

Word Representation

Given a Vocabulary V , we can use one-hot representation.

Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ \vdots \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ \vdots \\ 0 \\ 0 \end{bmatrix}$
\downarrow	\downarrow	\downarrow	\downarrow	\downarrow	\downarrow
e_{5391}	e_{9853}	e_{4914}	e_{7157}	e_{456}	e_{6257}

Vocabulary:

a	1
aaron	2
and	367
harry	4075
potter	6830
zulu	
<UNK>	10,000

$$\text{len(Vocabulary)} = 10^5$$

The problem with one-hot representation is that it cannot represent the "distance" or "similarity" between words. Applying dot product or subtraction won't work.

Eg:

I want a glass of orange juice.

I want a glass of apple .

If the answer to the first question is juice, then it's highly likely the second answer should also be juice because orange and apple are similar words.

Thus, a better representation is **Featurized Representation: Word Embedding**

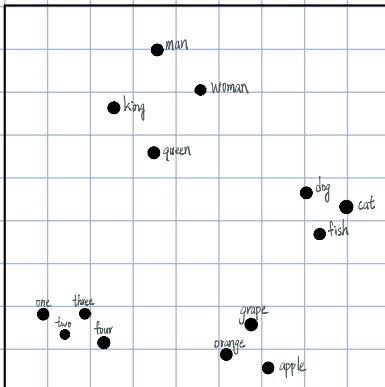
Feature	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.0	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.0
Age	0.03	0.02	0.7	0.69	0.03	-0.02
Food	0.09	0.01	0.02	0.01	0.95	0.97
Eg: 300	⋮					
		e_{5391}	e_{9853}		e_{456}	e_{6257}

With word embeddings, Apple and Orange are close.

Visualize Word Embedding

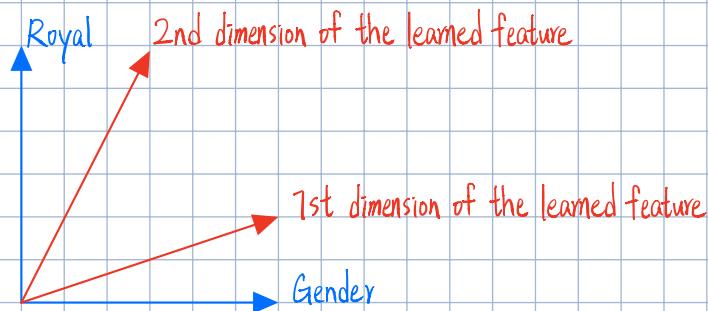
Read "Visualizing Data using t-SNE."

The word embedding vectors are 300 dimensional vectors. We can use t-SNE to project 300D into 2D plane.



The above example represents words by using meaningful features (Gender, Royal, Age, etc). However, features in word embeddings are automatically learned (learning word embedding will be discussed later). Thus, we **can't** really interpret meaningful information from automatically learned features.

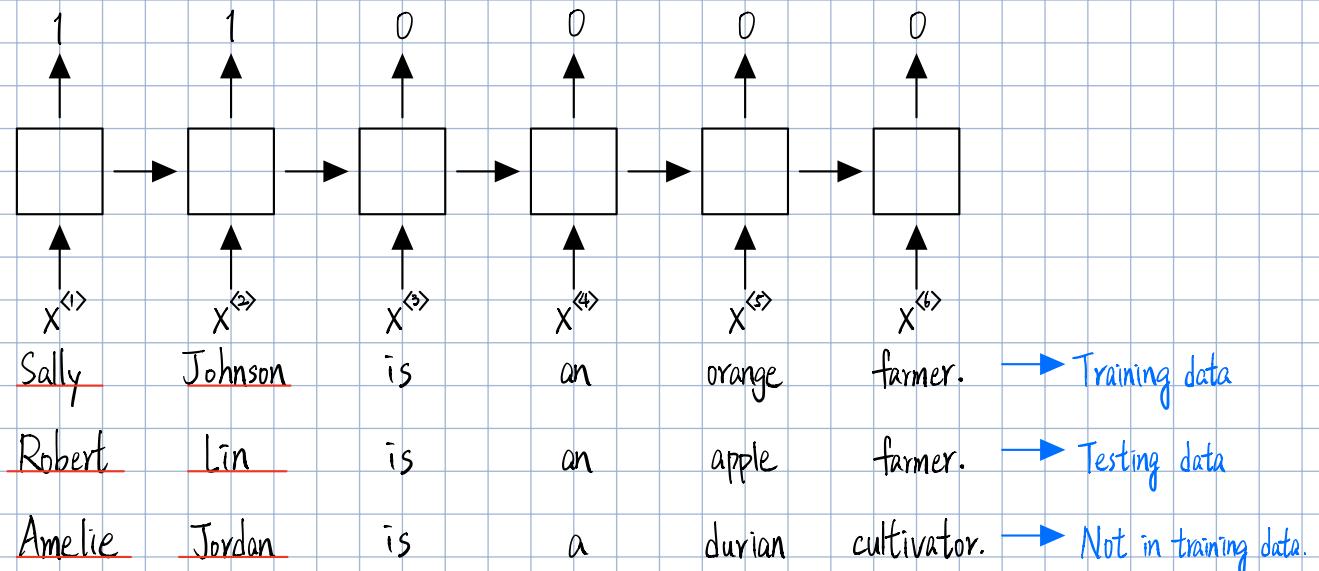
A note on the featurization view of word embeddings:



The learned features are not orthogonal.

Using Word Embeddings

Named Entity Recognition example



durian and cultivator are not seen in the labelled training set of Named Entity Recognition tasks, but it can still figure out Amelie and Jordan are names due to the words, durian and cultivator, are close to apple and farmer in the sense of word embeddings.

Word embeddings are learned from huge unlabelled corpus (1 billion to 100 billions) so durian and apple, and orange are grouped in high dimensional feature space.

Don't you see. It's **Transfer Learning**. Thus, I can do the following:

1. Learn word embeddings from large text corpus. (1 billion to 100 billions words).
Or download pre-trained word embeddings online.
2. Transfer word embeddings to new tasks (eg: Named Entity Recognition) with **smaller labelled training data.**
(eg: 100 k words)
3. Optional: Continue to finetune word embeddings with new data.

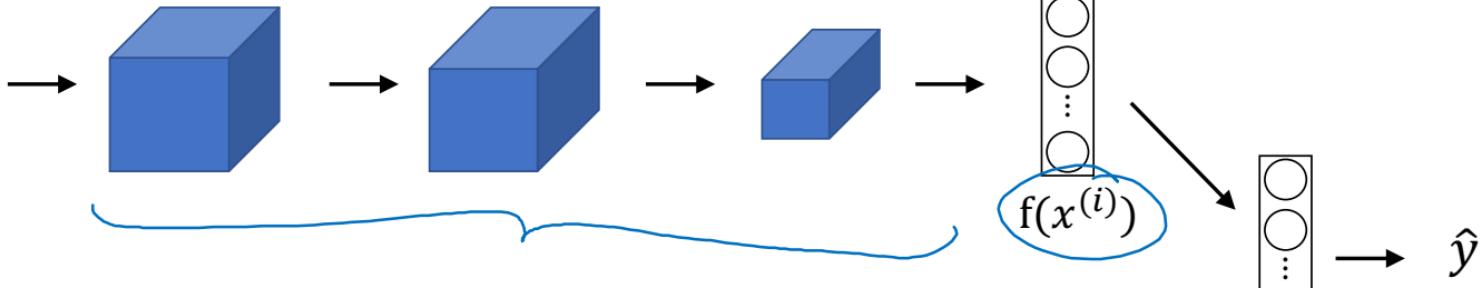
We can use dense 300D feature vectors, instead of sparse 10⁵D one-hot vectors.

Relation to face encoding (embedding) 128D

Unlimited faces

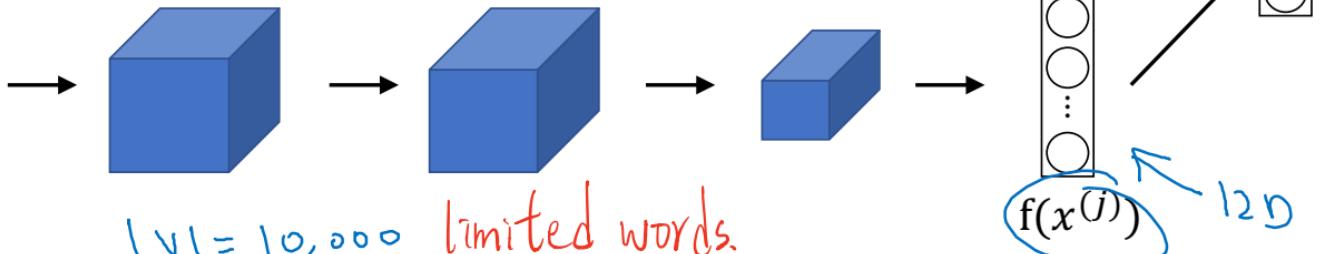


$x^{(i)}$



$x^{(j)}$

$|V|=10,000$ limited words
 $e_1, \dots, e_{10,000}$



Properties of Word Embeddings

Question: man is to woman as king is to ?

