

Week 2

* Word Embeddings:

Word Representation

$$V = [a, aarden, \dots, zalm, \text{}]$$

1 hot representation:

• Man (5391)

$$\begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

05391

↳ 0 of one-hot

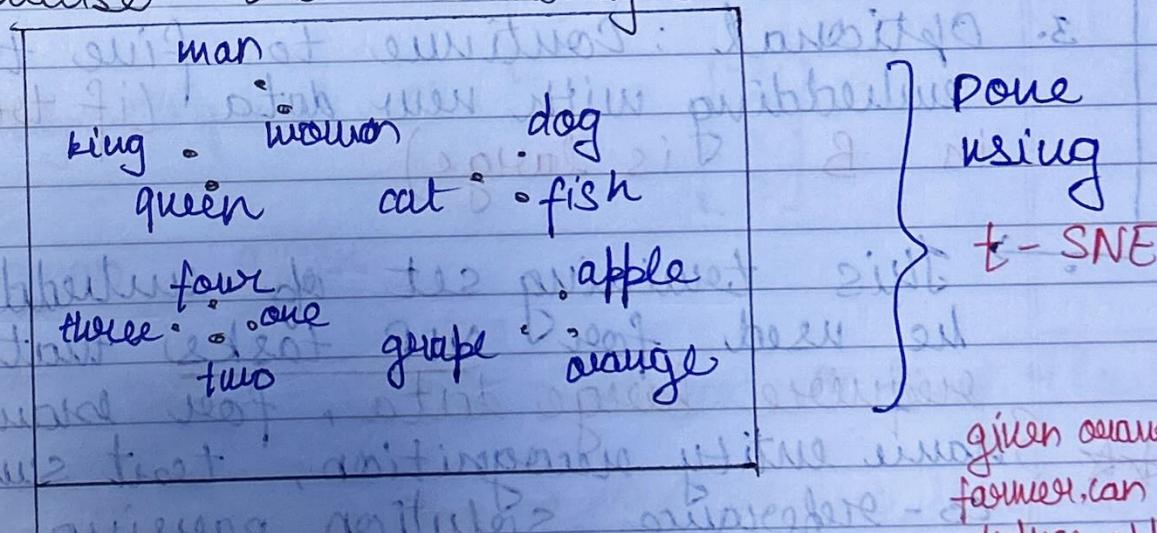
Disadvantage: Does not generalise well
 Eg: I want orange juice, it cannot
 guess I want apple juice here.

* Featureized Representation: Word Embeddings

	(5891)	(9853)	(4914)	(1157)	(456)	(6254)
Man	1	1	King	Queen	Apple	Orange
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	0.04	0.01	0.02	0.01	0.95	0.97
size						
cost						
verb						

300 dimensional vector to represent man,
 known as
 15391 → notation EMBEDDING

- Visualise Word Embedding S



Name Embedding Entity Recognition

Sally Johnson is an orange farmer
 Robert Linn is an apple farmer
 a durian cultivator

given orange
 farmer can
 deduce apple
 farmer.

- Using the features & similarity in word embeddings can also deduce that durian & cultivate are similar to fruits & farmer respectively.
- We also use transfer learning to then include the new deduced words from here, as part of vocabulary for other examples.

* Transfer Learning & Word Embeddings:

- Learn word embeddings from large text corpus (1-100B words)
→ or download pre-trained embedding online
- Transfer embedding to new task with smaller training set (say, 100K words)
- Optional : Continue to fine tune the embedding with new data (if training set in B is large)

→ This training set of embeddings can be used for NLP tasks that do not require large data, for example:

★ Name entity recognition, text summarization, co-reference resolution, parsing

- Implementing transfer learning on set 1, the advantage we get is to use less training data, smaller labelled text

We can use relatively lower dimensional feature vectors. Thus instead of 10,000 dimensional one-hot, use 300 dimensional dense vector.

