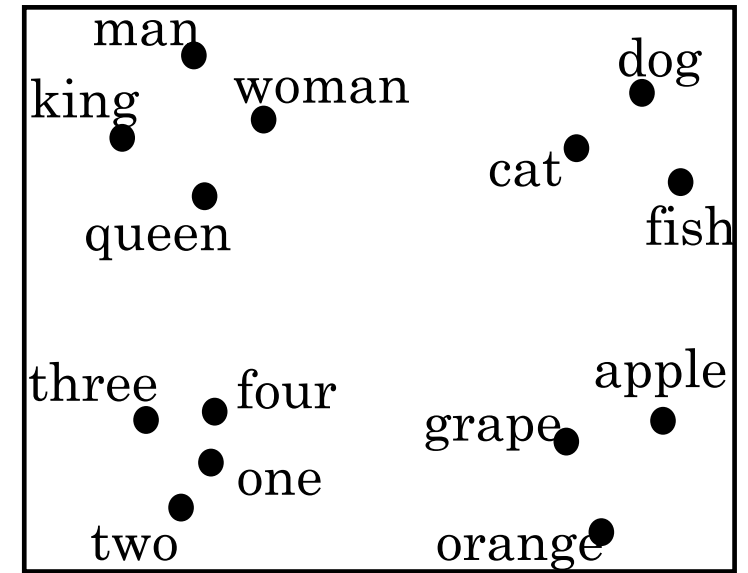
Analogies using word vectors



300 D

Find word w : $\arg \max_w$

3000 \rightarrow 20
↑



t-SNE

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

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

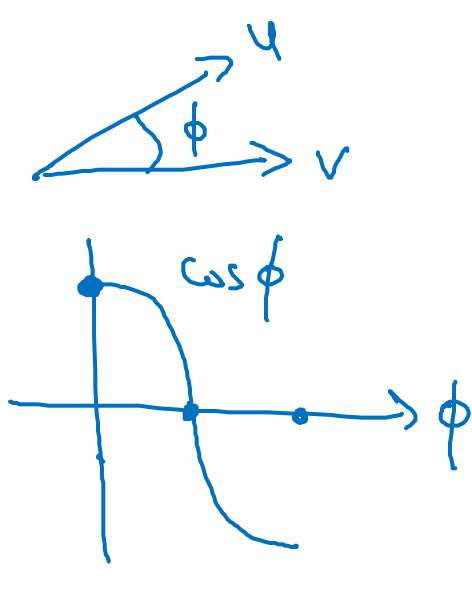
Similarity

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

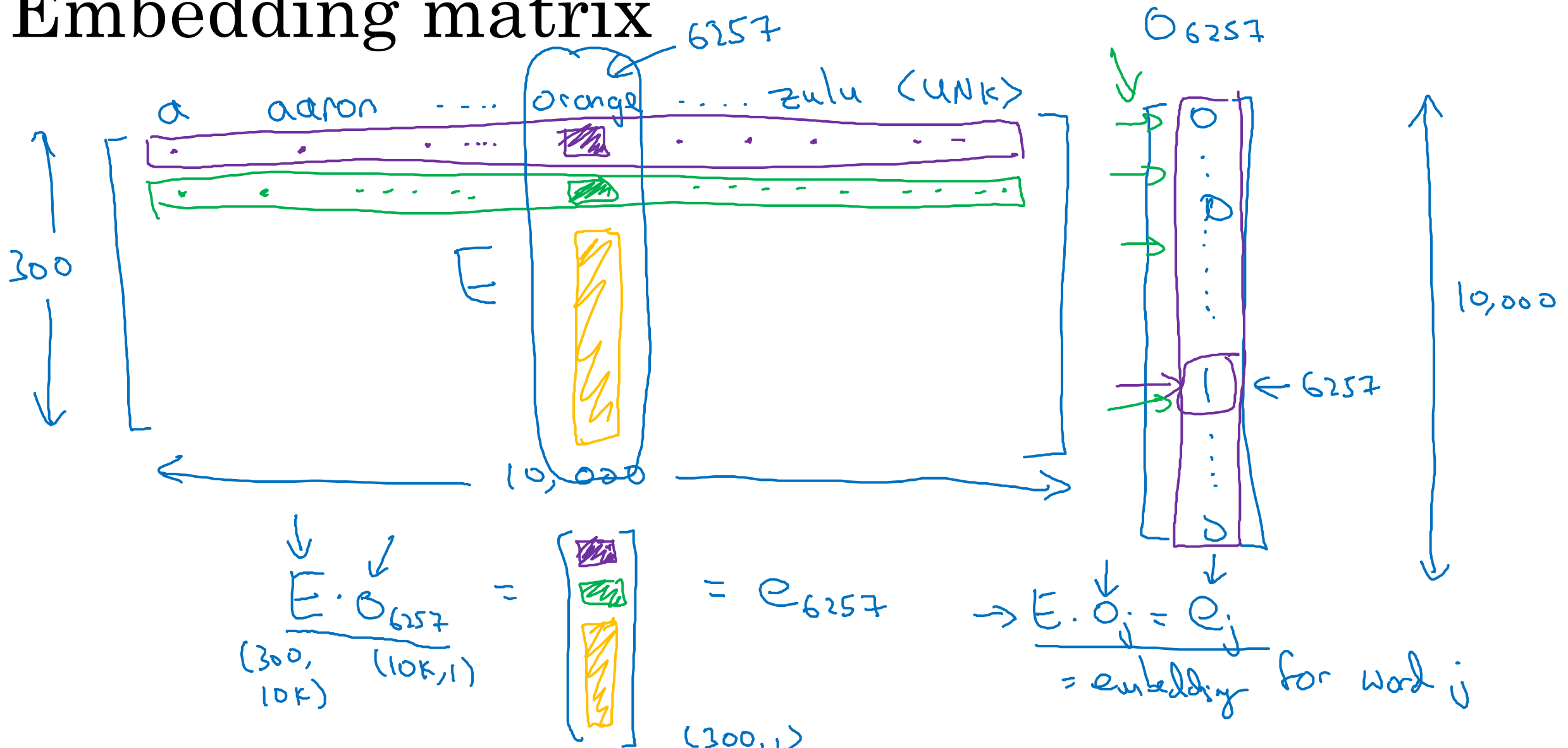


deeplearning.ai

NLP and Word Embeddings

Embedding matrix

Embedding matrix



