

Week 2: NLP & Word embeddings

Intro to word embeddings:

1) Word representation :

$$V = [a, \text{anon}, \dots, \text{zulu}, \langle \text{UNK} \rangle]$$

$$\|V\| = 10000$$

Man
(5391)

Woman
(9853)

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

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

I want a glass of orange juice

I want a glass of apple ?

e_{5391}
(One-hot vector)

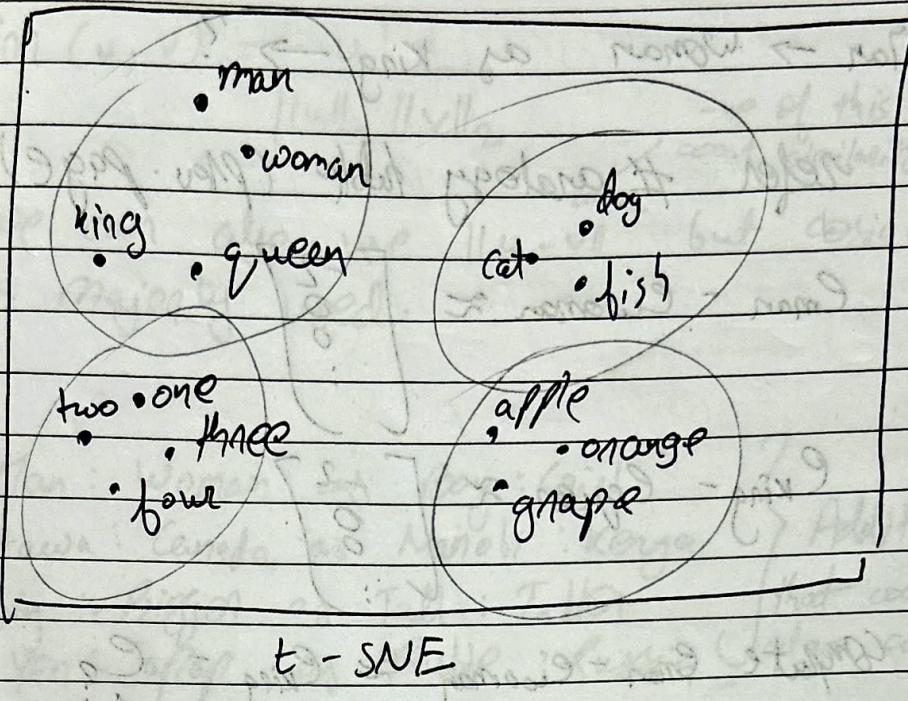
e_{9853}

Featureized representation : word embedding

Analogy table

	Man (5391)	Woman (9853)	King (5914)	Queen (7157)	Apple (456)	Orange (8297)
Gender	-1	1	-0.95	0.97	0.00	0.00
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	:	:				
hours						
verb...						
	$e_{5391} = e_{\text{man}}$	$e_{9853} = e_{\text{woman}}$				

Visualizing word embeddings:



2) Using word embeddings

Name entity recognition example

Sally Johnson is an orange farmer

Using a model that uses word embeddings, we will be able to predict that Robert Lin is also a person.

Robert Lin is an apple farmer

The algorithm will identify orange and apple being belonging to the same class hence output Robert Lin as a person.

Instead of apple farmer, if we use durian another fruit cultivation, our algo will still be able to give us correct outputs

3) Properties of most word embeddings:

e.g. Man \rightarrow Woman as King \rightarrow ?

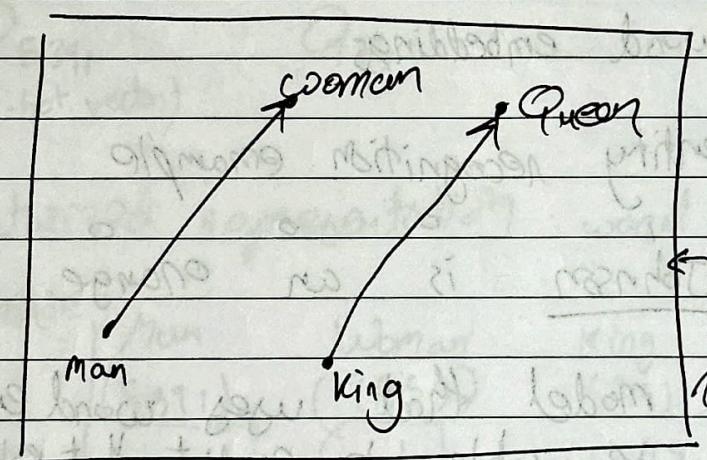
refer # analogy table (p16v. page)

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

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

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

↑
queen



← 300 Dimensional
where words are
represented as a point.

