

Week 2

Video 1: Word Representation

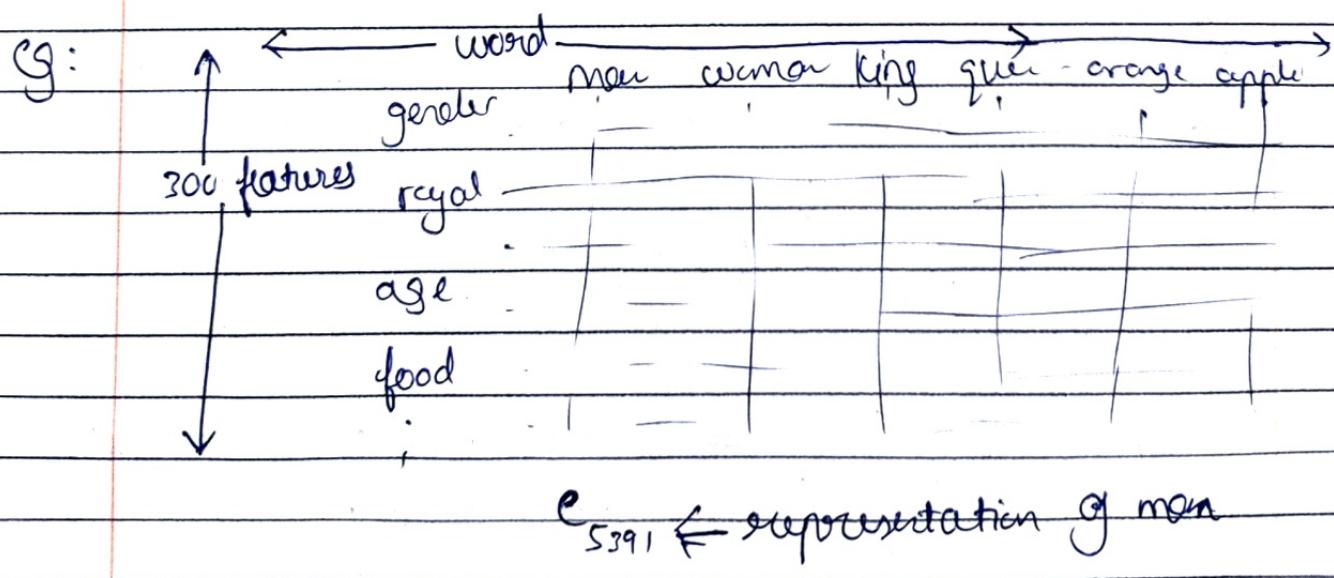
- * Word embedding is a way of representing words
- * We define language using vocabulary and represent our words with a one hot vector
- * Weakness - words are treated as objects and can't be generalised

Example:

"I want a glass of orange" → predict if juice found
 "I want a glass of apple" → can't (0.0)
 predict juice

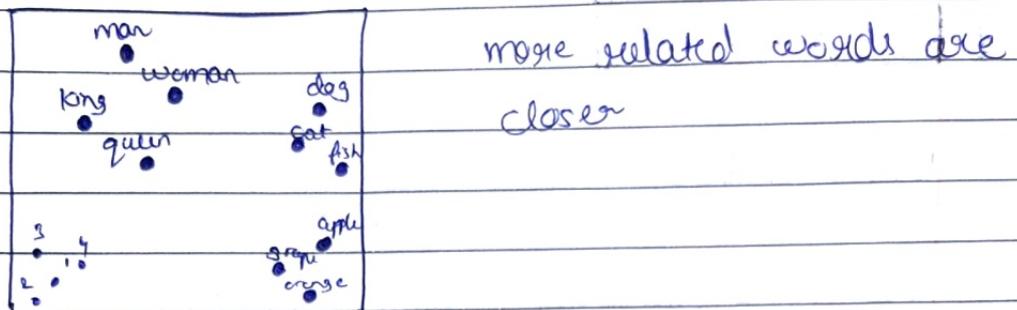
- * product of one hot encoding vectors is 0.

