

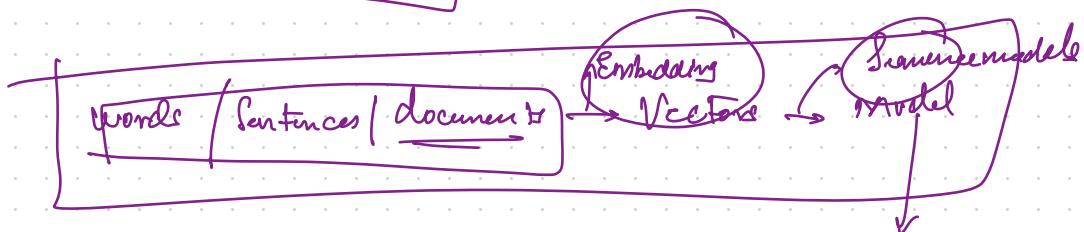
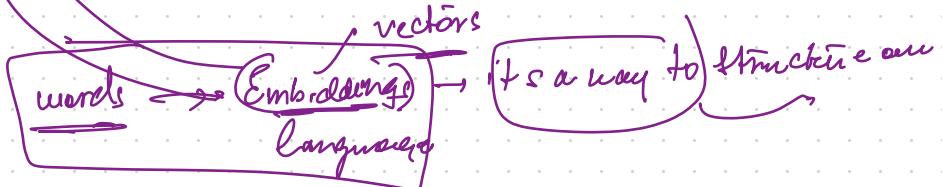
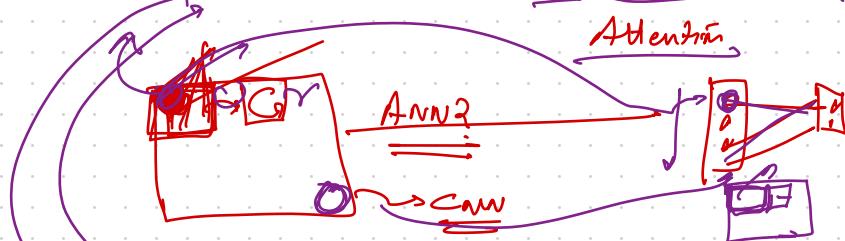

Natural language Processing

① NLP problems are easel.

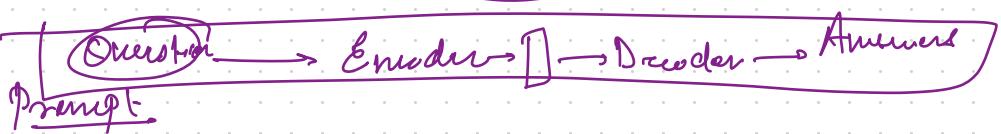
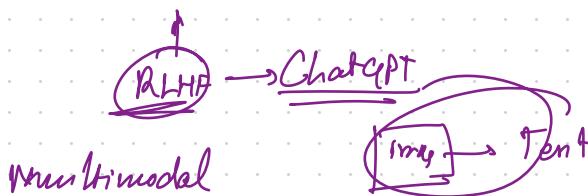
② How do we structure words / sentence

Feature engineering

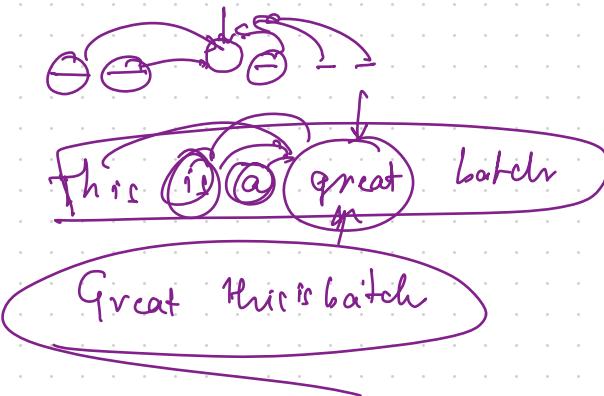
③ Same models → RNN, LSTM, GRU, Transformer



large language models



Natural Language - is a Sequence



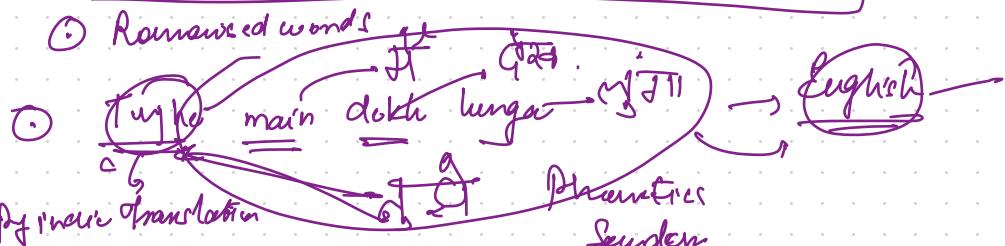
Natural language Processing

Understanding information from text

- ① Unstructured
- ② Context / Sentiment / Change / Sarcasm

③ Indigo -----

④ Romanized words



- ⑤ Structured

Use Cases:

Translation

Gmail auto type → auto fill

Search engines

Sentiment analysis

Small spam → Ham

Chatbot

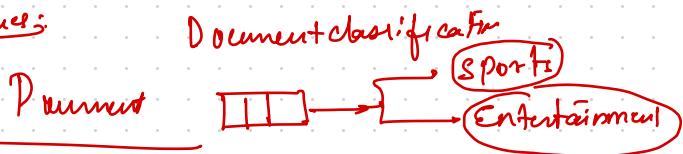
Text classification, Sum

WFR, information

$$\textcircled{1} = f(x) \quad \begin{matrix} \text{Elemental} \\ \uparrow \uparrow \uparrow \end{matrix} \quad \begin{matrix} \text{language model} \\ \boxed{\text{language model}} \end{matrix}$$

Sentences are very unclear

Preprocessing techniques:



Special characters

Removal of unwanted Characters or patterns of characters from our source sentence.

Ticket classifier

Sentiment classification

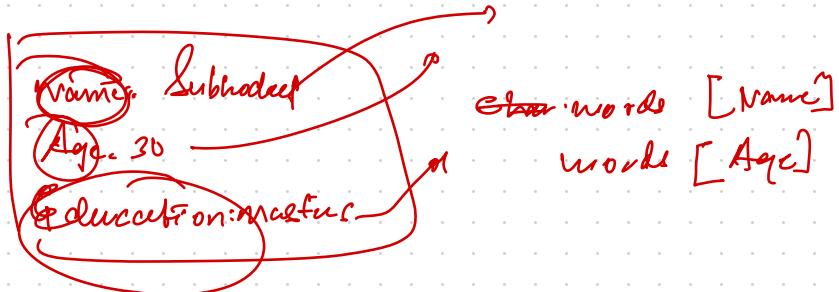
Ticket

Regex → a process of defining certain patterns using certain rules.

"I am very hungry. I will order and give to my friend!"

Region → Extract some defined patterns
from sentences

[w]
[s]



Stopwords

(an), a, the
because and Then

Translation

, what bot → Tomo tan 25°C
→ I am com - - - - -
What's the temp. phenomenon

It's time to go home → Tom (S) air going home

Stemming / lemmatization

Words → Embedding
→ Vectors → Model

Embedding

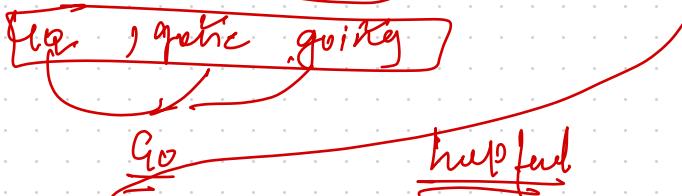
✓ Frequency based
✓ Count of the words
in the sentence

Context?

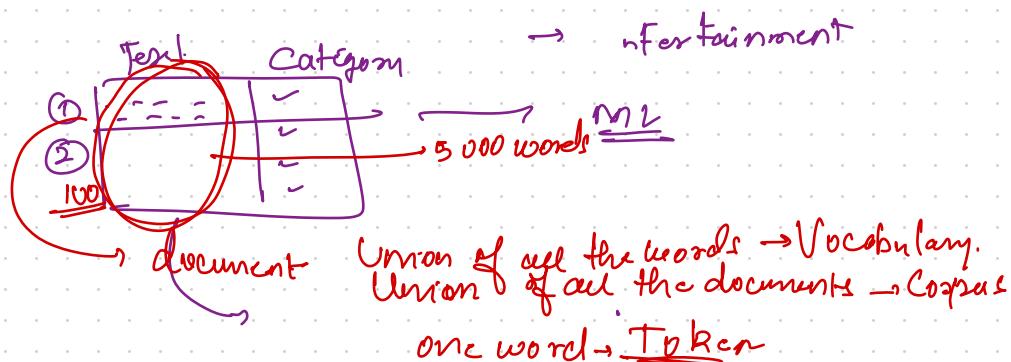
Randomness based:
vectors generated
through LNN.

If am going home, home is where heart is

g:1, am:1, going:1, home:2, ts:2
where:1, heart:1

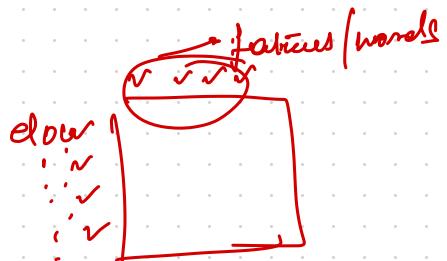
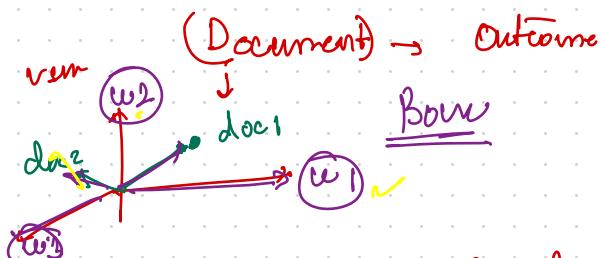


How to represent sentences/documents in a vectorised form using frequency embeddings.