- Transfer learning is not useful for language modelling and ML settings.
- Word embeddings show similarity to face encoding. Except here we may have a fixed vocab & there an unknown sea of judges.
- * Properties of Word Embeddings:

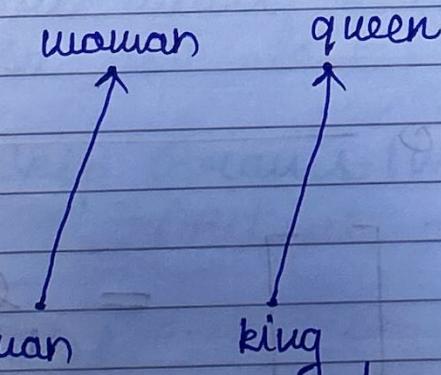
It can develop analogies:

Eg: man \rightarrow woman then king \rightarrow ?

(Take ref. from table 1 in word embedding)

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

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



$$\| \text{man} - \| \text{woman} \approx \| \text{king} - \| \text{queen}$$

$\| w \Rightarrow \max_{\text{maximum}} \sim (\| w, \| \text{king} - \| \text{man}) + \| \text{woman} \}$
similarity funcⁿ

easier visualisation of

t-SNE (300D \rightarrow 2D) for understanding, not common to have such easy 11gm like relation

* Similarity function: $\cos \theta = \frac{u \cdot v}{\|u\| \|v\|}$

Inference: COSINE SIMILARITY is based on the angle between vectors.

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

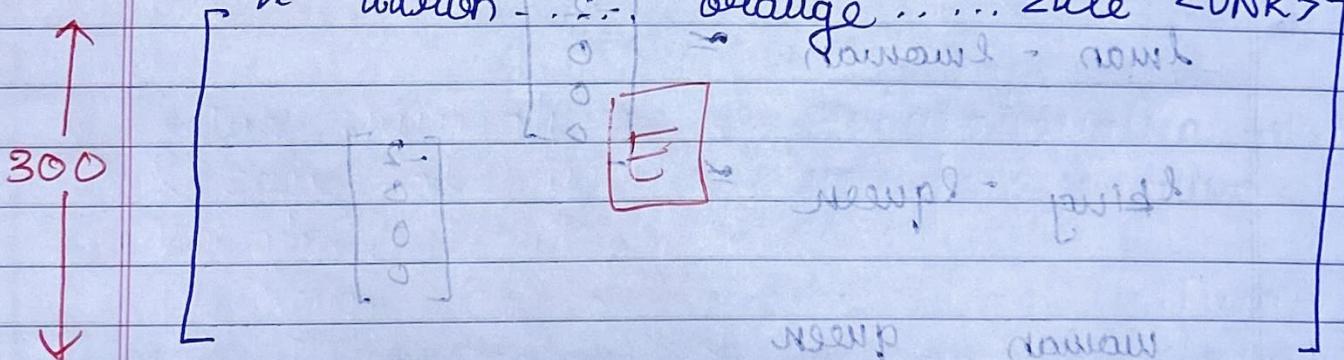
ref: $\cos \theta = \frac{u^T v}{\sqrt{\|u\|^2 \|v\|^2}}$

→ called cosine similarity as the formula is the cosine of the angle between vectors.

$\|u - v\|^2 \rightarrow$ Euclidean distance

→ rather a measure of dissimilarity

* Embedding matrix



(word - word, E - 10,000)