- * Feature representation (word embeddings)

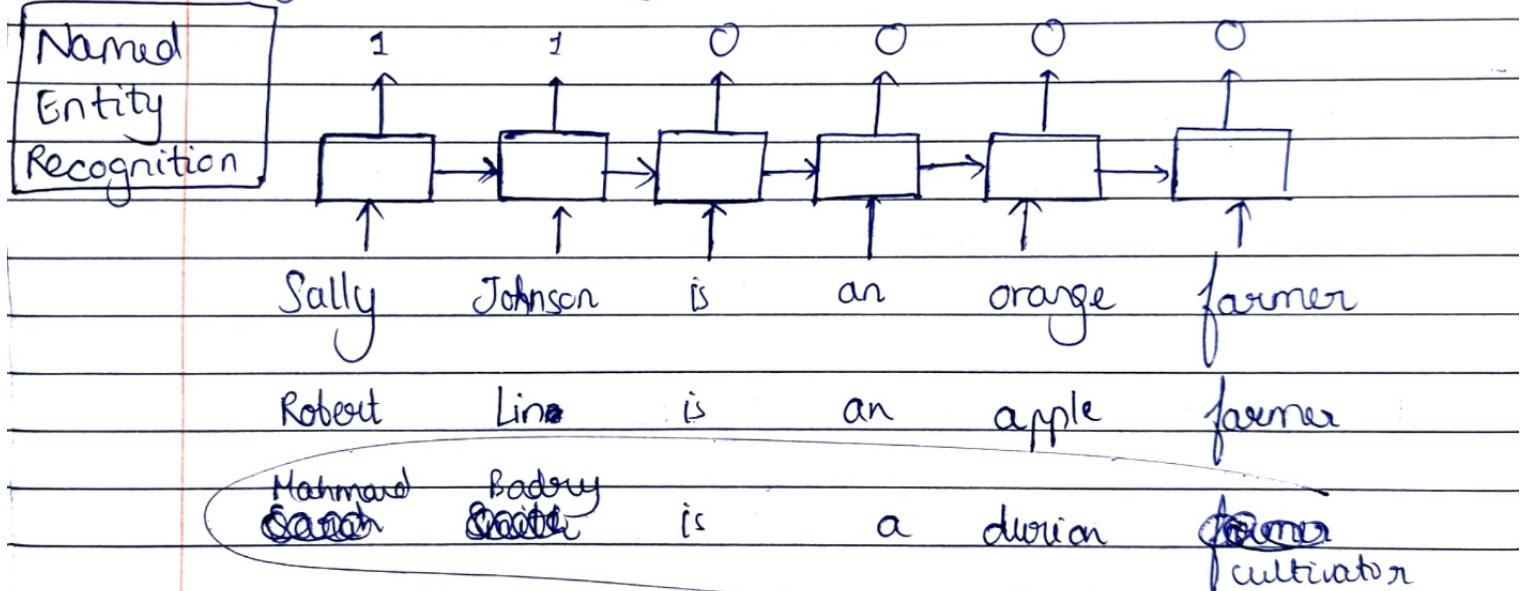


orange & apple share similar, hence algo can generalise between them.

* To visualise word embeddings we use a t-SNE algorithm to reduce features to 2 dimensions



Video 2: Using Word Embeddings



If the algorithm is trained on a smaller set of words it might not recognise durian as a fruit. But even if it does,

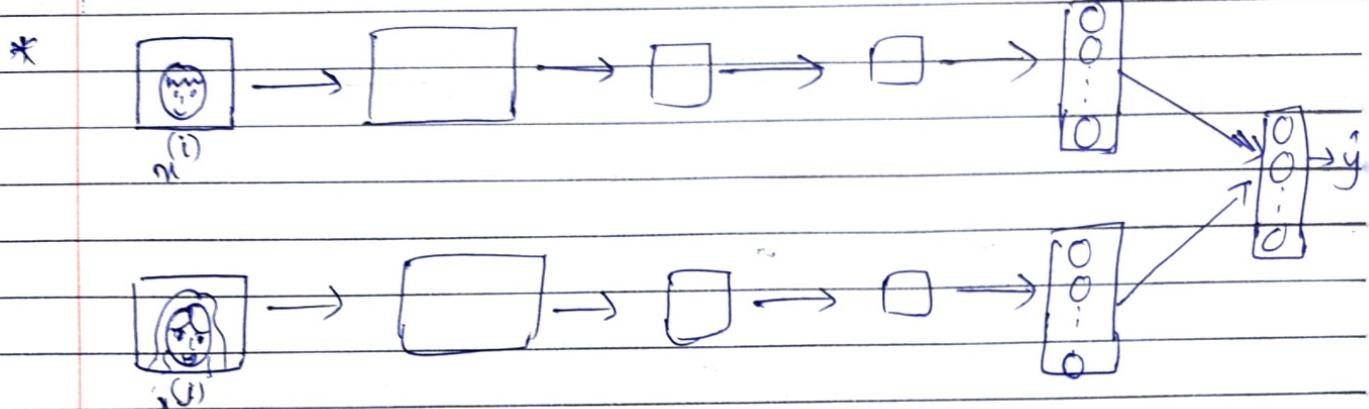
durian: fruit: orange
cultivator: farmer

Word embeddings are learned from a large text corpus probably 100 billion words.