In practice, use specialized function to look up an embedding.
 $\rightarrow \text{Embedding}$

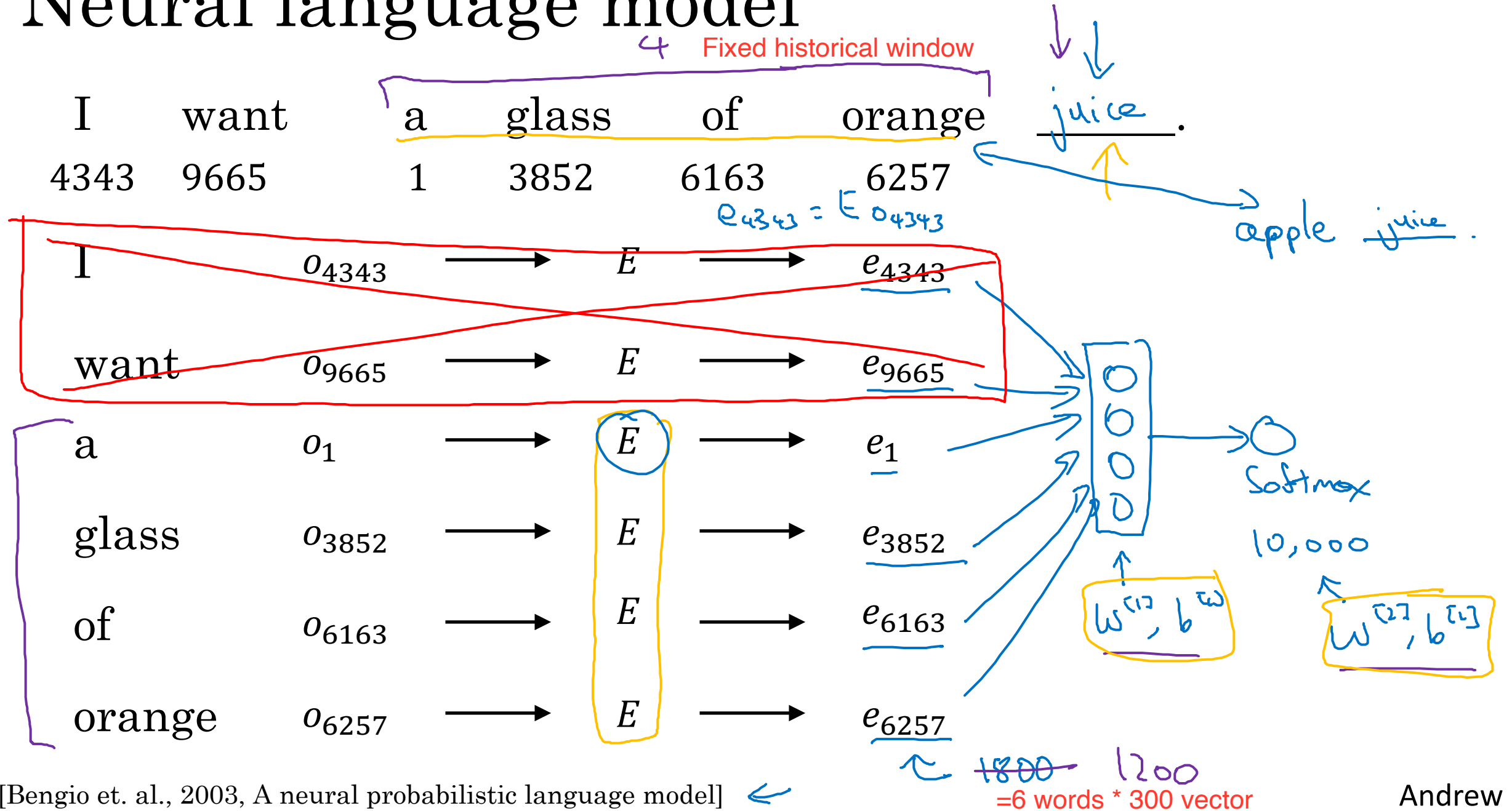


deeplearning.ai

NLP and Word Embeddings

Learning word embeddings

Neural language model



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 labeled 'context' spans the words 'a glass of orange'. A blue bracket labeled 'target' is positioned under the word 'juice'. A green arrow points from the word 'orange' to the word 'juice'.