$$(300, 10k) \quad (10k, 1)$$

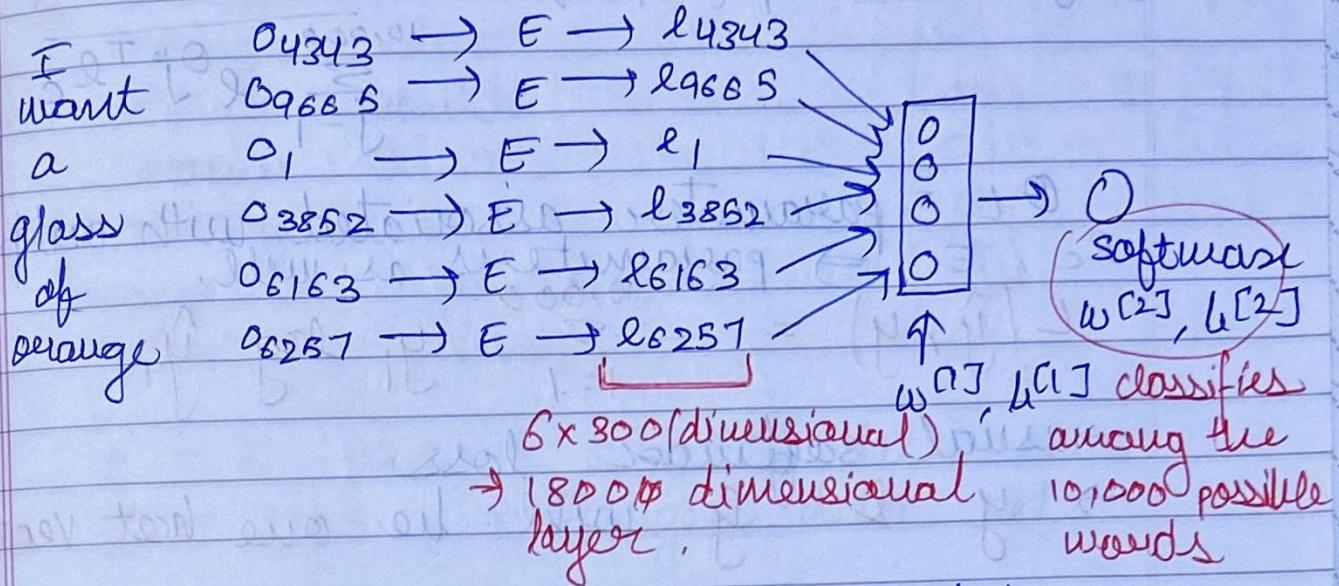
word - word

$$= \begin{bmatrix} & \\ & \\ & \end{bmatrix}_{\text{word}} = \begin{bmatrix} & \\ & \\ & \end{bmatrix}_{\text{embedding vector}}$$

E = Embedding matrix, O = One hot matrix

* Learning Word Embeddings :

I want a glass of orange juice.



if we always use a 4-word history, or a fixed history, then the dimensions of the layer will be $(4 \times 300) = 1200$ & thus it allows us to compute very long sequences as well.

If you want to learn a language model, you use the last 4 words. However if you want to learn word embeddings then you use last 1 word or nearly 1 word.

* Skip Grams (Word2Vec)

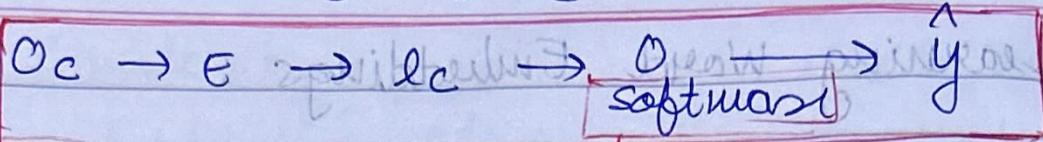
Context \rightarrow Target pairs

Model :

Vocab size = 10,000

Mapping : context c ("orange") \rightarrow target t ("juice")
 6257 4834

$$\mathbf{e}_c = \mathbf{E} \mathbf{o}_c$$



Software: $p(t|c) = \frac{e^{\mathbf{o}_t^T \mathbf{e}_c}}{\sum_{j=1}^{10,000} e^{\mathbf{o}_j^T \mathbf{e}_c}}$

\mathbf{o}_t = parameter associated with output t
 \mathbf{e}_c, \mathbf{E} → parameters as well.