Bag of words model

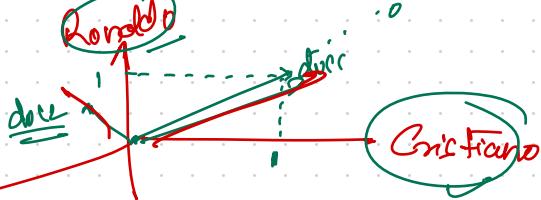
Document classification



- Christian Ronaldo got transferred
→ Sports

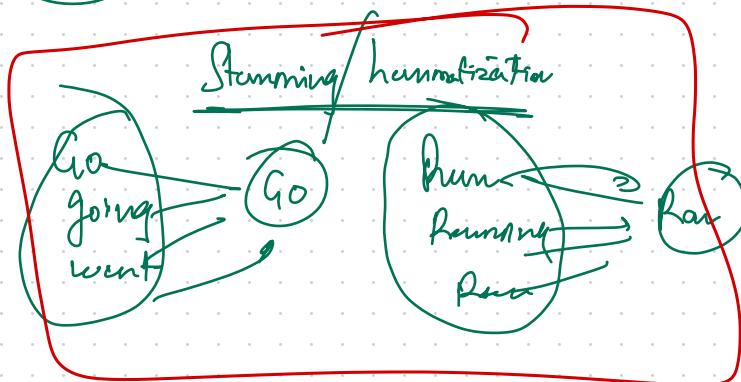
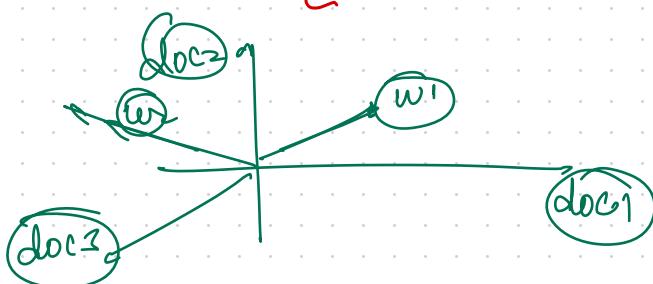
- Pierce Brosnan performs in Casino.

	Cristiano	Ronaldo	got	transferred	Dna	lipa	Performs in	Comme
①	1	1	1	1	0	0	0	0
②	0	0	0	0	1	1	1	1



Problem: ① semantic dimensions

- ① does not maintain context
 - order not maintained
 - By definition of semantics all words are independent





Lernstruktur

A handwritten diagram on lined paper. At the top right, the word "Verb" is written above the underlined word "Run". Below "Run", the word "Ran" is also underlined. To the left of "Run", the word "Running" is written above a bracket that spans both "Run" and "Ran". Below "Running", the word "Ran" is written above another bracket that spans both "Run" and "Ran". Red arrows point from the "Running" bracket to the "Run" and "Ran" underlines, and another red arrow points from the "Ran" bracket to the "Run" and "Ran" underlines.

Going to go

Spreading ✓
Specific ✓
Conscious ✓
Grey → Grey

The diagram illustrates the stemming process. At the top, a red bracket groups the sentence "I am great at public speaking". Below it, a green bracket groups the word "I". A red bracket groups "am", "great", and "at". Another red bracket groups "public" and "speaking". Below these, the word "I" is labeled "stemming I (word)". The word "am" is labeled "humantia". The words "great", "at", "public", and "speaking" are grouped under the label "sentences → words".

Sentence Tokenizer
Documents \rightarrow Sentence

• split('')

"I | am | going | to | my | home". Split("")

[I, am, going, to, my, home]

write \rightarrow word-tokenizer

tweet \rightarrow words

(count of the word in the tweets, count of the words in negative tweets)

+ve counts

-ve counts

tweet \rightarrow (Count of the word, Count of the -ve words, +ve words)

tweet: The Game

+ve 50 10
+ve 80 60
+ve 100 -ve 100

tweet (110, 110)

fear

$$\left\{ \begin{array}{l} (\text{game}, 0) : 20 \\ (\text{game}, 1) : 40 \\ (\text{awesome}, 0) : 50 \\ (\text{awesome}, 1) : 90 \end{array} \right. \quad \left. \right\}$$

fact: game is awesome

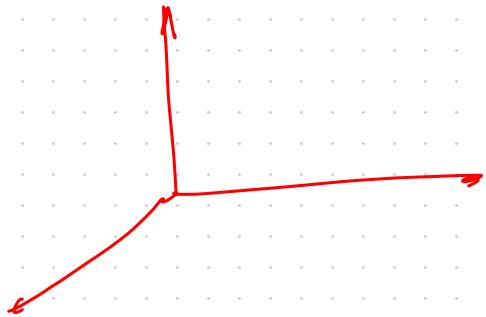
game awesome

$$\text{pos} = \text{fear}[(\text{game}, 1)] = 40 + \text{fear}[(\text{awesome}, 1)]^{, 90}$$
$$\text{neg} = \text{fear}[(\text{game}, 0)] = 20 + \text{fear}[(\text{awesome}, 0)]^{, 50}$$

$$\text{pos} = 130$$

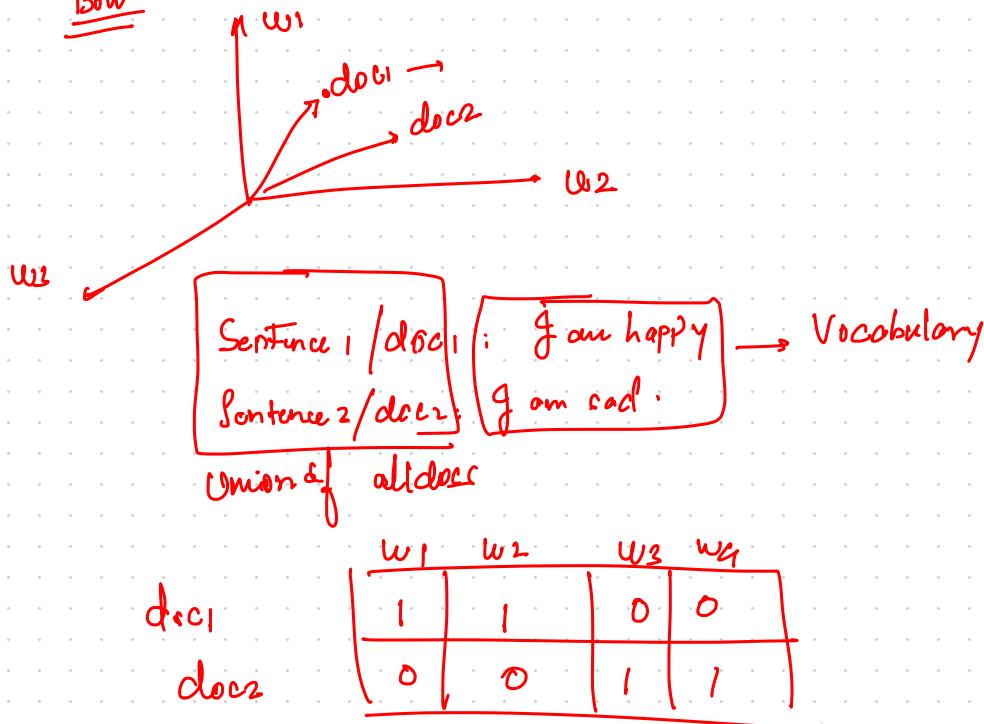
$$\text{neg} = 70$$

$$\text{fear} = (1, 130, 70)$$

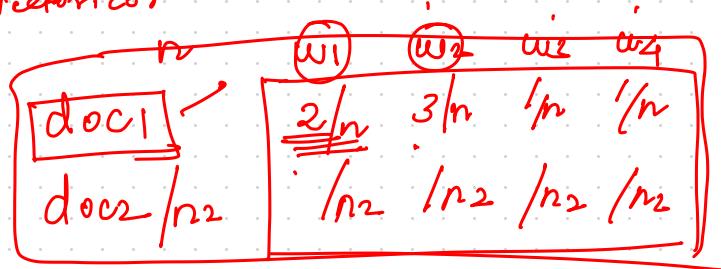


NLP \rightarrow words / docs into vector space models

BOW



Count vectorizer



① do not maintain any order

doc1 : I am happy today. today help I am
 I am happy today

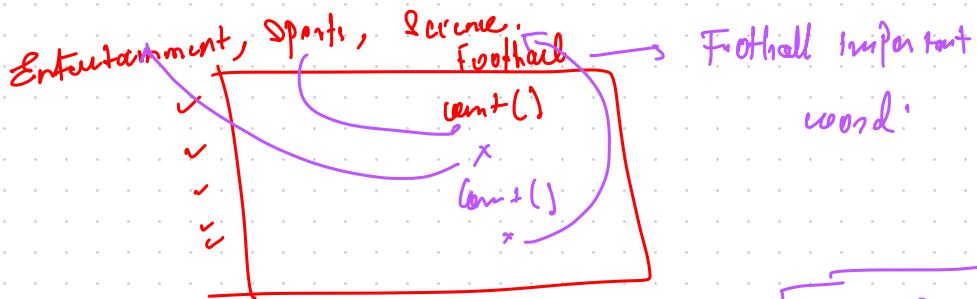
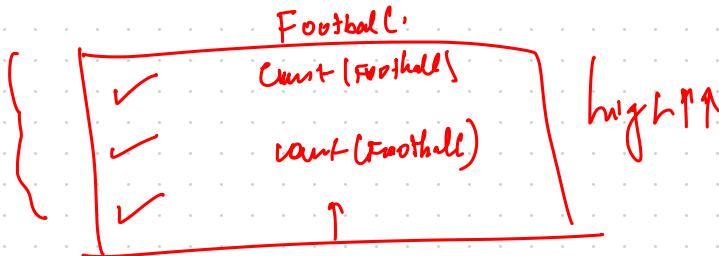
doc1 :