Context: Last 4 words.

- 4 words on left & right
of the blank
- Last 1 word before the blank
- Nearby 1 word

a glass of orange ? to go along with

orange ?

glass ?

skip gram model
simpler algorithm



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



randomly choose a target word

Target

juice

glass

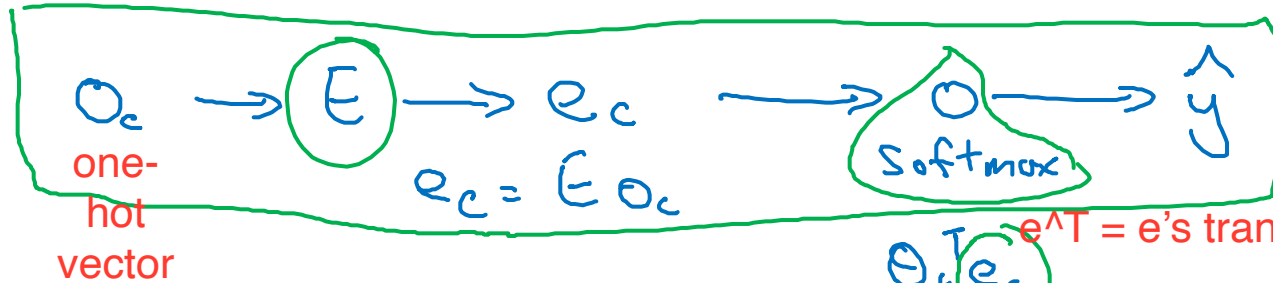
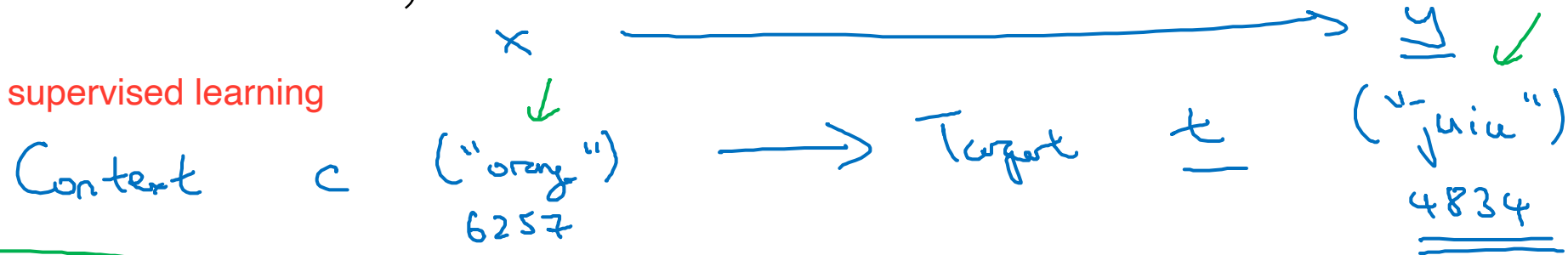
my



Model

Vocab size = 10,000k

based on supervised learning



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

where θ_t is the parameter associated with output t . Note: \mathbf{e}^T = e's transpose.

Loss function:

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

Output vector y :

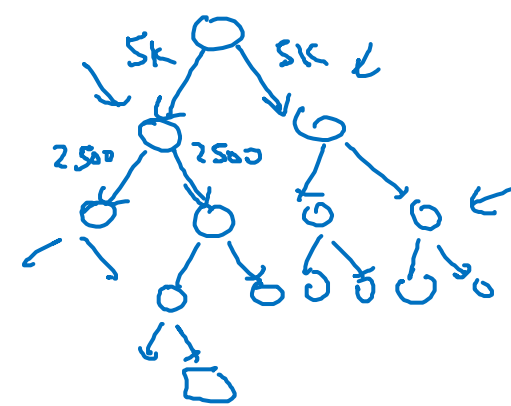

$$y = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \leftarrow 4834$$

Problems with softmax classification

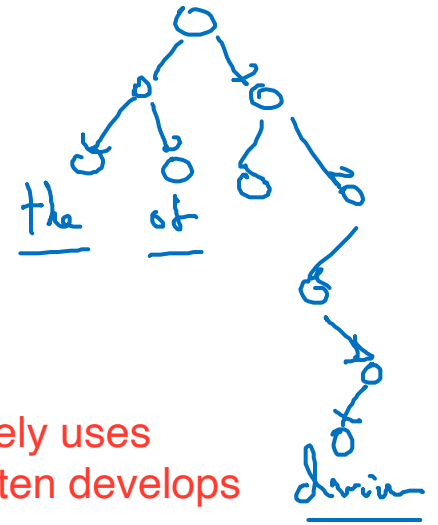
$$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$$

Hierarchical softmax to address computational cost

$\log |V|$
rather than
linear vocab size



In practice, Hierarchical softmax rarely uses balanced (symmetric) tree. Instead, it often develops trees of common words



How to sample the context c ?