Analogies

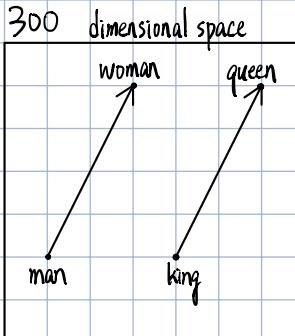
	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.70	0.69	0.03	-0.02
Food	0.09	0.01	0.02	0.01	0.95	0.97
	e_{5391}	e_{woman}	e_{king}	e_{queen}		

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

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

The solution is: first calculate $e_{man} - e_{woman}$, then find $e_?$ to make

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



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

Use similar vector to find the word w.

Find the word W : $\underset{w}{\operatorname{argmax}} \text{similarity}(e_w, e_{man} - e_{woman} + e_{king})$

Cosine Similarity

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

L2 norm can also be used.

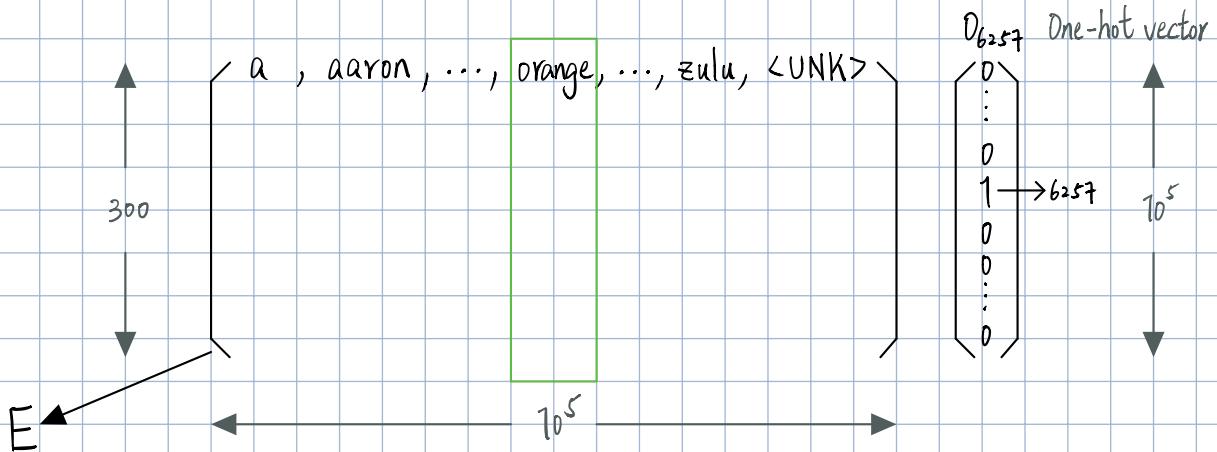
$$\text{sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2}$$

When $\text{sim}(u, v)$ is larger, then u and v are more similar.

Read 2013 "Linguistic regularities in continuous space word representation"