Find word $w : \arg \max_w \sin(e_w, e_{\text{king}} - e_{\text{man}} + e_{\text{woman}})$

↑
↑ maximizes similarity

Returns the word w which has maximum similarity to the equation $e_{\text{king}} - e_{\text{man}} + e_{\text{woman}}$ and it should be the word queen with accuracy of $70 - 75\%$.

Cosine Similarity

$$\text{Sim}(u, v) = u^T v$$

$$\|u\|_2 \|v\|_2$$

-up of this as we
want similarity not dissimilarity

We can also use $\|u - v\|^2$ but cosine similarity is majorly used.

Man : Woman as Boy : Girl

Ottawa: Canada as Nairobi: Kenya

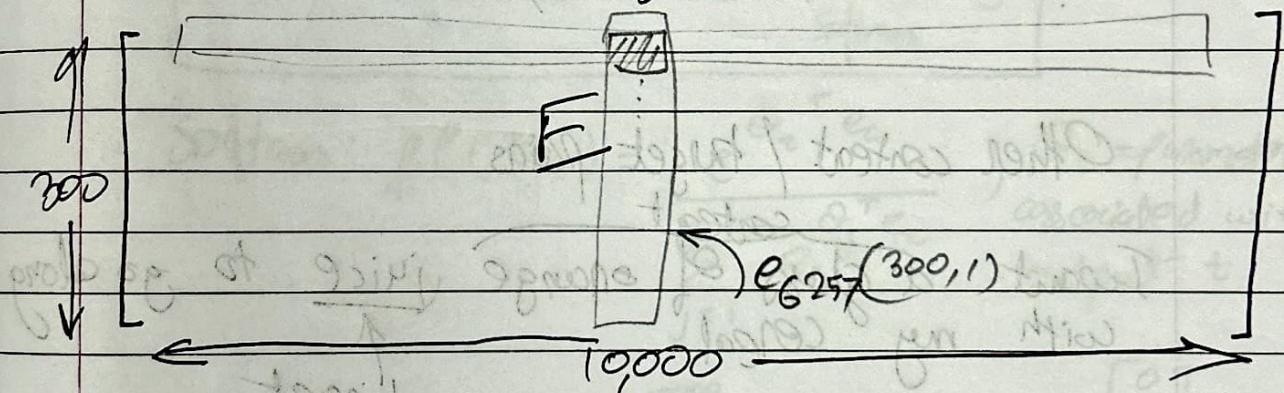
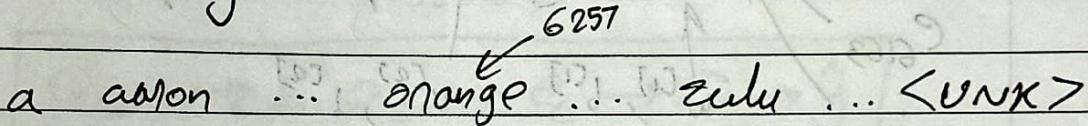
Big - Bigger as Tall : Taller

Yen: Jay Ray as Rubio? Russia's Jalgo can do.

Addition Examples

that could emb.