→ the, of, a, and, to, ...

→ orange, apple, lemon

Q. Durian

 t
$$C \rightarrow t$$
$$P(\omega)$$



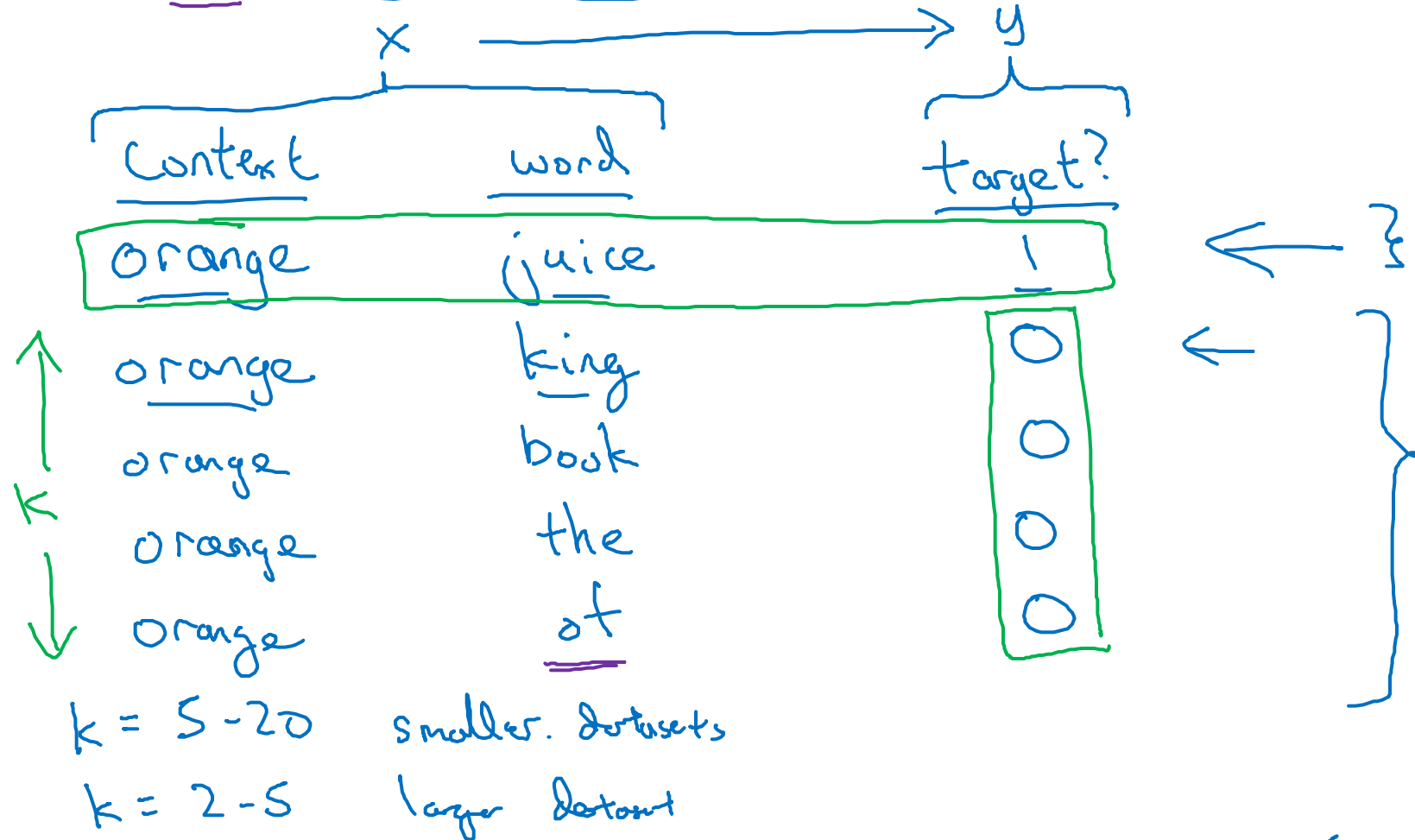
deeplearning.ai

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

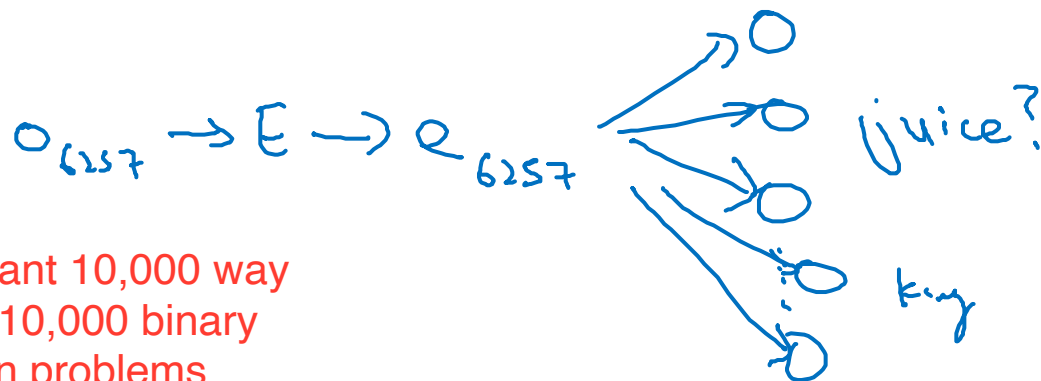
Softmax:
$$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$$

10,000-way softmax