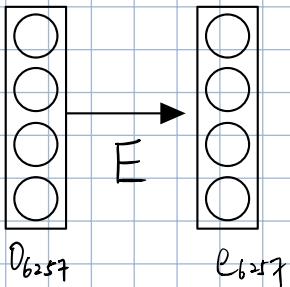
Embedding Matrix E

Vocabulary = a , aaron, ..., orange, ..., zulu, <UNK>. length = 10^5



$$E \cdot D_{6257} = \begin{bmatrix} \\ \\ \\ \end{bmatrix} \quad (300, 1)$$

Embedded matrix E is just like a fully connected layer. Recall $a = wx + b$

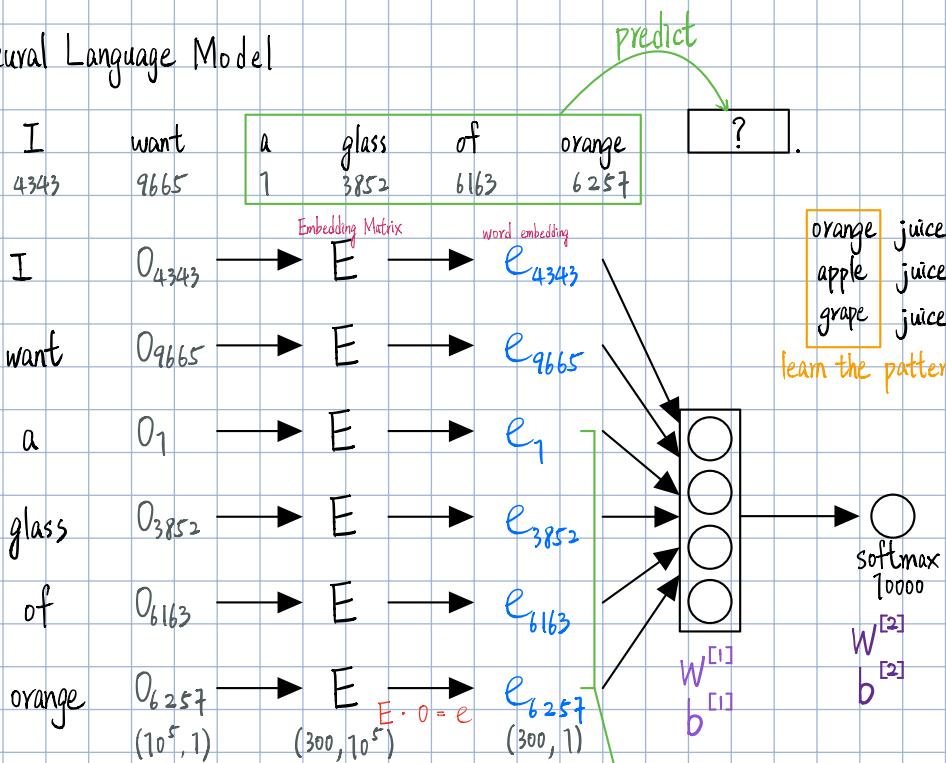


Learning Word Embeddings

Read 2003 "A Neural Probabilistic Language Model."

We've beening using word embeddings to solve various tasks, but we haven't talked about how to create word embeddings at the first place. Now, let's talk about it.

Neural Language Model



E is learned.

Input layer: In order to have a fixed size of an input layer, we can choose to use the fixed number of previous words. Eg, use 4. The input layer would be $(300 \times 4, 1)$.

Other context / target pairs.

Another example:

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

Context: Last 4 words.

4 words on left and right.

Last 1 word.

Nearby 1 word.

Vocabulary:

a	1
aaron	2
•	
•	
and	367
•	
•	
harry	4075
•	
•	
potter	6830
•	
•	
zulu	10,000

Word2Vec

Read 2013 "Efficient Estimation of Word Representations in Vector Space."

Skip-grams

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

Context orange orange orange	Target juice glass go	}	randomly chosen from a fixed size window
---------------------------------------	--------------------------------	---	------------------------------------------

Model:

Vocabulary size = 10000

Context c ("orange") $\xrightarrow{6257}$ Target t ("juice") $\xrightarrow{4834}$

$O_c \rightarrow E \rightarrow e_c \xrightarrow{\theta_t^T e_c} \text{softmax} \rightarrow \hat{y}$

$$\text{softmax: } p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10^5} e^{\theta_j^T e_c}}$$

target
context

θ_t : parameter associated with target output t

θ_t and $e_c \in \mathbb{R}^{300}$

$$L(\hat{y}, y) = \sum_{i=1}^{10^5} y_i \log \hat{y}_i$$

$$y = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \rightarrow 4834 \quad 10^5$$

Problems with softmax classification

$$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10^5} e^{\theta_j^T e_c}}$$

takes too much time to compute. Use Hierarchical Softmax

How to sample the context C?

Note: It's OK if we do poorly on this artificial prediction task. The more important by-product of this task is that we learn word embeddings.

Negative Sampling

Previously, skip-gram model constructs a supervised learning task, mapping from context to target and learning word embeddings. But the downside was that the softmax was slow to compute.