1	1	1	1	1
---	---	---	---	---

frequency based embeddings don't capture the context

Term Frequency vs Document Frequency

Football



$$\text{tf}(\text{Term frequency}) \times \text{idf}$$

$$\text{tf}(w_2, \text{doc}_1) \times \frac{\log\left(\frac{N}{n}\right)}{\text{Total number of tokens in doc}_1}$$

where $w_1, w_2, w_3, \dots, w_m$

$N = \{ \text{doc}_1, \text{doc}_2, \text{doc}_3, \dots \}$

number of documents in which w_i occurs.

doc1: Ronaldo is lovely, Ronaldo is rich.

doc2: Ronaldo is famous.

doc3: Messi is famous.

doc1: Mondo is lovely, Ronaldo is rich.

doc2: Ronaldo is famous.

doc3: Mussi is homely.

Normalized term frequency $\times \log \left(\frac{N+1}{n} \right)$

$\frac{1}{3}$

doc1

doc2

doc3

Mondo

Mussi

lovely

rich

famous

$\log \left(\frac{4}{2} \right)$

$\log \left(\frac{4}{2} \right)$

$\log \left(\frac{4}{2} \right)$

$\frac{2}{4} \times \log \left(\frac{4}{2} \right)$

$\frac{1}{4} \times \log \left(\frac{N+1}{n} \right) = \frac{1}{4} \times \log \left(\frac{3}{2} \right) = n \times \log \left(\frac{1}{n} \right) \rightarrow 0$

$\uparrow \log \left(\frac{N+1}{n} \right) \rightarrow n \approx N$

$\frac{3+1}{2}$

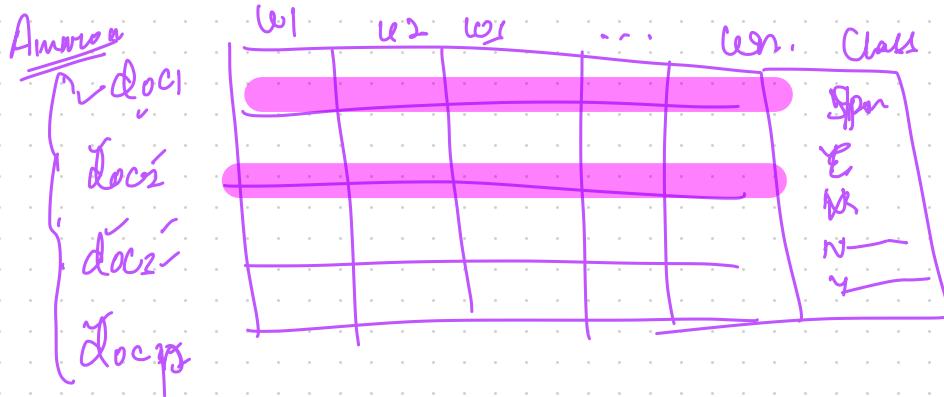
$\log \left(\frac{N+1}{n} \right) \approx 1$

TF-IDF (t_i)_{doc1}

= Normalised term frequency $\times \log \left(\frac{N+1}{n} \right)$

$N \rightarrow$ total number of docs in corpus.

$n \rightarrow$ total number of docs in which term occurs.



Find out Similar Documents

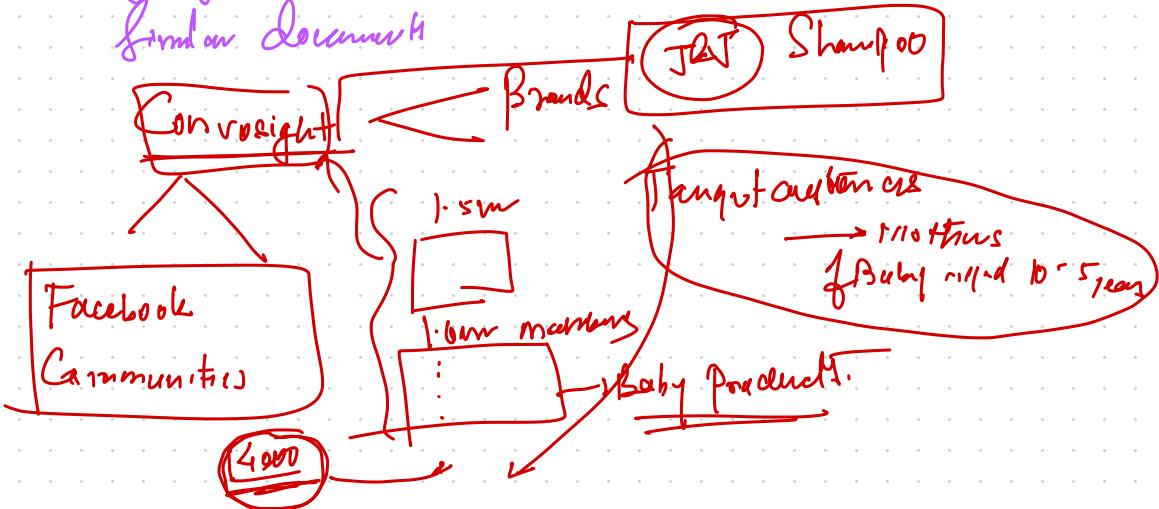
doc1 (Product)
doc3 (Product) ← → :

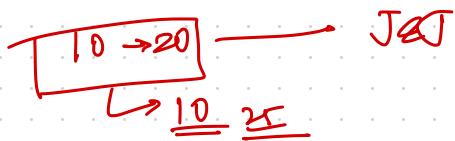
addresses



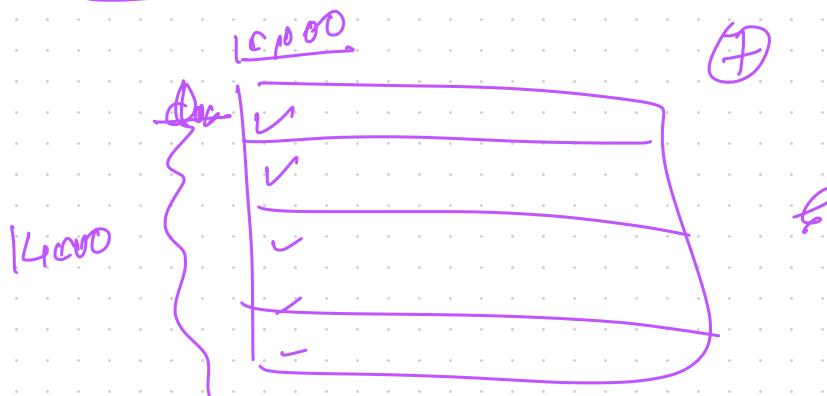
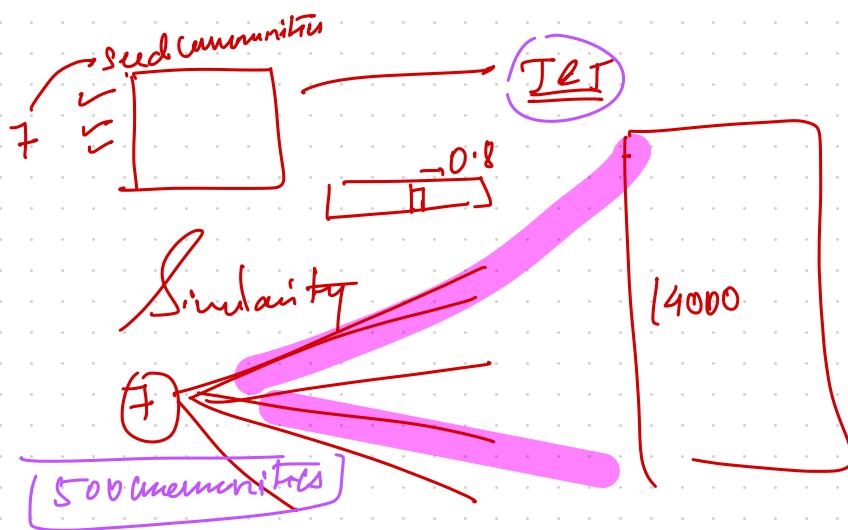
Using document based embedding of word2vec