These are transferred to a task which requires smaller

training set.

You must fine-tune the word embeddings with new data.



face \rightarrow vector \rightarrow check for similarity

Video 3: Properties of Word Embeddings

* Analogy Reasoning

	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

e_{5391}

e_{man}

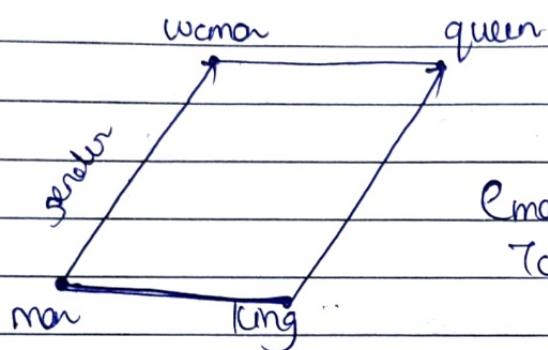
e_{woman}

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

$$e_{\text{kng}} - e_{\text{queen}} = \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

~~e_kng - ?~~

man : woman :: king : ?



$$e_{\text{man}} - e_{\text{woman}} = e_{\text{kng}} - e_{\text{queen}}$$

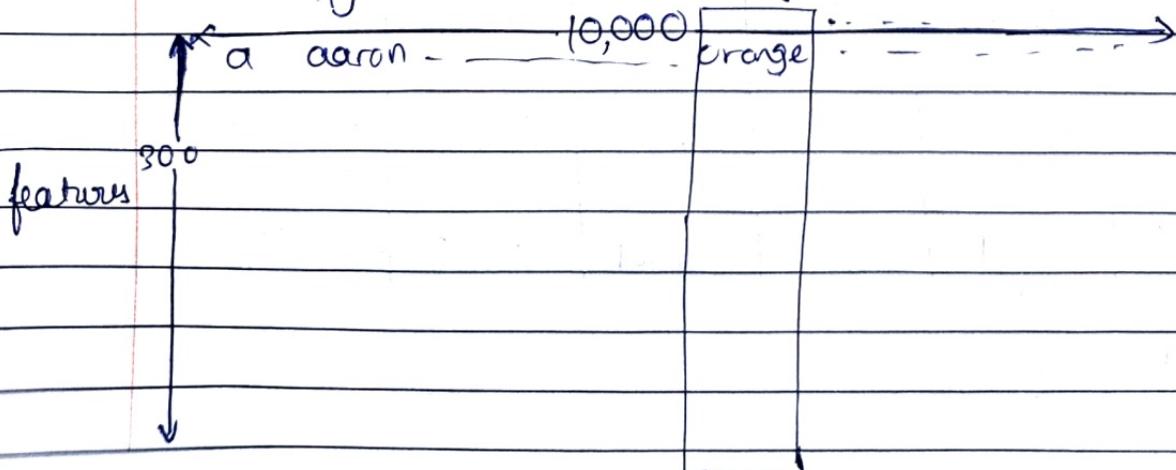
To find word w:

$$\arg \max_w \sin(e_w, e_{\text{kng}} - e_{\text{man}} + e_{\text{woman}})$$

↑
≈ 75% accurate

$$\sin(u, v) = \frac{u^T v}{\|u\| \|v\|} = \frac{\|u \cdot v\|^2}{\|u\| \|v\|}$$

Video 4: Embedding Matrix



E - Embedding matrix

O_{6257} - One hot vector

$$E \cdot O_{6257} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} = e_{6257}$$

(300, 1)

Learning Word Embeddings

Video 4:

Neural Language Model

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

$$I \quad O_{4348} \rightarrow E \rightarrow e_{4348}$$

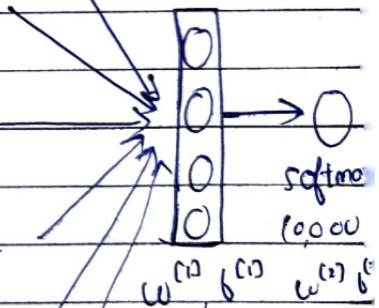
$$\text{want} \quad O_{9665} \rightarrow E \rightarrow e_{9665}$$

$$\text{a} \quad O_{\dots} \rightarrow E \rightarrow e_{\dots}$$

$$\text{glass} \quad O_{3852} \rightarrow E \rightarrow e_{3852}$$

$$\text{of} \quad O_{6163} \rightarrow E \rightarrow e_{6163}$$

$$\text{orange} \quad O_{6257} \rightarrow E \rightarrow e_{6257}$$



Context:

(i) last 4 words

4 → hyperparameters

(ii) 4 words on L & R

(iii) last one word

(iv) nearby 1 word

i.e skip grams model

Videos: Word to Vec

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

context	target
orange	juice
" "	glass
" "	my

Given a context word, you are asked to predict a randomly chosen word within $+/- 5$ / 10 window of input context word.

Model:

Vocab size = 10000

context c ("orange") $\xrightarrow{G_{257}}$ target t ("juice") $\xrightarrow{4834}$

$O_c \rightarrow E \rightarrow e_c \rightarrow O \rightarrow \hat{y}$
softmax

$$e_c = E_{0c}$$

softmax: $p(f|c) = \frac{e^{\theta_t \cdot j_{cc}}}{\sum_{j=1}^{10000} e^{\theta_t \cdot j_{cc}}}$

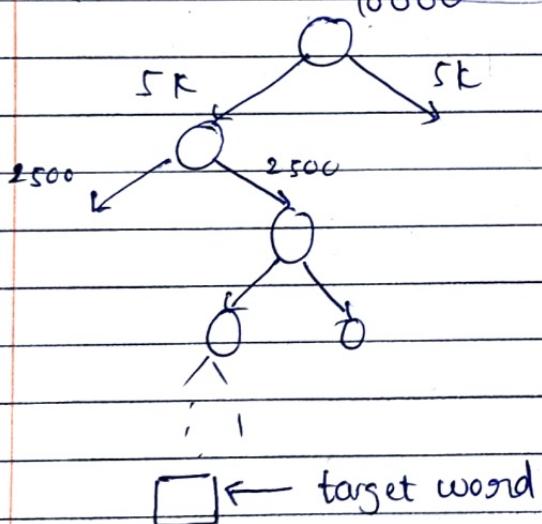
θ_t = parameter associated with output

$$\mathcal{L}(\hat{y}, y) = \sum_{i=1}^{10000} y_i \log \hat{y}_i$$

\hat{y} = possibility/probability of target output by softmax

$$y = \begin{bmatrix} ? \\ i \\ ? \\ ? \end{bmatrix} \quad \begin{bmatrix} 4.8357 \\ 0.001 \end{bmatrix}$$

Hierarchical softmax: uses log



Where does word lie?
first 5k or latter 5k

If first 5k:
then in first 2500/latter 2500

so on till it reaches the leaf of the tree-word

common | uncommon

How to sample context c?

- choose content randomly from corpus
- frequent words can dominate over uncommon words
- we don't do this normally.

Video 6: Negative Sampling

context (c) word (w) target (y)

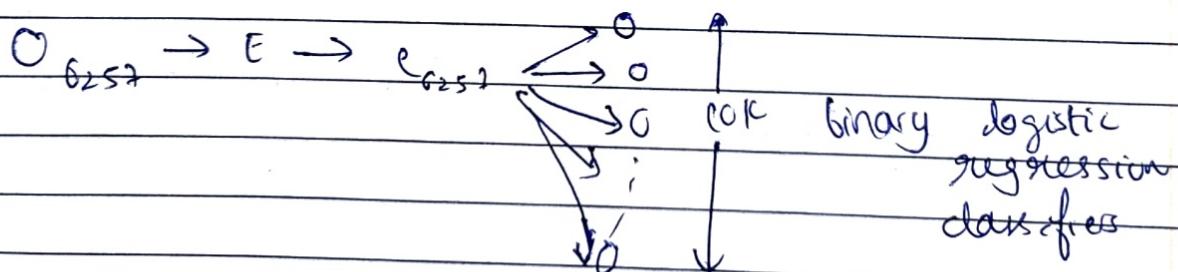
orange	juice	1
orange	king	0
orange	book	0
orange	the	0
Orange	of	0

define a logistic regression model

The chance of $y=1$, given input c, t pair

$$P(y=1 | c, t) = \alpha(\mathbf{g}_t^\top \mathbf{e}_c)$$

\uparrow
parameter \nwarrow embedding



Video 7: GloVe Word Vectors

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

c, t :

$$x_{ij} = \# \text{ times } j \text{ appears in context of } i$$

↑
 r_i
 c
 t
 ↑
 c

$x_{ij} = x_{ji}$ - can't that measure how often i & j appear close to each other.