Read "Distributed Representation of Words and Phrases and their Compositionalities"

Define a new learning problem.

Training sentence:

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

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

↑ ↑ ↑

c t y

One positive sample from the training sentence.

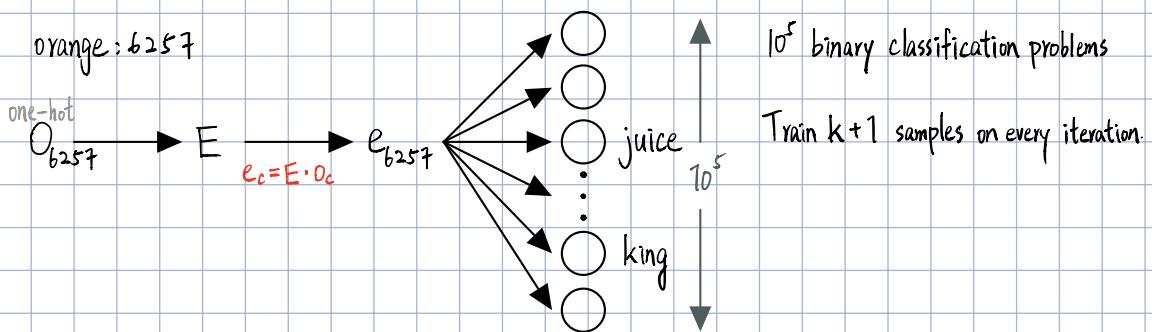
k negative samples sampling from Vocabulary.

$k = 5 \sim 20$ smaller dataset
 $2 \sim 5$ large dataset

Vocabulary:	
a	1
aaron	2
and	367
harry	4075
potter	6830
zulu	10,000

Model:

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



How to select negative samples?

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

Sampling averagely will get a lot of "the", "of", "and", etc, which is bad.

$$P(freq_i) = \frac{P(freq_i)^{\frac{2}{d}}}{\sum_{j=1}^{10^5} P(freq_j)^{\frac{2}{d}}}$$

GloVe Word Vectors

Read 2014 "GloVe: Global Vectors for Word Representation."

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

x_{ij} : the number of times that word j appears in the context of word i .

Model:

$$\text{minimize} \sum_{i=1}^{10^5} \sum_{j=1}^{10^5} f(x_{ij}) (\theta_i^T e_j + b_i + b_j' - \log x_{ij})^2$$

weighting term

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

$$\theta_i^T e_j + b_i + b_j' - \log x_{ij}$$

We want to know how related are words i and j as measured by how often they occur with each other, which is affected by x_{ij} .

$$\theta_t^T \text{ and } e_j \text{ are symmetric. } e_{\text{word}}^{(\text{final})} = \frac{e_{\text{word}} + \theta_{\text{word}}}{2}$$

θ_t and e_c should be initialized randomly at the beginning of training.



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.