4) Embedding matnix



O	6257	?	o
o		o	o
?		!	← 6257
o		o	v

$$E. O_{6257} = \boxed{71} - C_{6257}$$

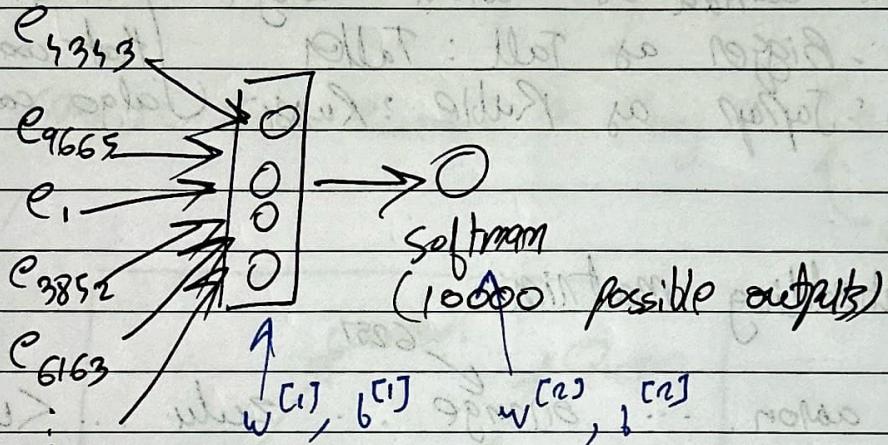
1) Learning word embeddings:

I want a glass of orange .

$I \rightarrow O_{4343} \rightarrow E. O_{4343} \rightarrow e_{4343}$ embedding vectors (300, 1)

want $O_{9665} \rightarrow E. O_{9665} \rightarrow e_{9665}$

a $O_i \rightarrow E. O_i \rightarrow e_i$



Other context / target pairs

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

Content: 1) Last 4 words.

- Feed into a NN ← 2) 4 words on left ? a glass of orange ?
- 3) & right to go along with
- 4) Last 1 word → orange ?
- 5) Nearby 1 word → glass ?

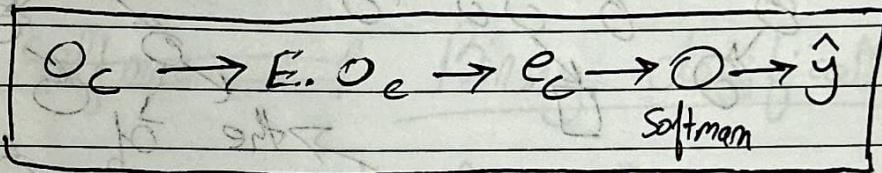
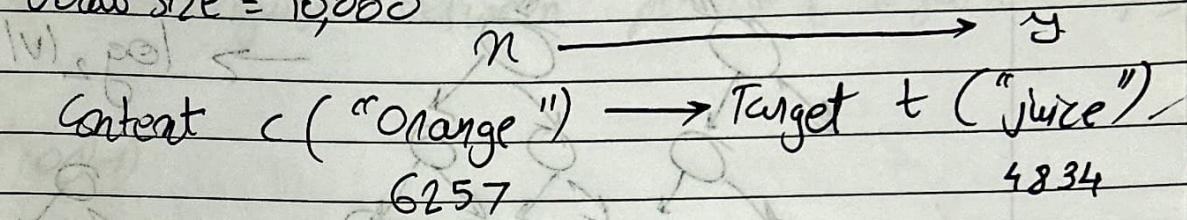
2) Word2Vec Skip-grams Model:

~~Skip-grams~~

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

	Content	Target
Model	Orange	juice
	orange	glass
	orange	my

Vocab size = 10,000



Softmax: $P(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10000} e^{\theta_j^T e_c}}$, θ_t = parameter associated with output t .

$L(\hat{y}, y) = - \sum_{i=1}^{10000} y_i \log \hat{y}_i$; $y \begin{bmatrix} 0 \\ i \\ 0 \end{bmatrix} \leftarrow 4834$

Then compute Back prop (One-hotvector)