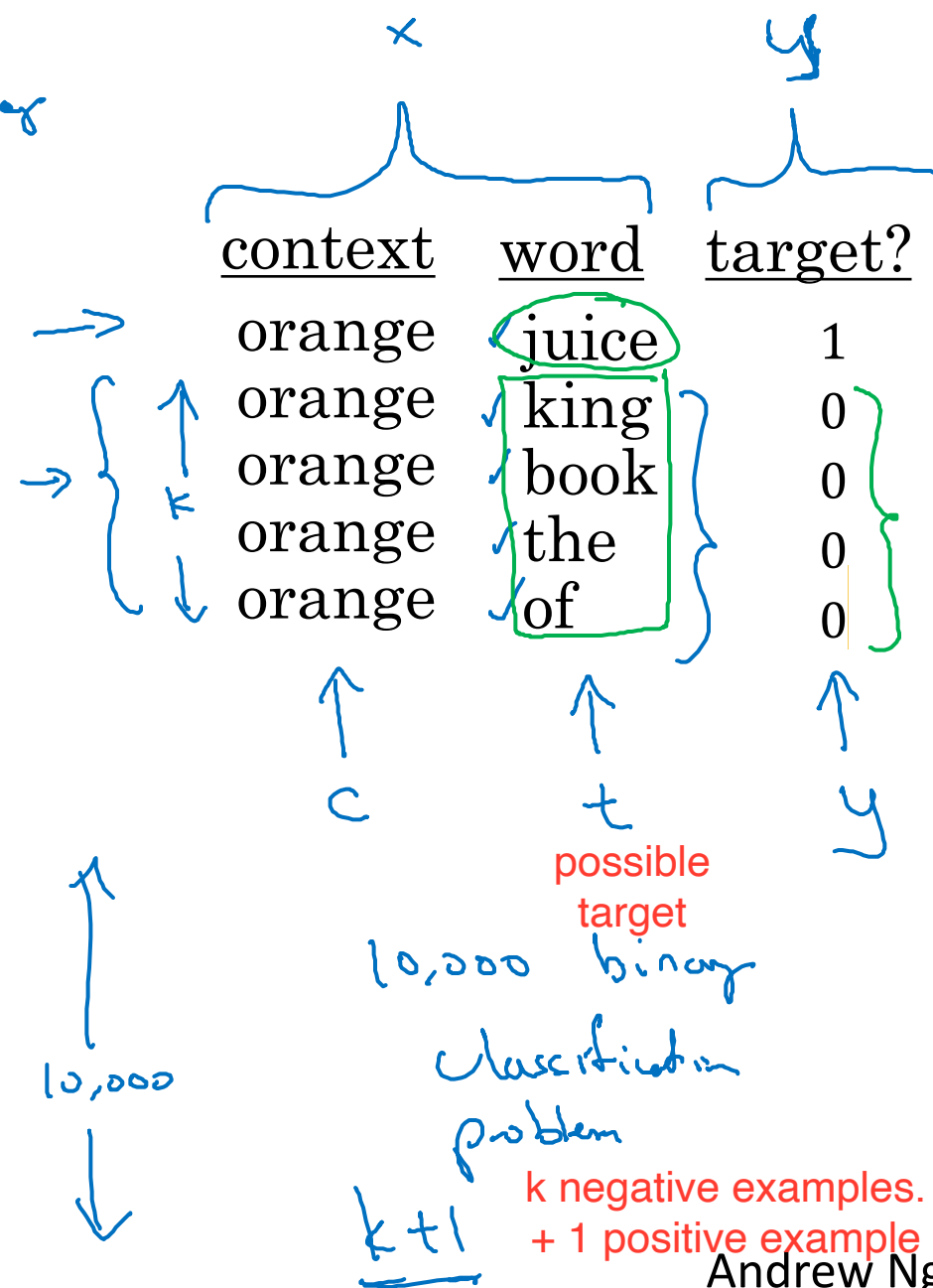
$$P(y=1 | c, t) = \sigma(\theta_t^T e_c) \leftarrow$$

basically logistic regression model

Orange
6257



instead of having a giant 10,000 way softmax, turn it into 10,000 binary logistic regression problems



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

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

the, of, and, ...