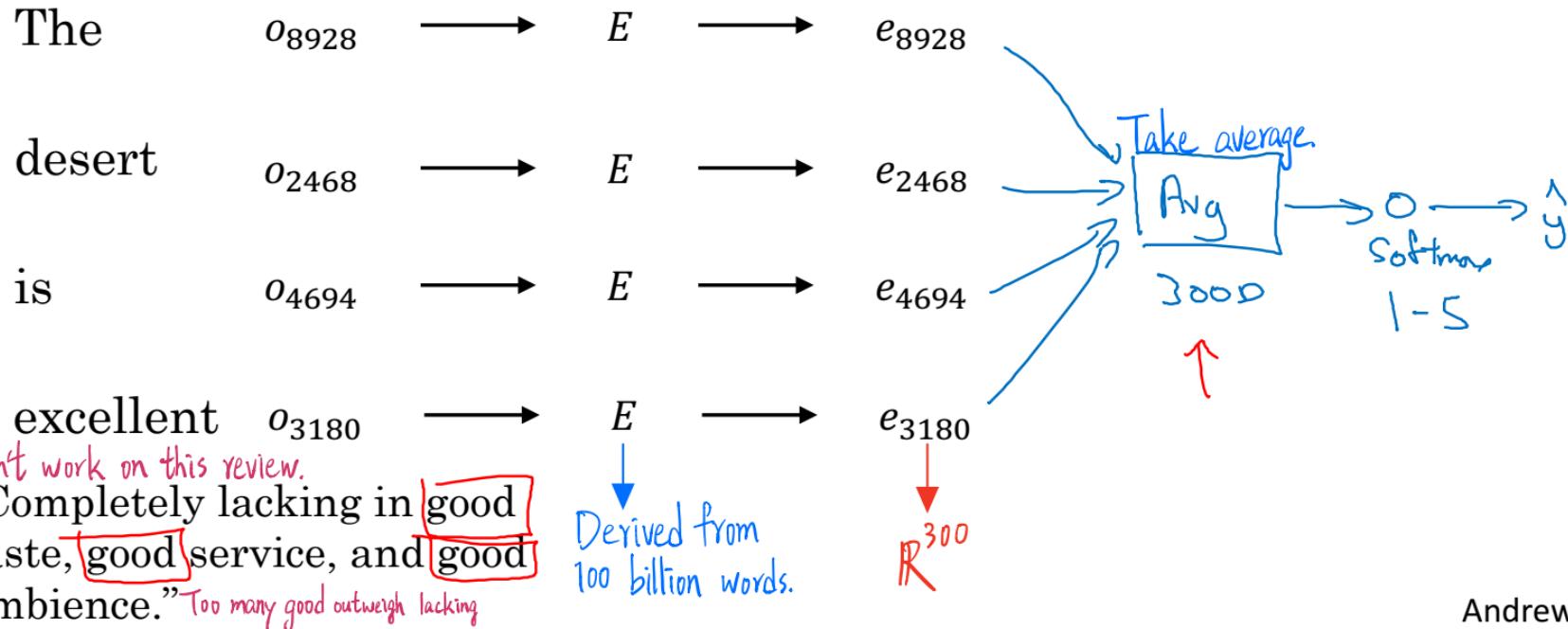
Completely lacking in good taste, good service, and good ambience.