Find our documents





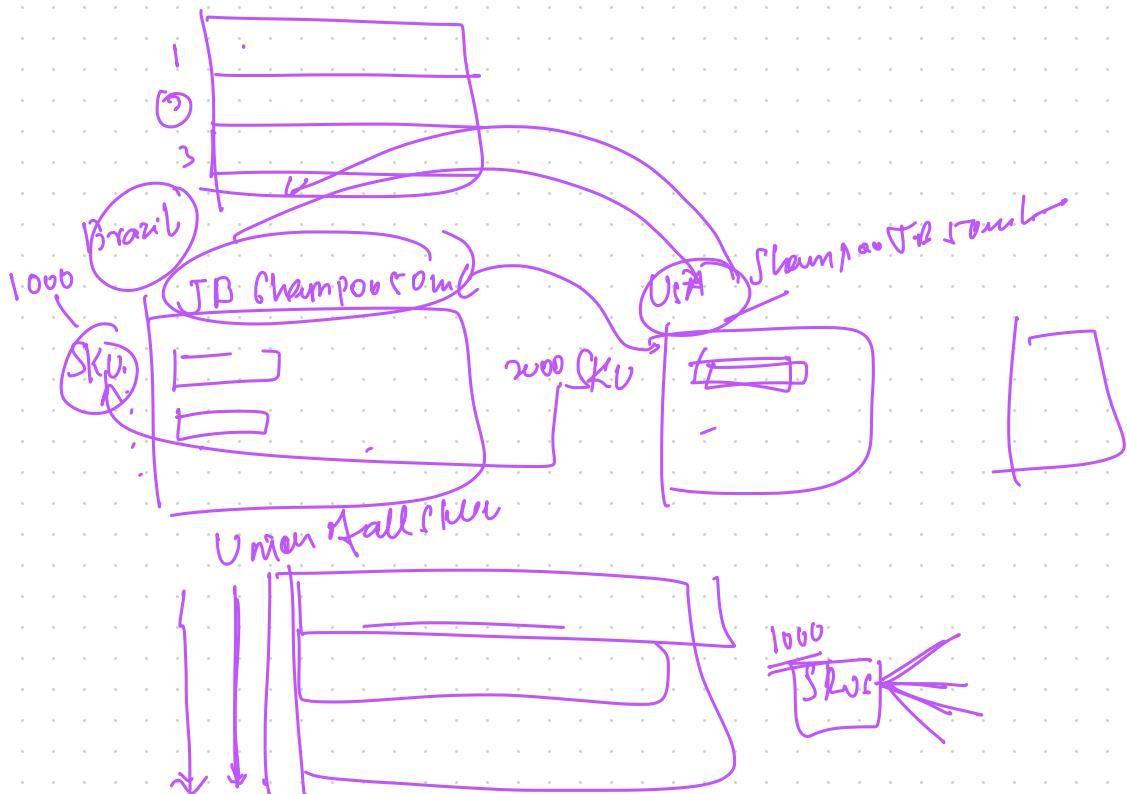
1000 users
↳ 1.6 billion conversations

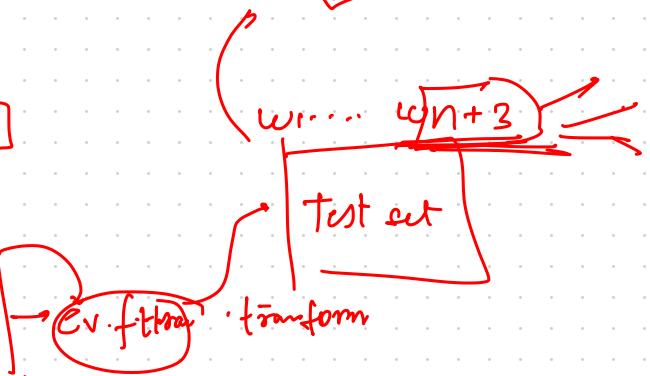
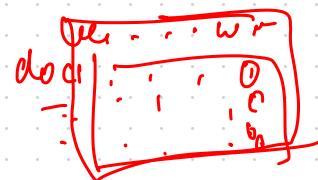
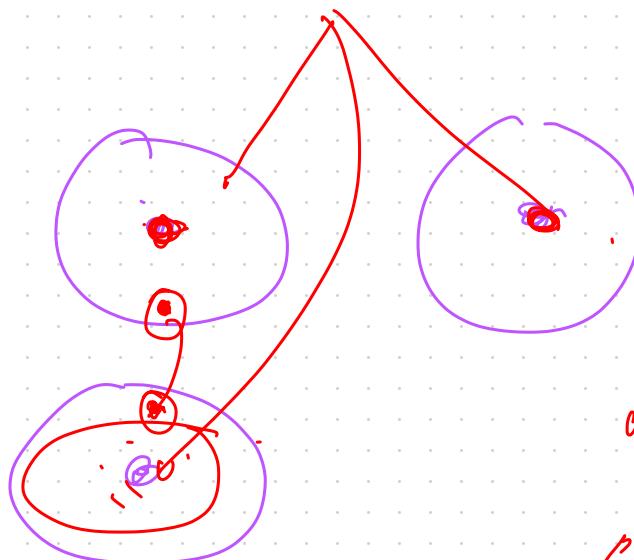


fan  happy to teach you

Food and drops for
weather, nutrients for babies

  Hey mother, see this new recipe
of created





corpus: [" Jan - . , " - ...]

2 for sentence in corpus

log(n/m)

sentence

two happy friends

doc1

$\times \log(n/m)$

doc2

$\times \log(M/m)$

|ʃ| |ən| |'gret|, thank you.

Undergroup:

|ʃ| |ən| |'gret| |θank| |yου|

Bigram: |ʃ| |ən| |'gret| |θank| |yου|

|ʃ| |ən| |ən| |'gret|

Embedding \rightarrow wordvec

Frequency based embedding

words

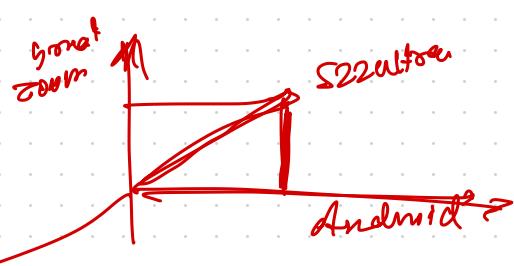
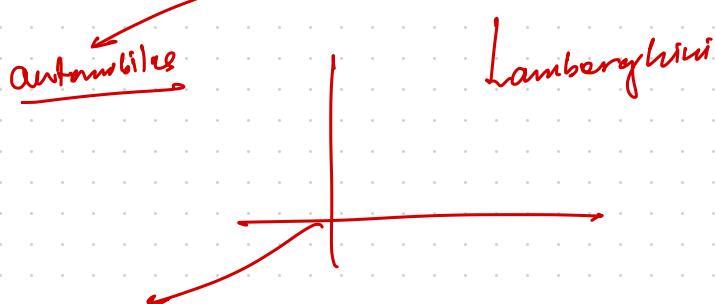
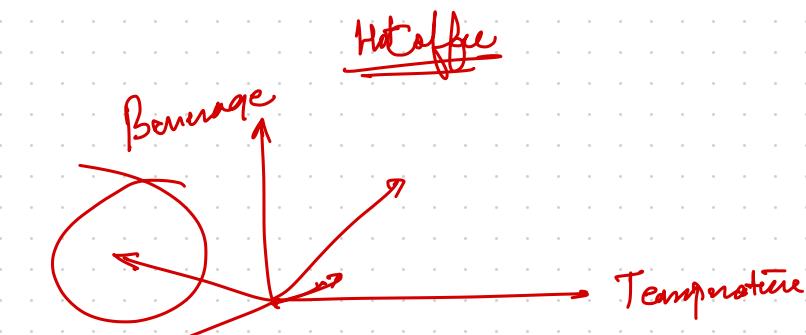
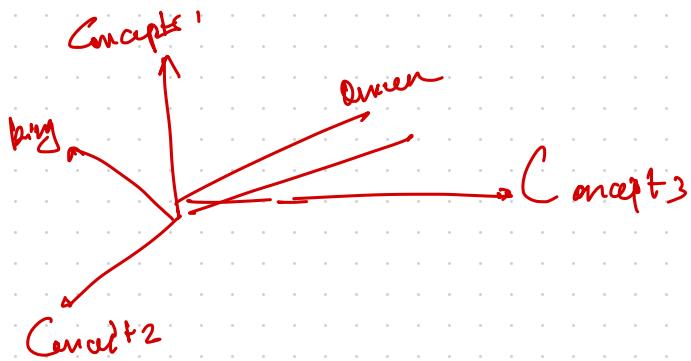
docs

v	v	v	f
v	-	-	-
.	-	-	-
.	-	-	-

Prediction based Embedding



Prediction based embedding



Embeddings are dense representation of concepts

Abstract nonlinear combinations
of multiple words.