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

→ usual software loss

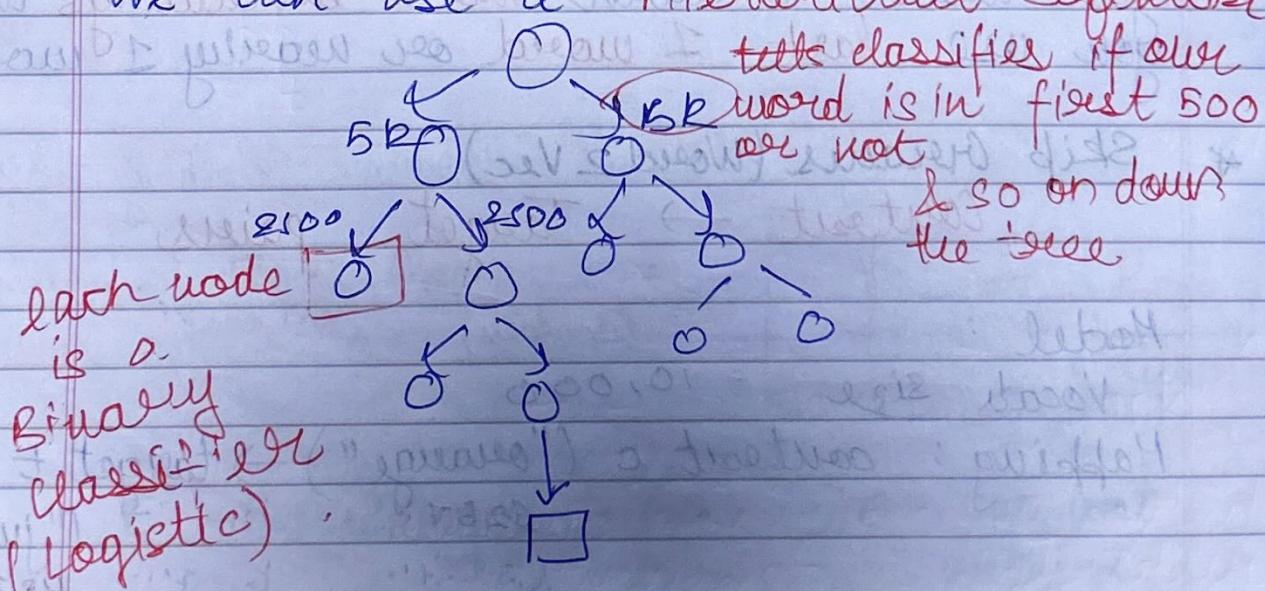
→ y & \hat{y} will be one hot vectors

* Problems with Software:

$$L_{\text{soft}}(p(t|c)) = \frac{e^{\mathbf{o}_t^T \mathbf{e}_c}}{\sum_{j=1}^{10,000} e^{\mathbf{o}_j^T \mathbf{e}_c}}$$

this summation is computationally expensive
& can take very long for bigger vocables.

We can use a Hierarchical Software:



* Negative Sampling :

Contest word	target ?
orange juice	1
orange king	0
orange book	0
orange the	0
orange of	0
(orange)	(juice / king / etc)
contest	→ target → label

$K = 5 \rightarrow 20 \Rightarrow$ smaller data-sets

$K = 2 - 5 \Rightarrow$ larger data sets.

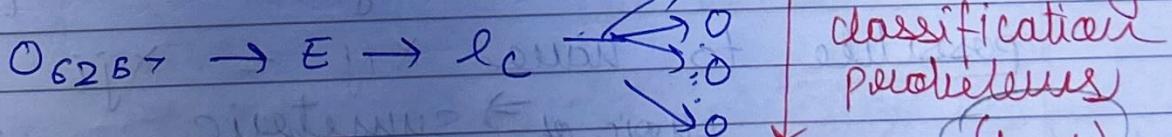
So, the problem statement is:

Given a pair of words, like orange & juice, do you think they appear together? Or we got them by sampling 2 words close to each other or as one word from test & the other chosen at random.

Model

$$P(y=1 | c, t) = r(\theta^T e_c)$$

think as $k+1$, negative to positive



→ On every iteration, we use only $(k+1)$ words, where 1 word is the answer & there are K randomly close - ve values.

How to choose -ve samples?

1. Choose according to empirical

~~very non-selective~~ importance, but that may words like: the, and, a, for, etc gain a lot of importance.