GloVe Model:

$$\text{minimise} : (\theta_i^T e_j - \log x_{ij})^2$$

$$\text{minimise} : \sum_{i=1}^{10000} \sum_{j=1}^{10000} f(x_{ij})(\theta_i^T e_j + b_i + b_j - \log x_{ij})^2$$

↑
weighing term

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

$$\log 0 = 0$$

Weighing factor gives meaningful amount of computation even to the less frequent words like duration & gives more weight to words like this, a, a, of which appear too many times

θ_i & e_j are symmetric

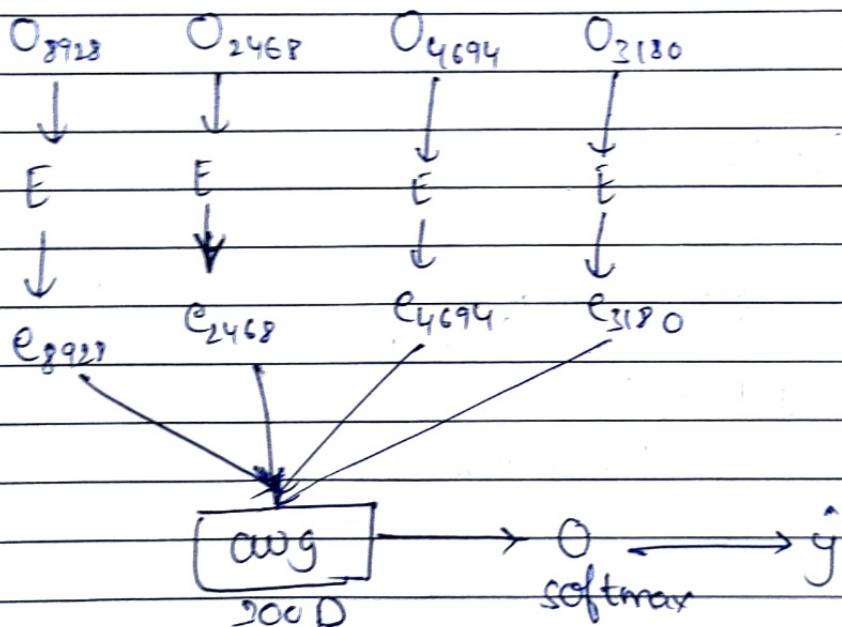
Video 8: Sentiment Classification

- * Huge label data set may not be available
- * Training sets are from 10000 to 100000.
- * That's why we use word embeddings

Example

The dessert is excellent

8928 2468 4694 3180



Drawback: ignores word order.

e.g. Completely lacking in service, good taste & good ambience.

algo will measure 3 "good" and ~~thank~~
evaluate the review to be good

Video 9: Debiasing Word Embeddings

- * bias here refers to undesirable bias like gender bias, ethnicity biases & soon.

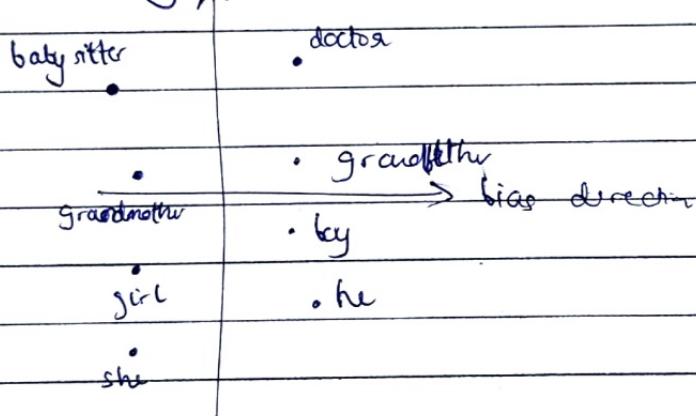
Example:

man: Computer programmer as woman: homemaker

father: doctor as mother: nurse

- * biases originate from the biases in text written by people.

Addressing bias in word embeddings



1. Identifying bias direction

$$\begin{pmatrix} \text{e_he} - \text{e_she} \\ \text{e_male} - \text{e_female} \end{pmatrix}$$

avg.

2. Neutralise: for every word that is not defined, project to get rid of bias