One extreme case is to sample the words in the middle (whatever was the observed distribution in the training set) but it's not representative; Another extreme, which is to take uniformly random ($1/v$), is also not representative of english words. So, the authors take the heuristic value (between two extremes of sampling from the empirical frequencies). They sampled proportional to the frequency of a word to the power of three forth. $f(w_i)$ is the observed frequency of a particular word in the english language

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



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 .



$X_{ij} = X_{ji}$ ←

X_{ij} is a count that captures how often do words i and j appear with each other, or close to each other.

Depending on the definition of context and target words, X_{ij} might be equal X_{ji} . However, if the choice of context is always the word immediately before the target word, X_{ij} and X_{ji} may not be symmetric. For the purpose of the GloVe algorithm, we can define context and target as whether or not the two words appear in close proximity. Say within plus or minus 10 words of each other.

Model

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

If x_{ij} is zero, its log is undefined and negatively infinity.. So we add an extra weighting term too sum over the terms where x_{ij} is zero.

Think of i and j as playing the role of c and t .
 $\Theta_t^T e_c$

How related are words i and j (measured by how often they occur with each other) is affected by x_{ij} .

If $x_{ij} = 0$, $0 \log 0$ is not relevant. The sum is sum only over the pairs of words that have co-occurred at least once in that context-target relationship.

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

" $0 \log 0$ " = 0 by using a convention

Stop words: some words that appear very often in the English language.

this, is, of, a, ...
derian

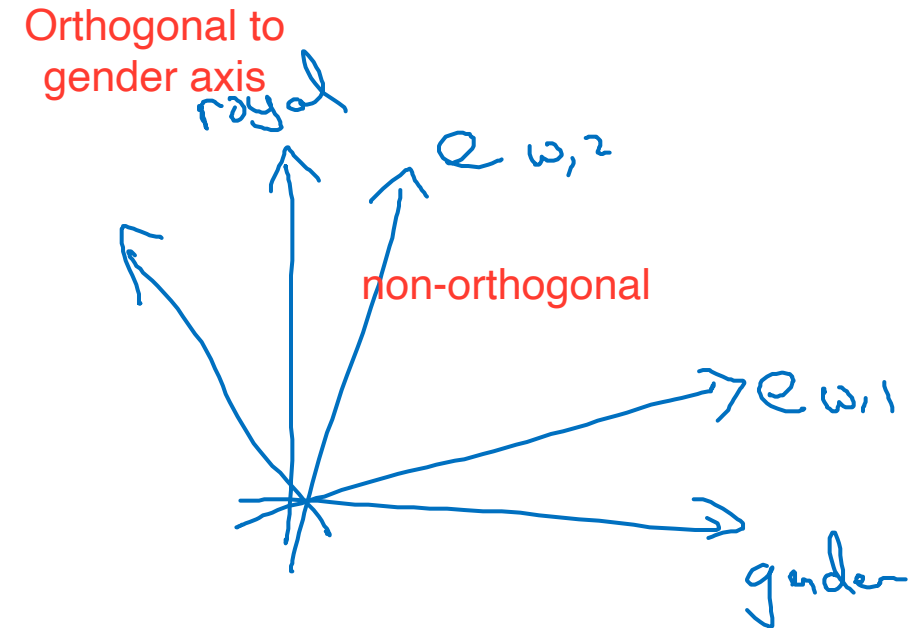
Θ_i, e_j are symmetric

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

One way to train the algorithm is to initialize θ and e both uniformly around gradient descent to minimize its objective, and then take the average when being done for every word. It is possible because θ and e play symmetric roles in this particular formulation.

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	←



We cannot guarantee that individual components of the embeddings are interpretable: the axis used to represent the features may not be well-aligned with what might be easily humanly interpretable axis.

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

If there was some invertible matrix A , then this could easily be replaced with the following:

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

Despite this type (potentially arbitrary) of linear transformation, the parallelogram map still works



deeplearning.ai

NLP and Word Embeddings

Sentiment classification

Sentiment classification is the task of looking at a piece of text and telling if someone likes or dislikes the thing they are talking about.

Sentiment classification problem

One of the challenges for sentiment classification is you might not have a huge label training set .

With word embeddings, you can build good sentiment classifiers even with only modest-size label training sets.

x

y

The dessert is excellent.



Service was quite slow.



Good for a quick meal, but nothing special.



