

Week 2

Word embeddings:

Trying to find relationships between words like man:woman::king:queen

word representation:

$V = [a, aaron, \dots, zulu, \langle \text{UNK} \rangle]$ $|V| = 10,000$

1-hot representation

Man Woman Apple Orange
(5391) 09853 0456 6257

$\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \end{bmatrix} \leftarrow (5391)$

I want a glass of ~~orange juice~~
orange juice

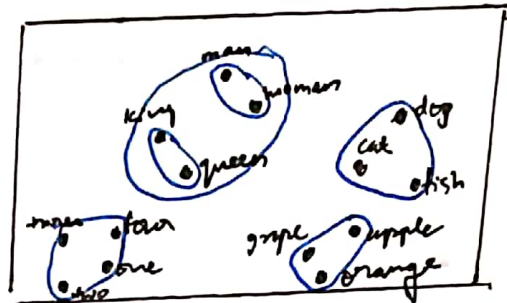
05391 I want a glass of apple ?

It can't decide apple juice even though it has learned to recognise orange juice since it doesn't know apple and orange are related so we use featurewise representation:

	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	0.09	0.01	0.02	0.01	0.95	0.97
...						
Size						
Cost						
Verb						
etc						

Since apple and orange are similar, it'll now be able to find out I want a glass of apple juice

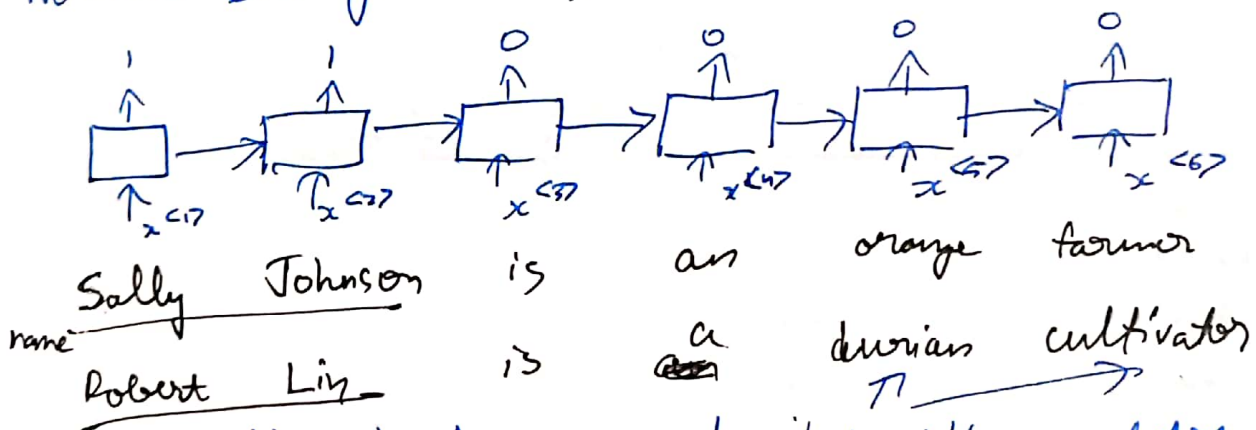
You can use the t-SNE algorithm to map the 300 dimensions into 2D.



It is called embedding since each word is similarly embedded into a 300D space.

Using word embeddings:

Named entity recognition example:



How to learn words it hasn't seen before and associate it with orange farmer?

- Use ~~1B~~ 1B - 100B words (unlabelled) from the internet along with your 100K labelled words

Transfer learning and word embeddings:

1. Learn word embeddings from large text corpus (1-100B onwards)
(or download pre-trained embedding online)
2. Transfer embedding to new task with smaller training set (say, 100k onwards)

Now you can use a smaller 10k or 300 word vocabulary (the training set).

3. Optional: Continue to finetune the word embeddings with new data.

Properties of word embeddings:

They can help with analogy reasoning.

Man \rightarrow Woman as King \rightarrow ?

	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

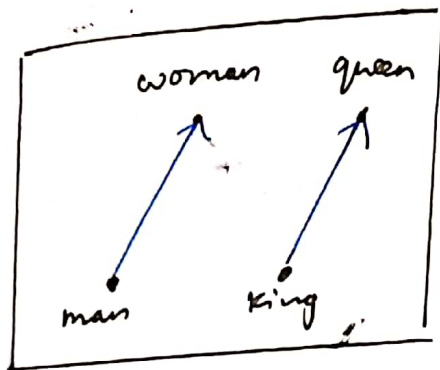
$$e_{\text{man}} - e_{\text{woman}} \approx \begin{bmatrix} 2 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad e_{\text{king}} - e_{\text{queen}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

\therefore Gender is the relation

\therefore It's queen

Analogies using word vectors:

So the algorithm we use to find the word w (queen in this example) is as follows:



Find a word w such that $e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - e_w$

$$\therefore \arg \max_w$$

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

Find the word with the maximum

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

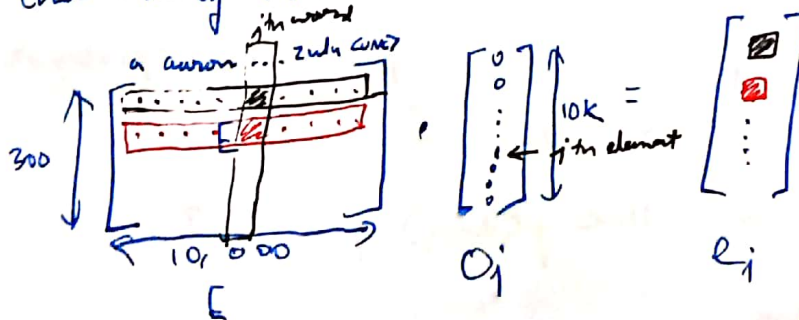
sim is called cosine similarity

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

It finds the cosine of the angle ϕ between the 2 vectors.

Embedding matrix:

embedding for word i , $e_i = E \cdot O_i$



In practice we use a specialized function (embedding() in keras) to look up an embedding since matrix multiplication is costly since most are 0s (zeros)

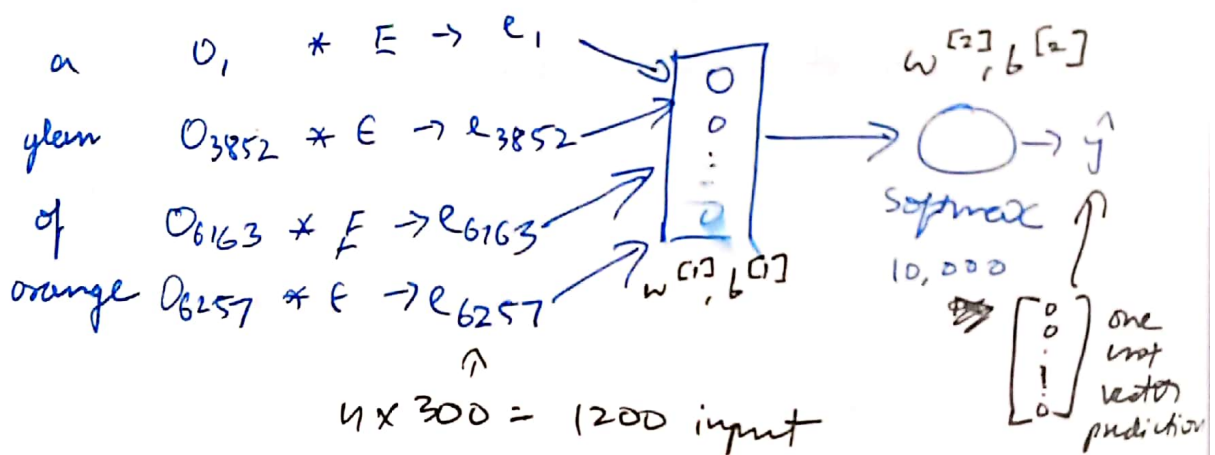
Learning word embeddings:

It started off with complex algorithms, but over time simpler algorithms have been found that does the job better. We will start learning from the complex algorithms to develop a good intuition.

I) A normal neural network:

I want a glass of orange —
4343 9665 1 3852 6163 6257

We take the last 4 words into consideration.



Here we used a 'last 4 words' context, other contexts that work better are:

- 4 words on left & right a glass of orange? — to go along with
- Last 1 word orange?
- Nearby 1 word (like glass) glass? — skip grams

↑ context word ↑ Target word

Word2Vec:

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

- Choose a random context word
- Then choose a target word randomly in a ~~4 or 5~~ word range (coming before)

Context	Target
orange	juice
orange	glass
orange	my

model:

Vocab size = 10,000

Context c ("orange") \rightarrow Target t ("juice")
 6257 4834

$$O_c \rightarrow E \rightarrow e_c \xrightarrow{\text{Softmax}} O \rightarrow \hat{y}$$

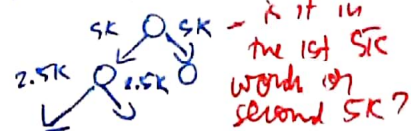
Softmax: $p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{i=1}^{10,000} e^{\theta_i^T e_c}}$ θ_t = parameter associated with output t

$$\rightarrow L(\hat{y}, y) = - \sum_{i=1}^{10,000} y_i \log \hat{y}_i \quad y = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \leftarrow 4834$$

Problem with softmax classification:

- It gets computationally expensive to sum up over 10,000 elements everytime. So we use hierarchical softmax

Common words are on top while rare words are deeper \rightarrow



However ~~also~~ negative sampling is a better method to make softmax faster.