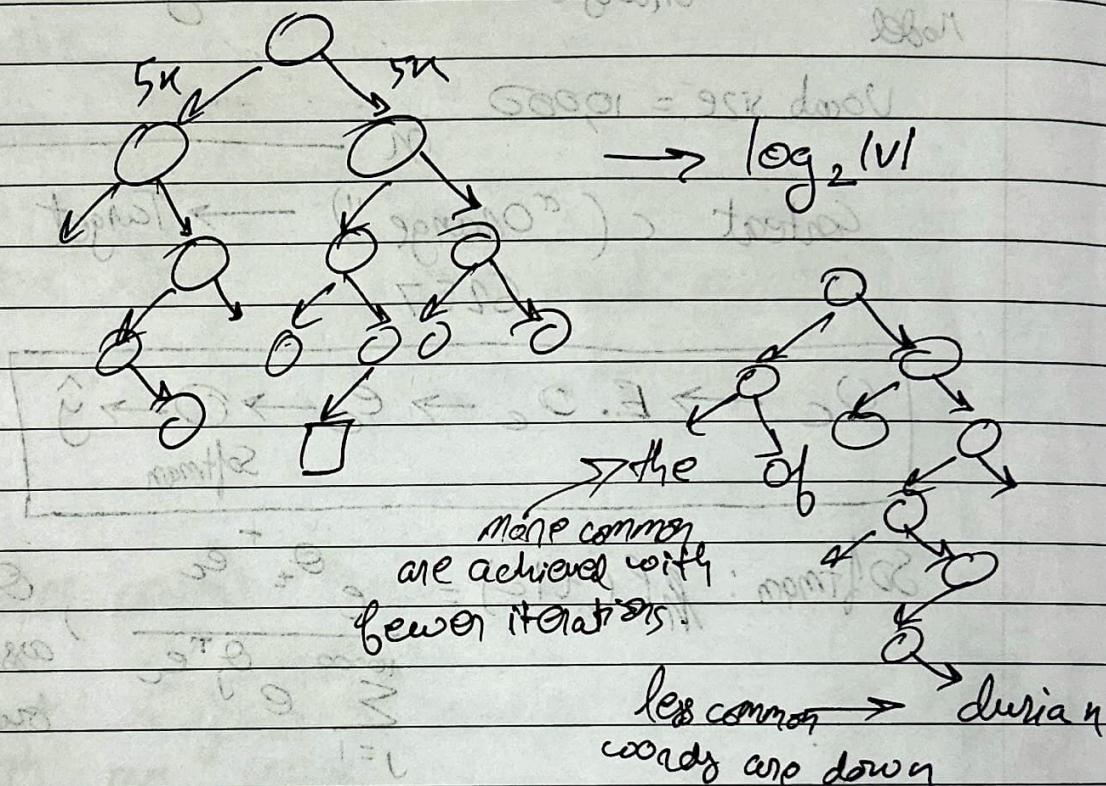
Problems with softmax classification:

$$P(+|c) = e^{\theta_i^T c}$$

$$\sum_{j=1}^{10000} e^{\theta_j^T c}$$

because of the denominator, this could become computationally time consuming as 10000 iterations make it slow. 1M on 1B words would make it even slower.

Hierarchical softmax classifier:



3) Negative Sampling

Defining a new learning problem

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

Pick a content-target pair from scratch
and set the Date now, target to 1

Content word	target?
orange juice	1
orange of	0
Orange king	0
orange the	0
orange book	0
Random words from the dictionary	0

positive example
set these to 0 even if any word appears in the sentence (negative examples)

Create a supervised learning problem, where the algorithm inputs (c, w) (content and word) and has to predict the target $y_{(c,w)}$ i.e. make to distinguish between positive & negative example

Model

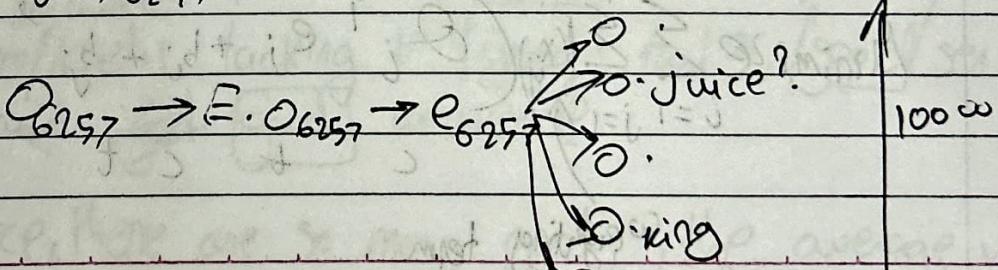
Softmax \rightarrow Binary classification