2. There is to divide by $1/V$, sampling -ve examples uniformly, $1/V$.

Thus middle ground is:

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

* Thus softmax vs -ve sampling
→ summation over 1000 binary classification.

* GloVe (global vectors for word represent.)

→ count that captures how often they appear together
 x_{ij} → times j appears in context of i

possible to have $x_{ij} = x_{ji}$

choice of \rightarrow symmetric
context is
→ target words appear in
+/- 10 words of each other

→ not symmetric for choice of context
words being immediately before/after.

Model :

$$\text{minimise} : \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(x_{ij}) (\underbrace{\theta_i^T \ell_j + b_i + b_j - \log x_{ij}}_{\text{weighted terms}})$$

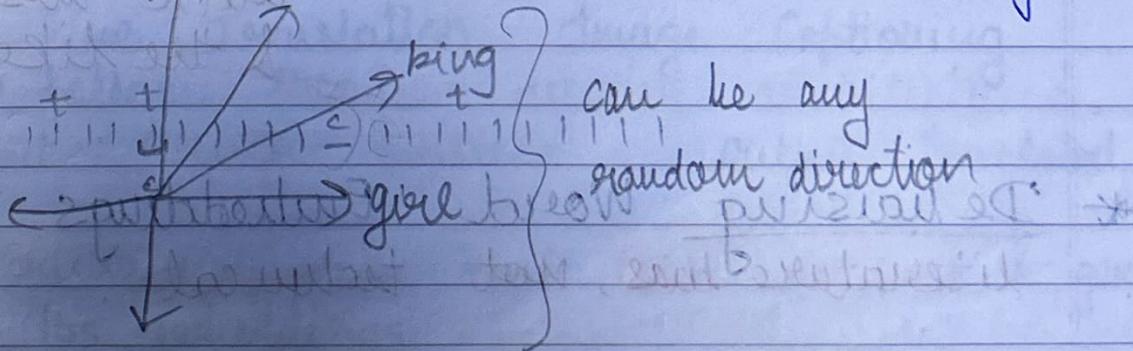
$f(x_{ij})$ weighted terms,

- to give $f(x_{ij}) = 0$ for $x_{ij} = 0$
- does not give too much importance to θ_i and the etc.
- does not give too little importance to dimension etc.
- thus valid only when the words have co-occurred atleast once.
- here, θ_i & ℓ_j are symmetrical, as they play pretty much the same role & we could reverse them or swap them & achieve same objective.

$$\ell_{\text{final}} = \ell_w^{\text{final}} = \ell_w + \theta_w.$$

With an algo like this, we cannot guarantee direction of axes & minimise priority of features as:

$$\theta_i^T \ell_j \Rightarrow (\theta_i^T)^T (A^{-T} \ell_j) = \theta_i^T A^T A^{-T} \ell_j = \theta_i^T \ell_j$$



* Sentiment Classification :

The dessert is excellent

8928

2468

4694

3180



The 08928 → E → l8928
 dessert 02468 → E → l2468
 is 04694 → E → l4694
 excellent 03180 → F → l3180

3000

Avg

 software
1-5

→ decent algo, however fails when:

The hotel lacked good audience, good taste and good service.

→ Acc to avg, good will give this 3/5 star while it is 1 star.

→ Use RNNs with word embeddings.

to deal with avg's disadu. & consider lacked to deal with sequences that replace 'lacked' with 'absent' & the like.

* Debiasing Word Embeddings literature bias, not technical

man : woman as king : queen
 man : computer-programmer as X
 woman : homemaker

7 Father: Doctor as Mother: Nurse X,
 reflect biases of the text used to
 train \Rightarrow word embeddings,

* Addressing bias in word embeddings
 (not bias direction)

1. Identify bias direction \rightarrow doctor

{ he - she
 female - female }

{ male - female }

→ average

babysitter

grandmother
 girl

she

bias

grandfather

- boy
- he

2. Neutralise :

For every word that
 is not definitional, \rightarrow not affected
 by bias.
 project to get rid of bias.

3. Equalise pairs