How to sample context c ?

→ the, of, a... ∈ common words

→ orange, apple... ∈ less frequent

we try to train on these less frequent words

Negative Sampling: 2 words

Given a ~~target and context~~ pair we will predict if they are a target and context pair.

I want a glass of orange juice to go along with my cereal
Context c word t target? y

orange	juice	1
orange	king	0
orange	the	0
orange	of	0
orange	...	0
orange	...	0
orange	...	0

← even though it comes beside orange

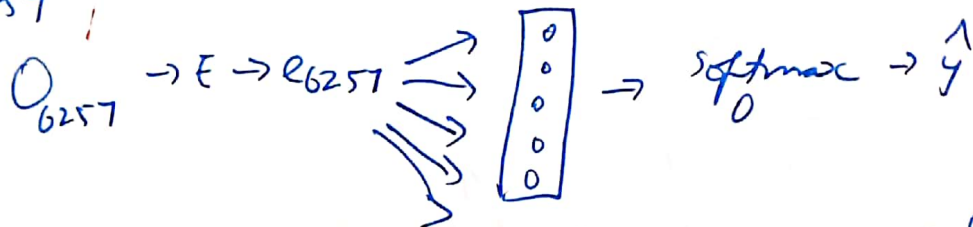
← pick k random words from the dictionary with orange and put 0

$k = 5-20$ smaller datasets
 $k = 2-5$ larger datasets

Model:

context	word	target	
orange	joke	1	- 1 correct example
orange	king	0	} 4 randomly picked words to make incorrect examples
orange	book	0	
orange	the	0	
orange	of	0	

Orange
6257



\therefore Rather than training all 10,000, we only train 5 (1 + 4)

How to choose the negative examples:

\rightarrow Choose words that have a high probability of occurring.

Problem: words like the, of, and will come

\rightarrow The other extreme is to choose rare words but that isn't good either

\rightarrow The best is to choose in between given by

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{i=1}^{10000} f(w_i)^{3/4}}$$

Global Word Vectors:

Global \rightarrow global vectors for word representation

Here we use

X_{ij} = # times j appears in context of i
 $\swarrow \quad \uparrow$ $\quad \uparrow$ $\quad \uparrow$
 $c \quad t$ \quad like t \quad like c

model:

$$\text{minimize } \sum_{i=1}^{10000} \sum_{j=1}^{10000} f(x_{ij}) (\underbrace{\Theta_i^T e_j}_{\substack{\uparrow \\ \text{weighting} \\ \text{term}}} + b_i + b_j - \log x_{ij})^2$$

\uparrow
we train this

\rightarrow Here we use $f(x_{ij})$ to prevent computation if $\log x_{ij} = 0$

i.e. $f(x_{ij}) = 0$ if $x_{ij} = 0 \rightarrow 0 \log 0 = 0$
like it's anyways gonna be 0, so skip this training example

We can even use it to give less priority to words like this, is, of, a ... and more priority to words like durians so they kind of balance out

$\rightarrow \Theta_i, e_j$ are symmetrical so at the end we can take their average

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

Applications using Word Embeddings:

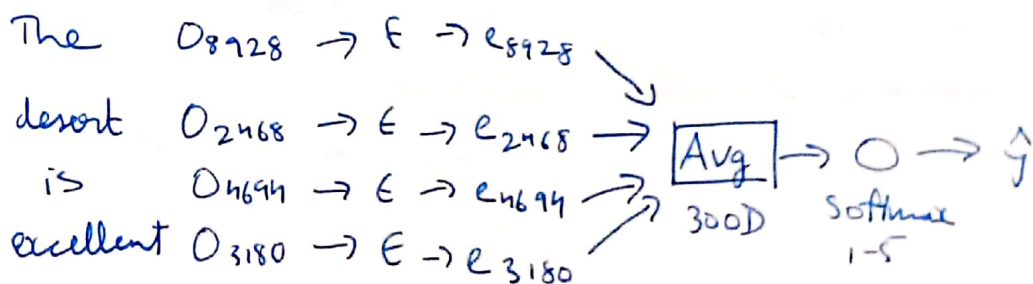
Sentiment Classification:

The task of looking at a piece of text and telling if someone likes or dislikes the thing they're talking about.

→ Challenge is that we may not have a huge labelled training set for it but with word embedding we can build a good classifier.