Embedding's Dimensionality

Word are an model

Embedding SVD $\propto \propto x$

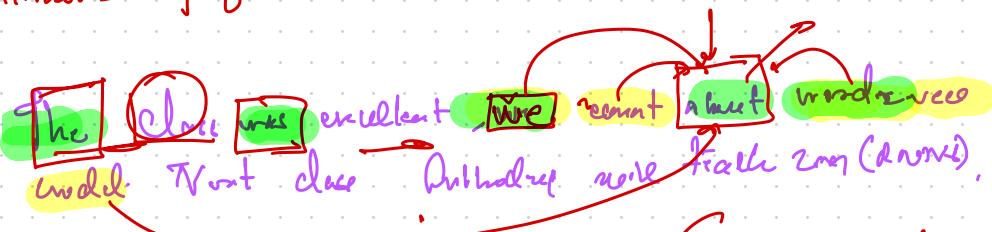
Embedding from sequences models \rightarrow (Supervised way)

Lm. Seqs (p_{lm}, h_{lm})

Word2vec

cbow
continuous Bag of words

Skip gram model



Context words
The, was
close, excellent
was, we

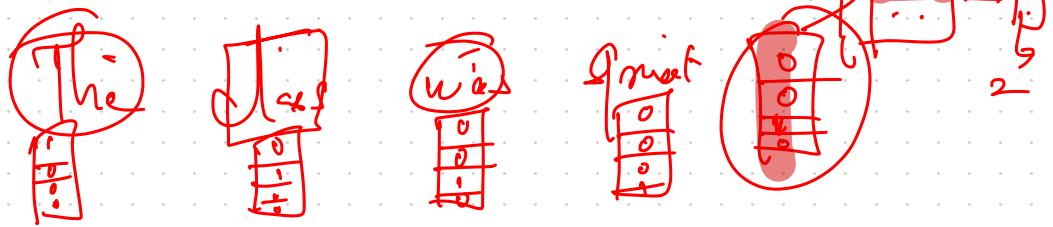
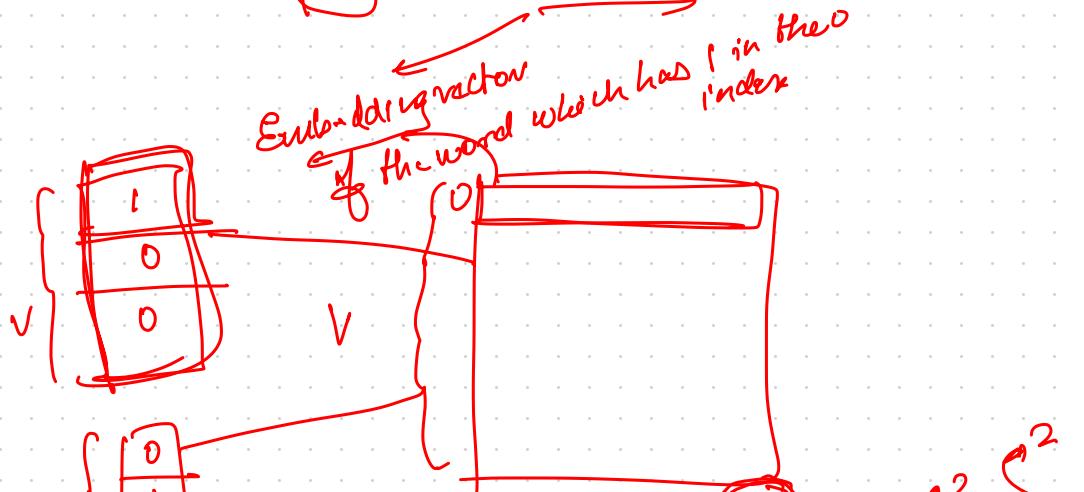
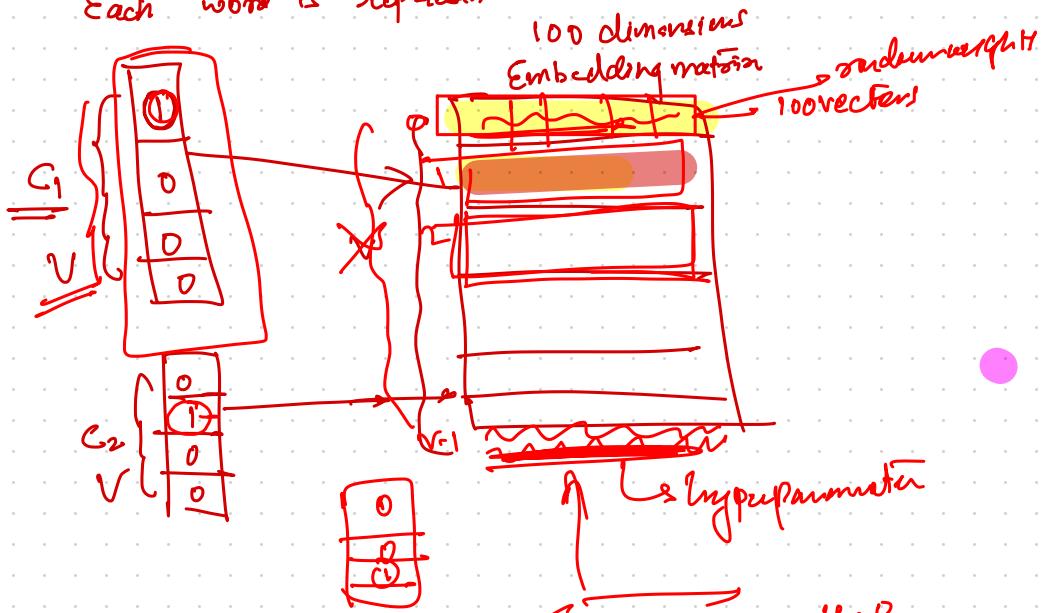
Target word
close
was
excellent

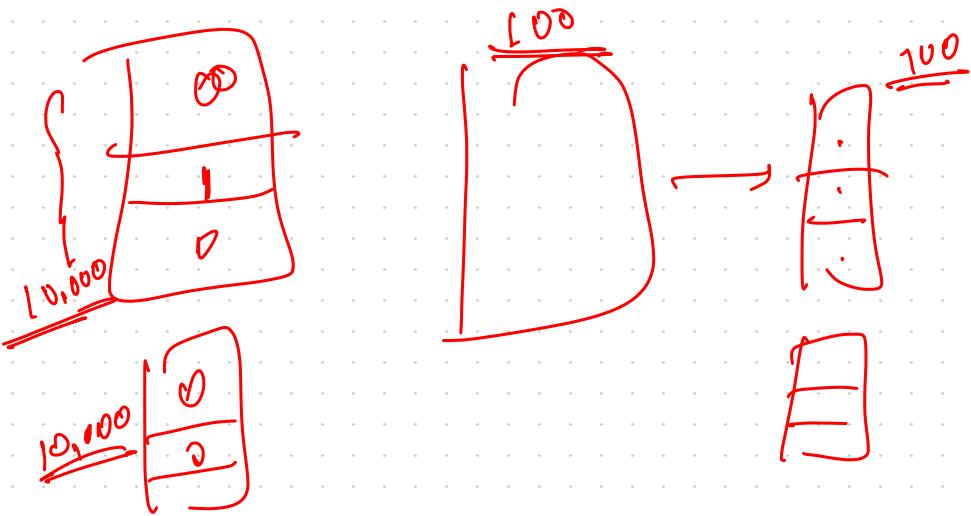
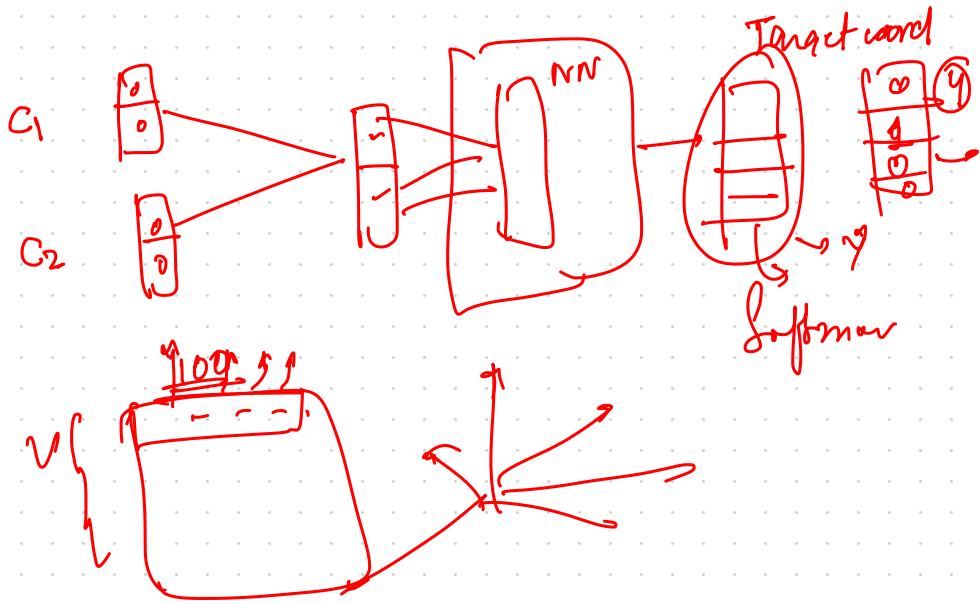
Context window
 $= 1$
Context window
 ≤ 2