Using logistic regression,

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

↑ ↑
parameter vector θ_t embedding vector e_c

Orange \rightarrow 6257



* instead of 10,000 we can update k+1 units

Selecting negative example:

$$P(w_i) = \frac{f(w_i)}{\sum_{j=1}^n f(w_j)}$$

$f(w_i) \rightarrow$ observed frequency

4) GloVe word vectors

→ global vectors for word representations

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

c, t

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

$x_{ij} = x_{ji} \leftarrow$ may / may not be true, depends on the no. of words you're considering i.e. either in the range of ± 10 words or the word immediately next

Model:

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

Weighting term

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

$f(x_{ij}) \rightarrow$ this, is, of, a, ... frequent words
 weight not too big
 Julian, ... infrequent words
 weight is not too small.
 given by $f(x_{ij})$

Q_{is} and e_j are symmetric

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

Applications of Using Word Embeddings

1) sentiment classification (like or dislike)

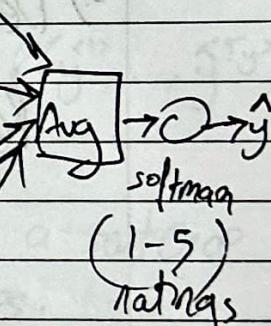
eg: The dessert is excellent 4-★ Review
 8928 2468 4694 3180

The $O_{8928} \rightarrow E. O_{8928} \rightarrow e_{8928}$

desert $O_{2468} \rightarrow E. O_{2468} \rightarrow e_{2468}$

is $O_{4694} \rightarrow E. O_{4694} \rightarrow e_{4694}$

excellent $O_{3180} \rightarrow E. O_{3180} \rightarrow e_{3180}$

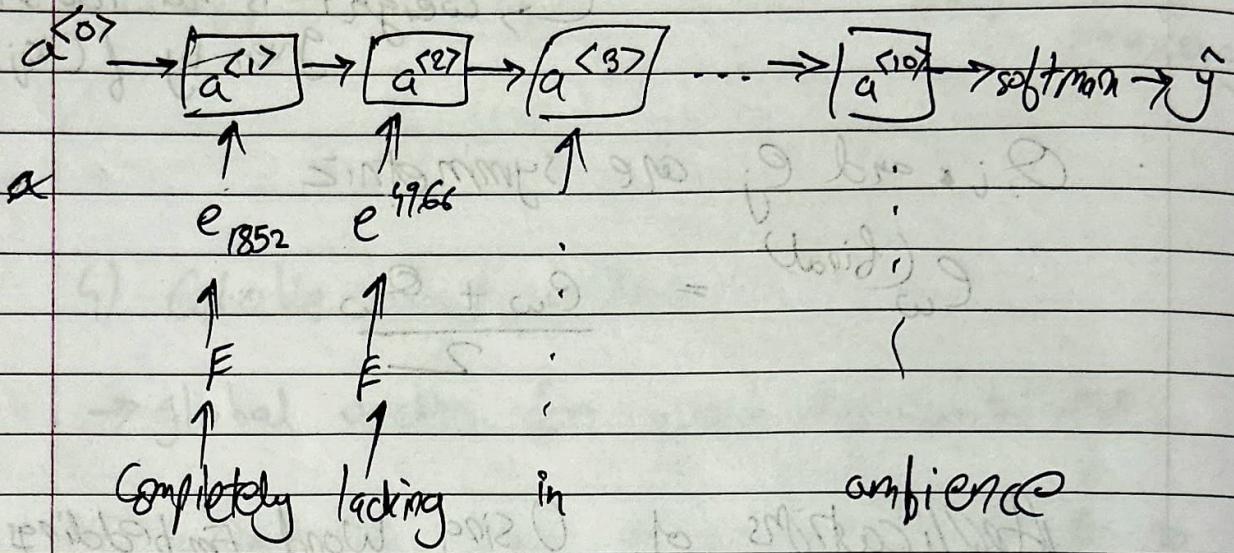


Drawback:

completely lacking in [good] taste, [good] service,
 and [good] ambience

Since there are so many good, the average will
 orient itself to good but in reality it's a 1★ review

RNN for sentiment classification



This algorithm will do a better job at classification as it can work on words that are not in the training set. This is a type of many-to-one architecture.

2) Debiasing word embedding

Refer WA