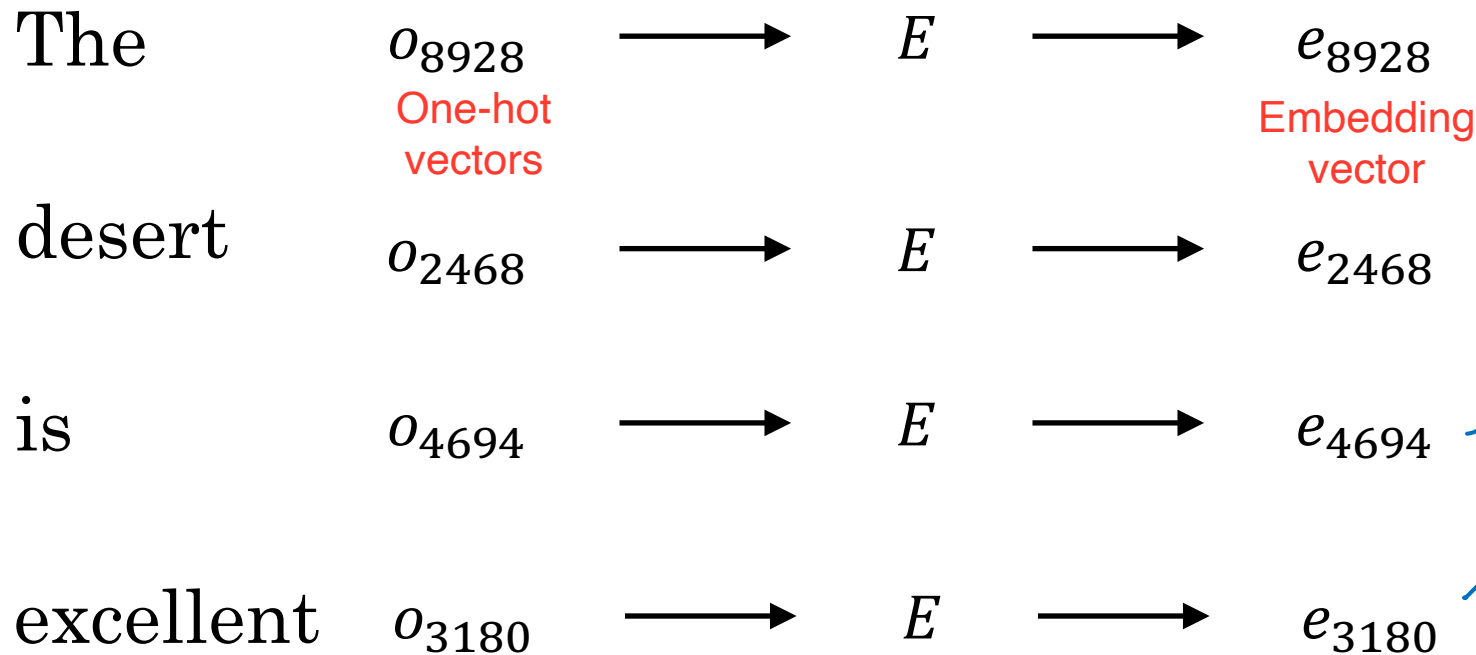
Completely lacking in good taste, good service, and good ambience.



10,000 \rightarrow 100,000 words

Simple sentiment classification model

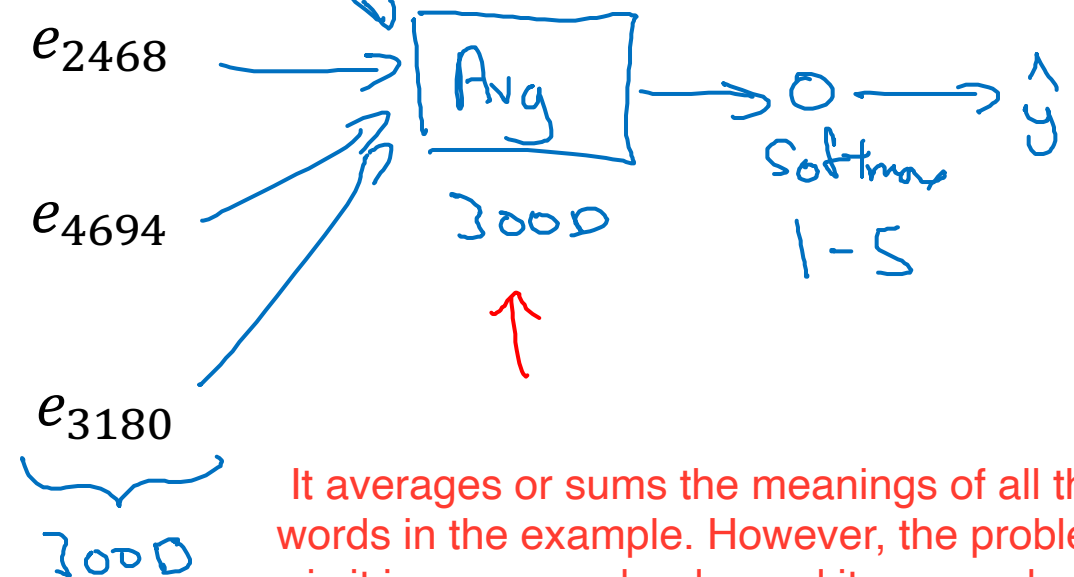
The dessert is excellent
8928 2468 4694 3180



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

↑
100 B
words

★★★★☆
Notice that by using the average operation here,
this particular algorithm works for reviews that
are short or long because you can just sum or
average all the feature vectors for all 100 words
so that it gives you a representation, a 300-
dimensional feature representation, which you
can then pass into your sentiment classifier.