close

The

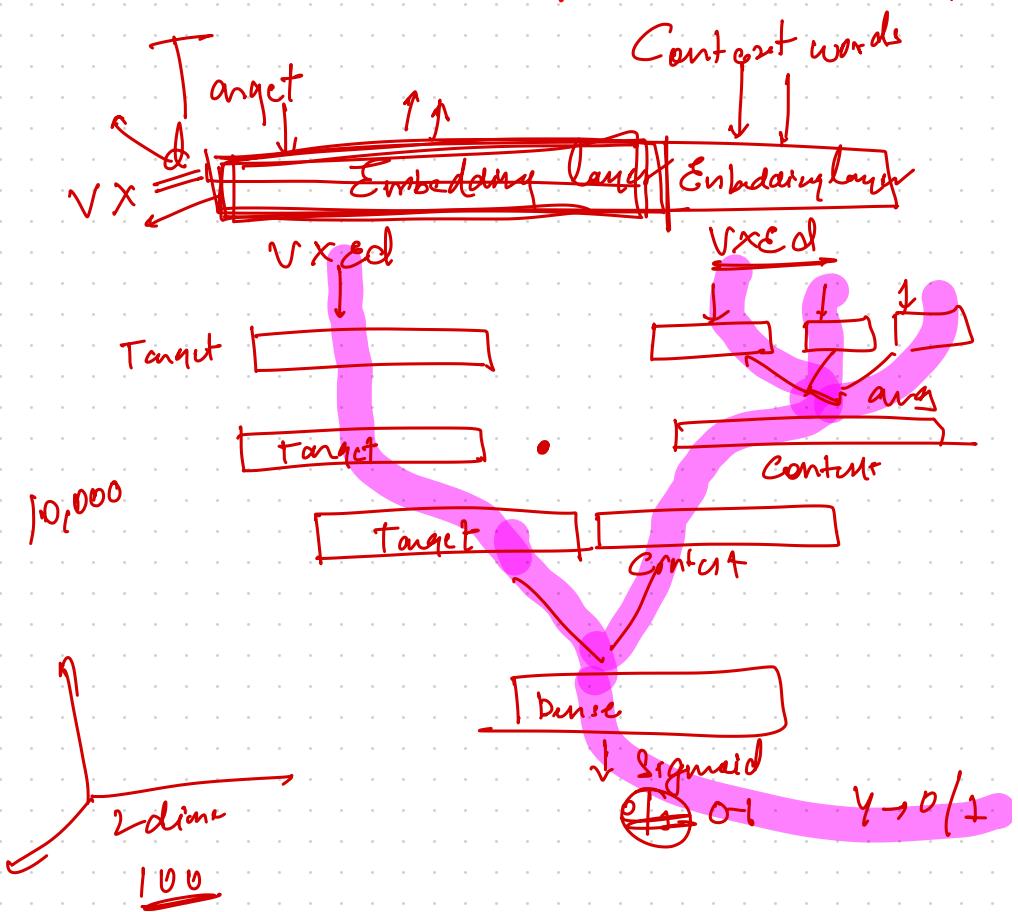
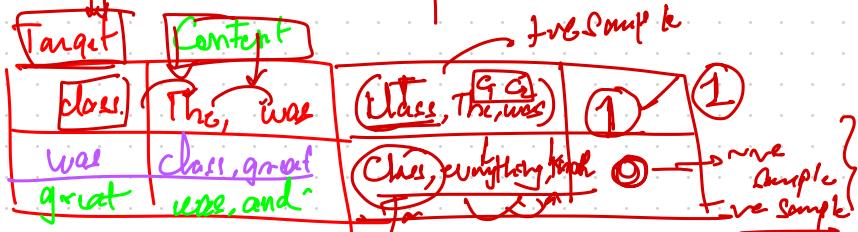
Each word is represented as a one-hot encoded vector

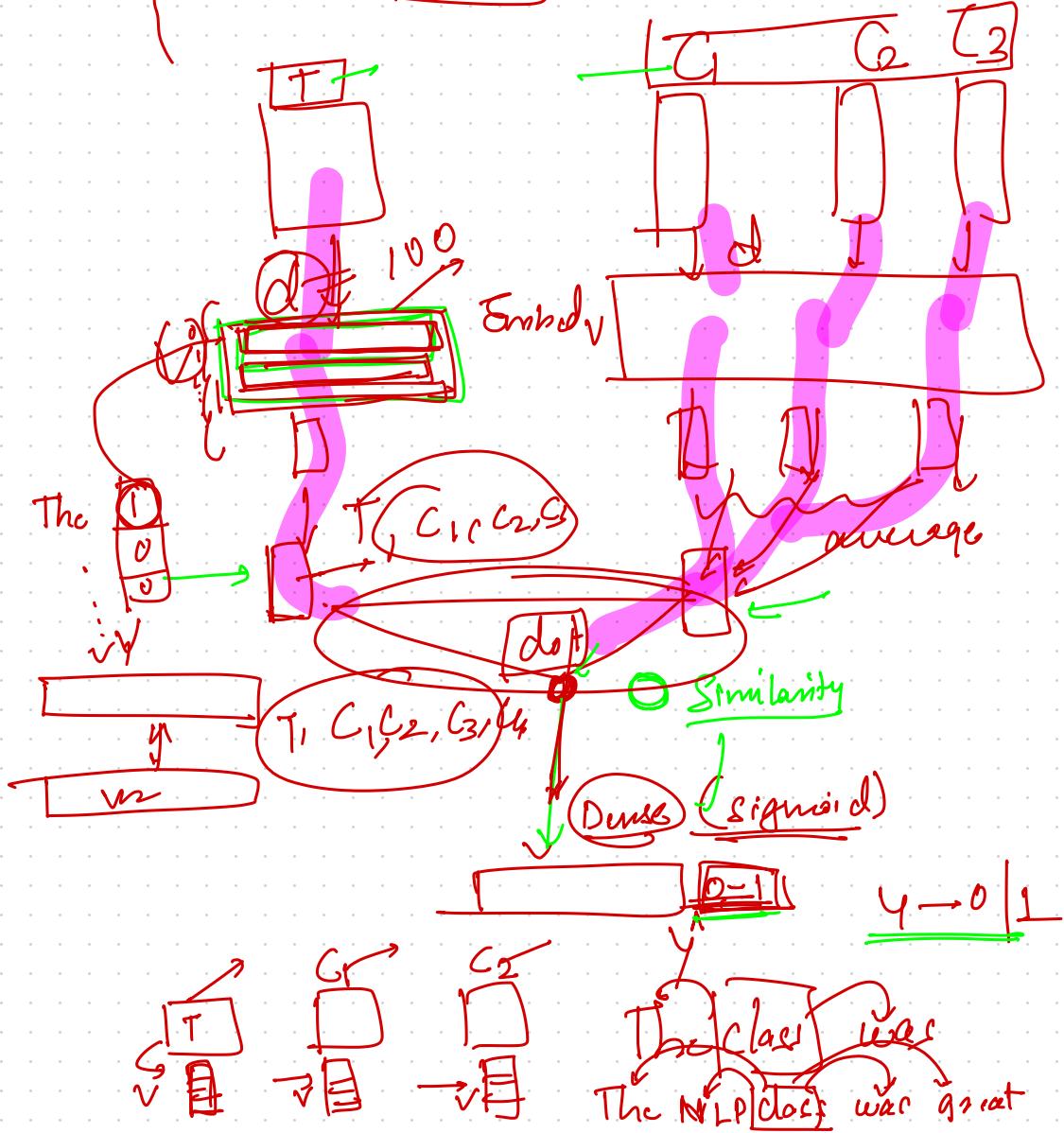
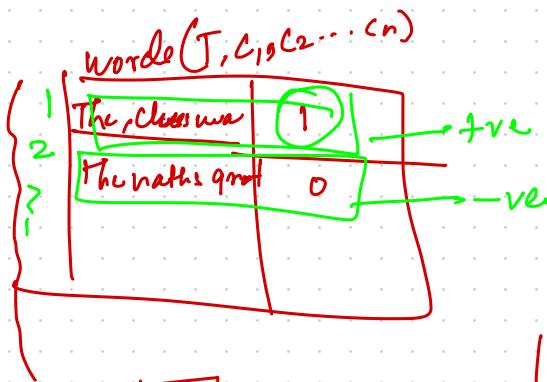




Skip-gram model

The class was great and we got to know new concepts





Target

Class

Content words

The NLP was great

v

v

v

v

v

(SxV)