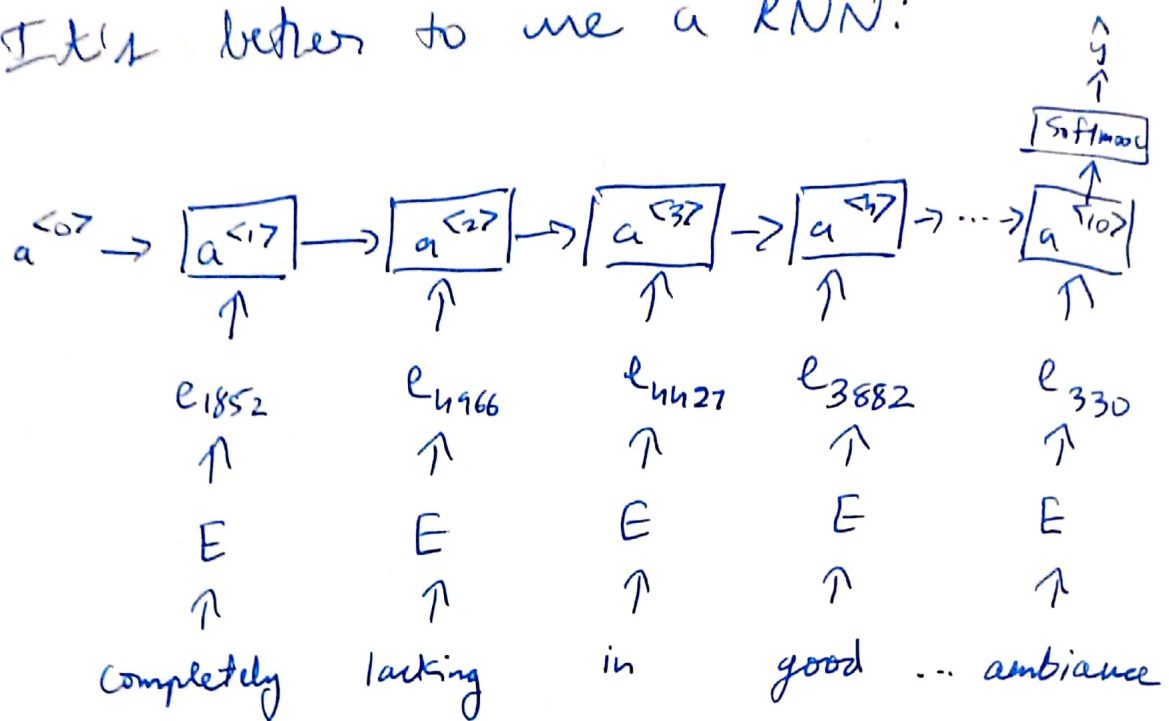
x	y
1) The desert is excellent	★ ★ ★ ★ ★
2) Service was quite slow	★ ★ ☆ ☆ ☆
3) Good for a quick meal, but nothing special	★ ★ ★ ★ ☆ ☆
4) Completely lacking in good taste, service & ambience	★ ☆ ☆ ☆ ☆

Simple model: (for example)



For a sentence like "completely lacking in good taste, good service and good ambience", it'll think it's a positive message if we use the above model.

It's better to use a RNN:



many-to-one model

Even if a word wasn't in the training example (like we use absent instead of lacking), it'll still perform good.

Debiasing word embeddings:

A few researchers found that the AI learnt biased stuff like:

→ Man: computer programmer as Woman: Homemaker

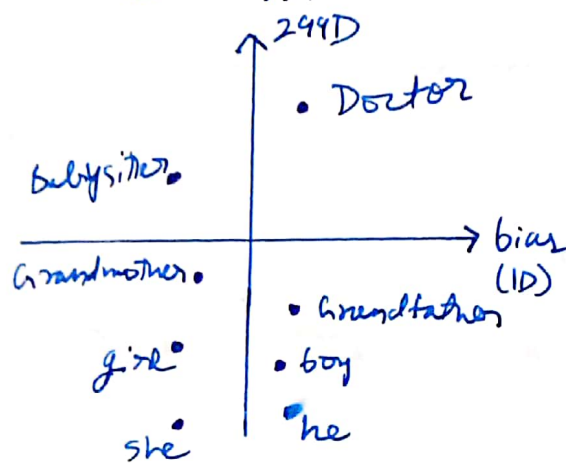
→ Father: Doctor as Mother: Nurse

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

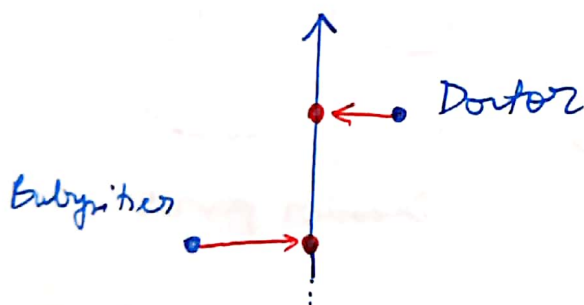
Steps to address it (gender as example):

1. Identify bias direction

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



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



3. Equalize points