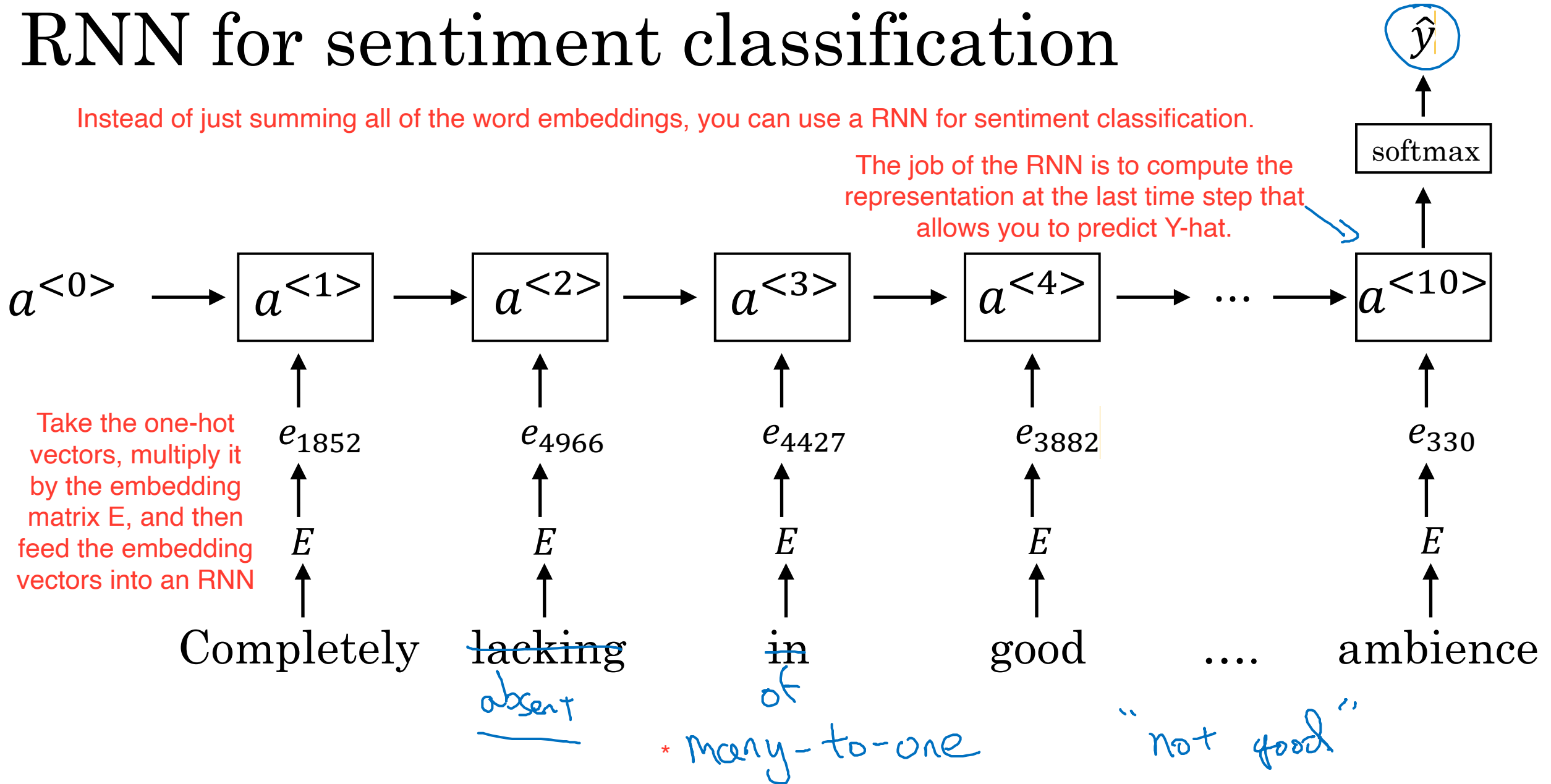


It averages or sums the meanings of all the
words in the example. However, the problem
is it ignores word order and it may end up
having a lot of the representation of good.

RNN for sentiment classification

Instead of just summing all of the word embeddings, you can use a RNN for sentiment classification.

The job of the RNN is to compute the representation at the last time step that allows you to predict \hat{Y} .





deeplearning.ai

NLP and Word Embeddings

Machine learning and AI algorithms are increasingly trusted to help with, or to make important decisions.

Debiasing word embeddings

The problem of bias in word embeddings

Bias here is not bias variants but gender bias or ethnicity bias.

Man:Woman as King:Queen

Unhealthy gender stereotype

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.

Bias relating to socioeconomic status



Addressing bias in word embeddings