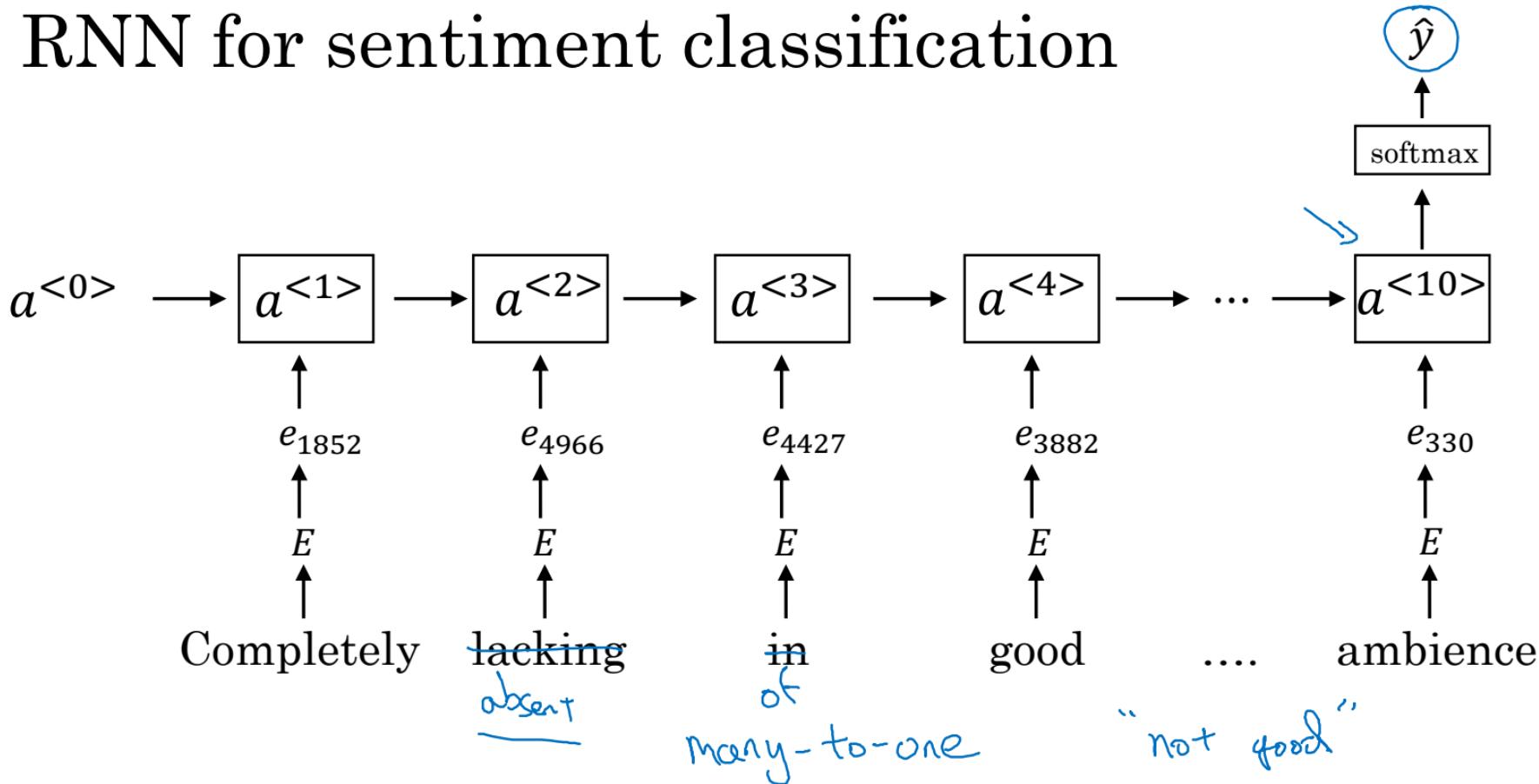
10,000 \rightarrow 100,000 words

Simple sentiment classification model

The dessert is excellent
8928 2468 4694 3180



RNN for sentiment classification





deeplearning.ai

NLP and Word Embeddings

Debiasing word embeddings

The problem of bias in word embeddings

Man:Woman as King:Queen

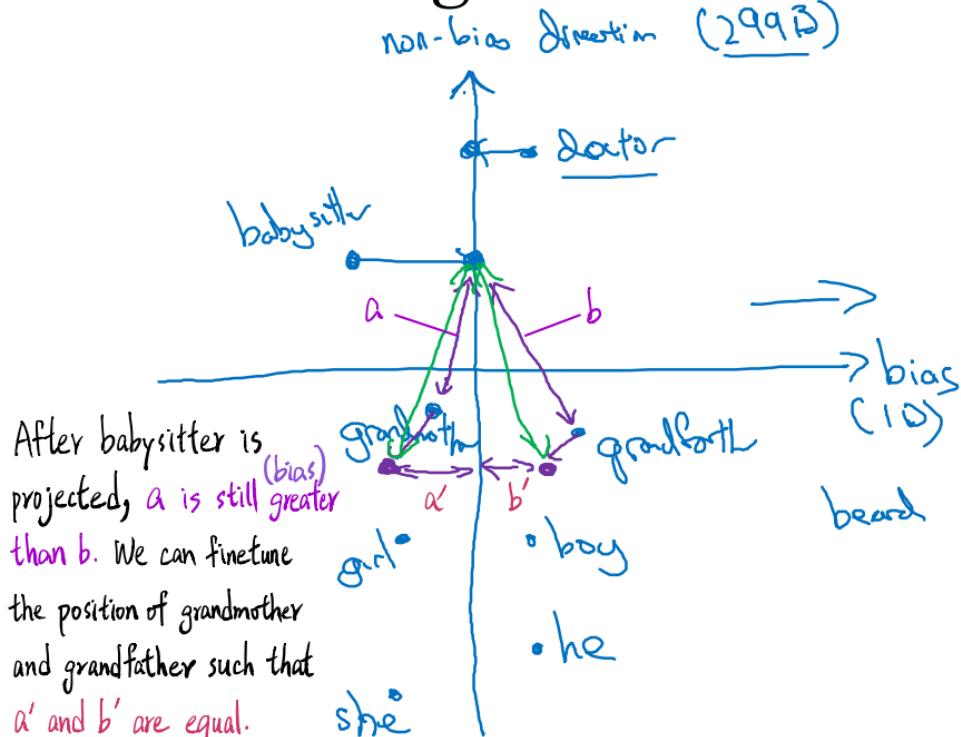
Man:Computer_Programmer as Woman:Homemaker 

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.

$$\begin{cases} \text{he} - \text{she} \\ \text{male} - \text{female} \\ \vdots \\ \text{average} \end{cases}$$

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

3. Equalize pairs. Make the distance equal so there's no bias.
 $\rightarrow \text{grandmother} - \text{grandfather}$