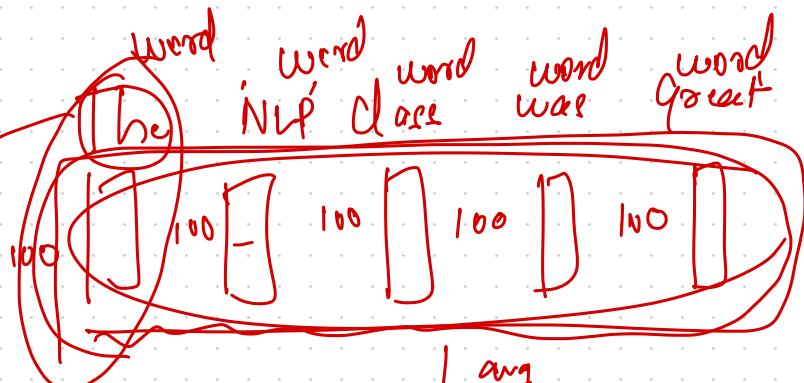


The class was great

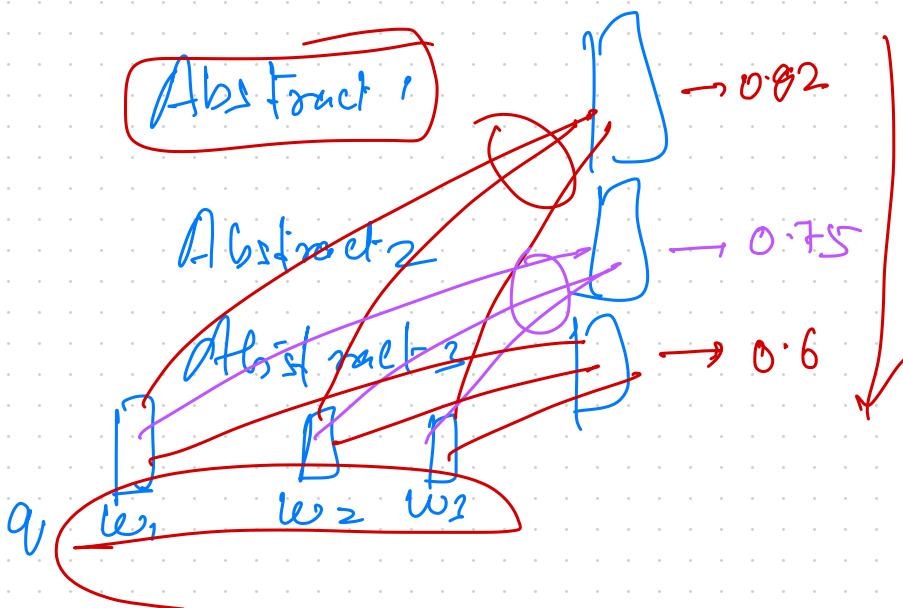
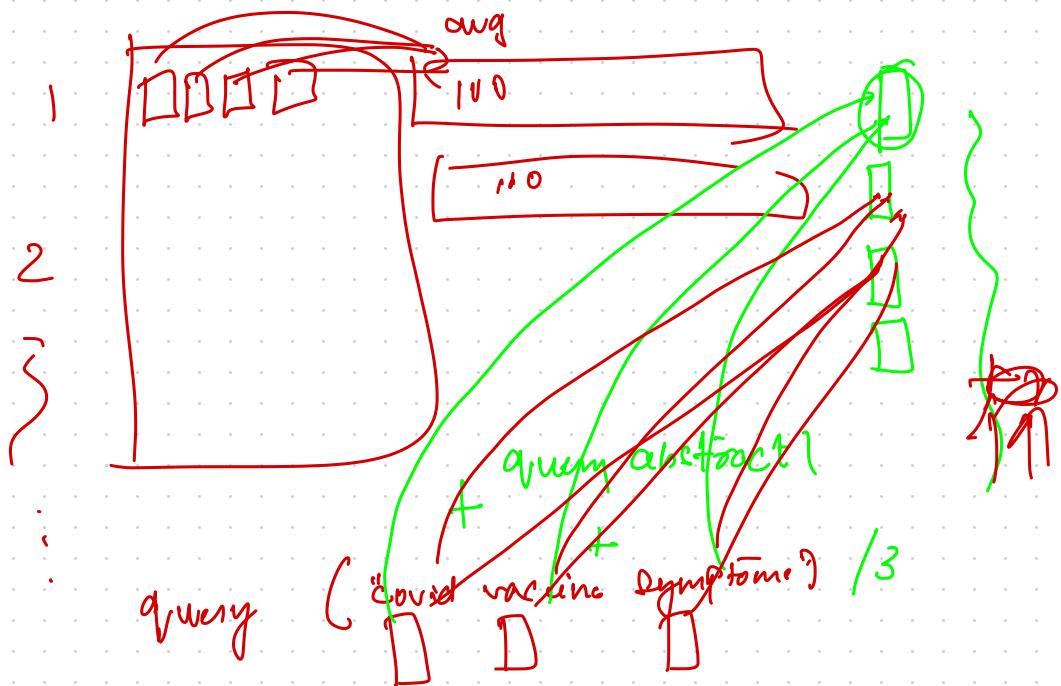
The session was great

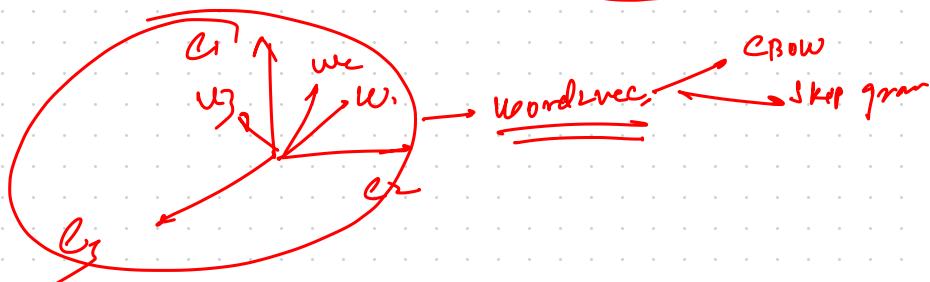
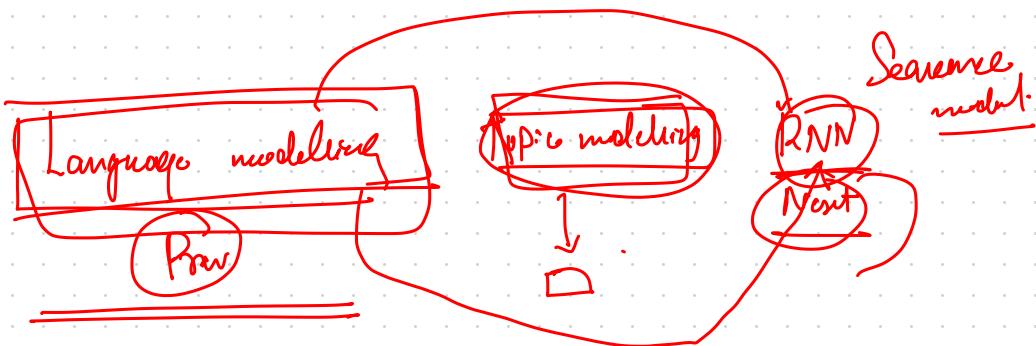
"The class was great."
The NLP is interesting

[the, class, was, great], [the, NLP, is, interesting]

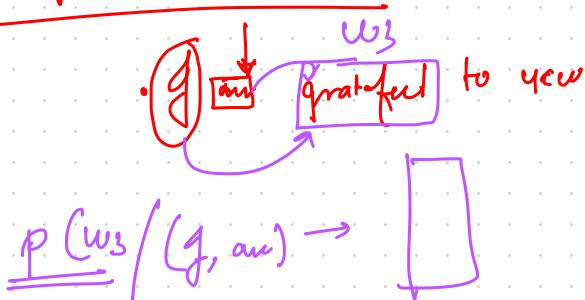


model. wv ['the'] 100 → $\sum_{i=1}^n \text{word}_i$ → avg → Representative vector for the sentence

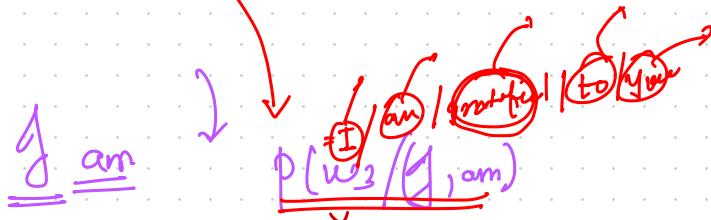




language model



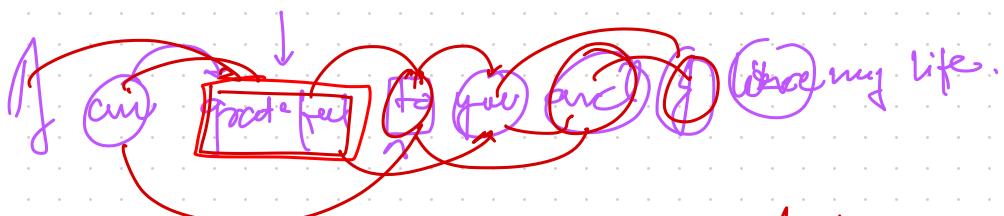
I am grateful to you



I am thankful for this life

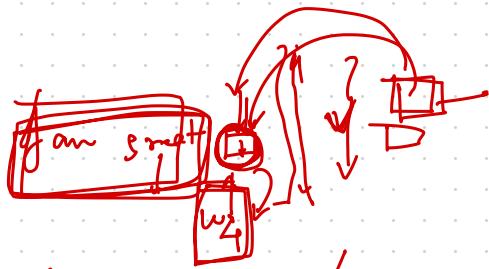


Marcovian A sequence



$K=1$ $K=2$

If a particular word at a position is dependent on previous two words



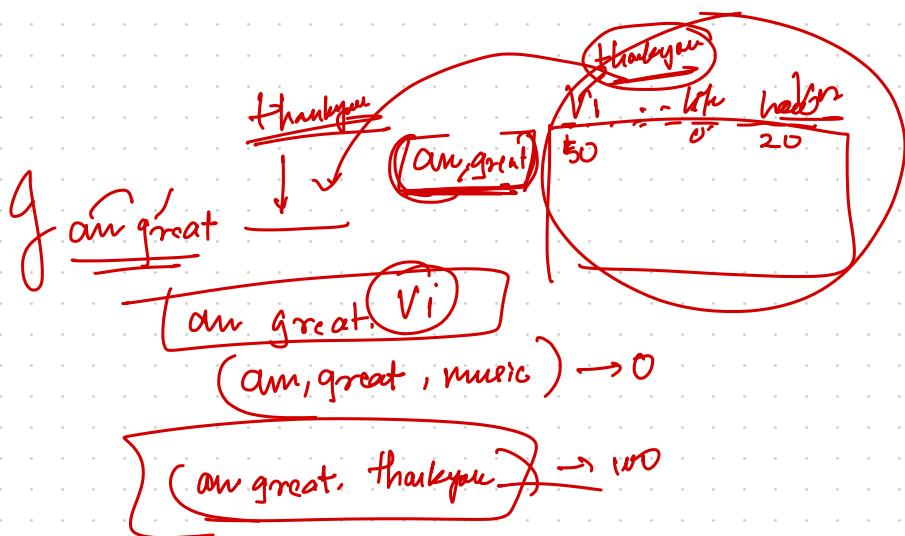
$$P \left(\text{---} \right) / \underline{(G, \text{am}, \text{great})}$$

$$K=0 \quad \sim P \left(\text{---} \right) \text{ Unigram}$$

$$\boxed{K=1} \quad \sim P \left(\underline{w_{43}} \text{ } Vi \text{ } \right) / \underline{\text{great}} \text{ Bigram} \rightarrow$$

$$K=2 \quad \sim P \left(w_4 = Vi / \underline{\text{great}}, \underline{\text{am}} \right) \text{ Trigram}$$

$$K=3 \quad \sim P \left(\underline{w_4: Vi} / \underline{\text{great}}, \underline{\text{am}}, \underline{f} \right) \text{ Quadgram}$$



The image shows handwritten notes from a child's notebook. On the left, there is a drawing of a television set with a large 'X' through it, followed by the word 'bad'. To the right is a circle containing the words 'our good', with a large 'X' through the word 'good'. Below this circle is a smaller circle with the number '2' inside. Next is the phrase 'good how' with a bracket underneath and the number '1' below it. On the far right, there is a question 'how are you?' with a large bracket underneath, followed by three numbers: '1', '1', and '1'.

good

g x
are x
good
have x
are x
your x
two x

if an — p (G/an)

P ($w_3 = v_i$ / (Jan))

I want good people to love my country. Jhore India

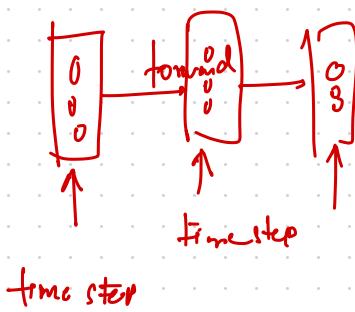
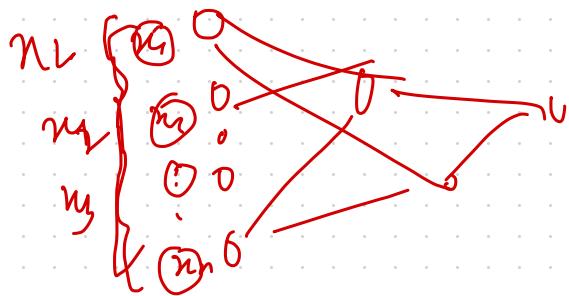
K gran

P C

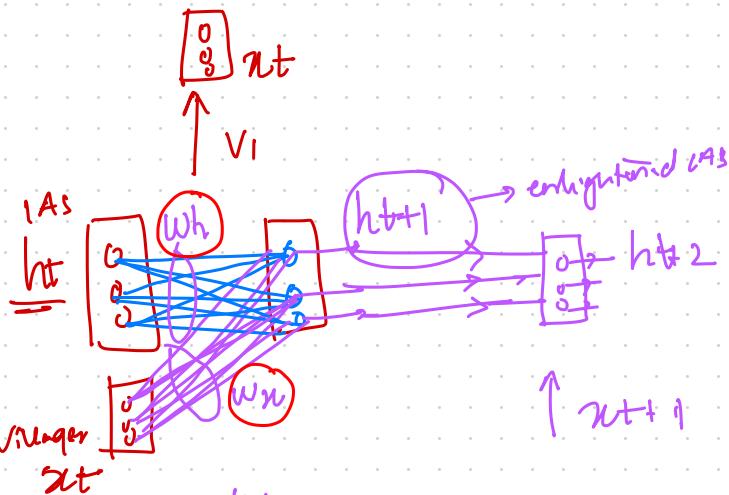
100

Recurrent Neural Networks

9:51 → 9:55 10:00 a.m.



1AS



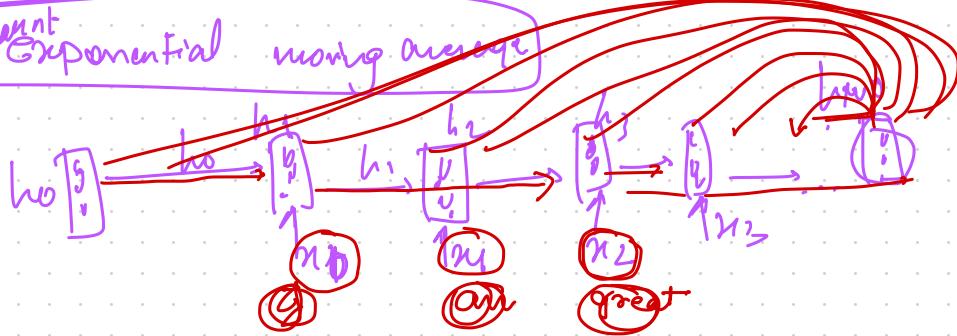
$$\underline{h_{t+1}} = \tanh(\underline{W_h} \times \underline{h_t} + \underline{W_x} \underline{x_{t+1}} + b)$$

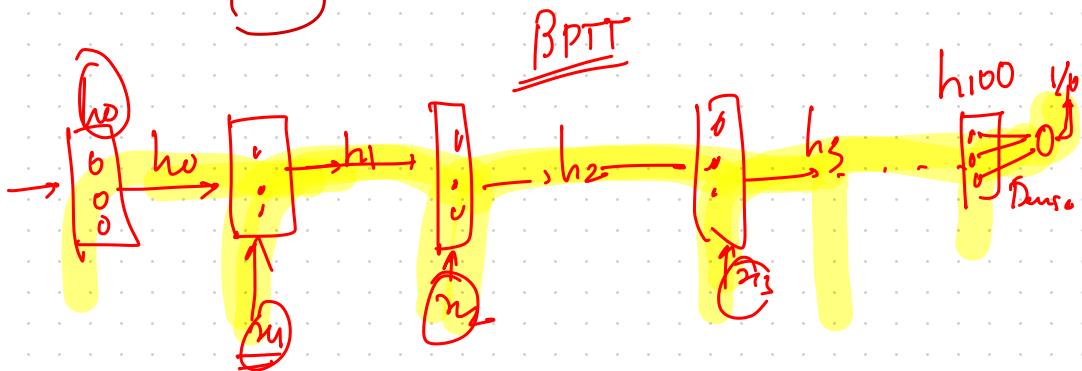
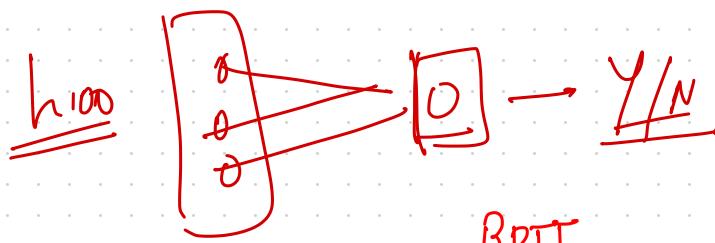
$$\underline{h_{t+2}} = \tanh(\underline{W_h} \times \underline{h_{t+1}} + \underline{W_x} \underline{x_{t+2}} + b)$$

$$\underline{h_{t+3}} = \tanh(\underline{W_h} \times \underline{h_{t+2}} + \underline{W_x} \underline{x_{t+3}} + b)$$

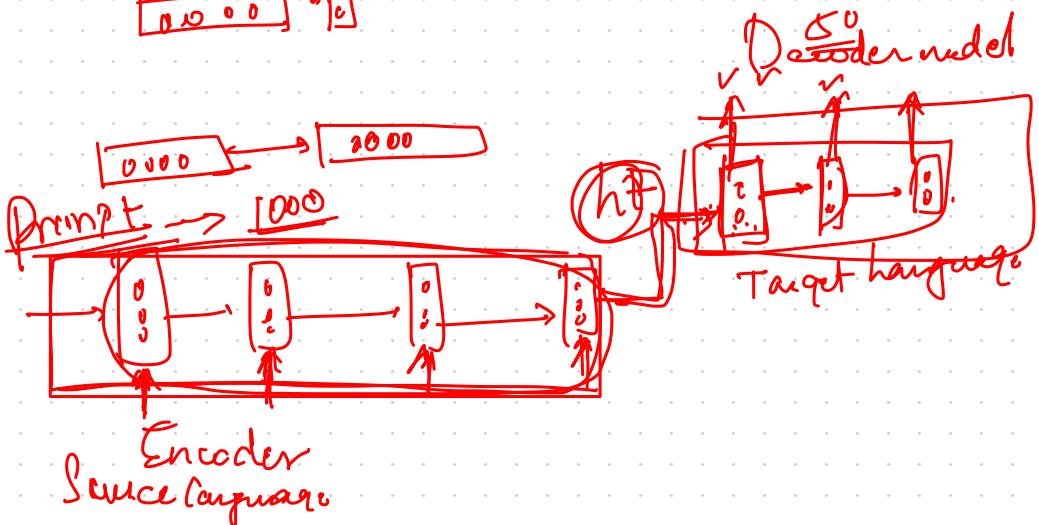
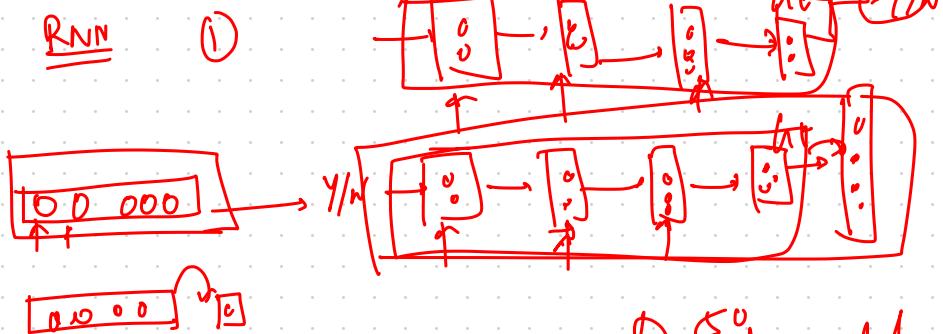
$$\underline{h_{t+4}} = \tanh(\underline{W_h} \times \underline{h_{t+3}} + \underline{W_x} \underline{x_{t+4}} + b)$$

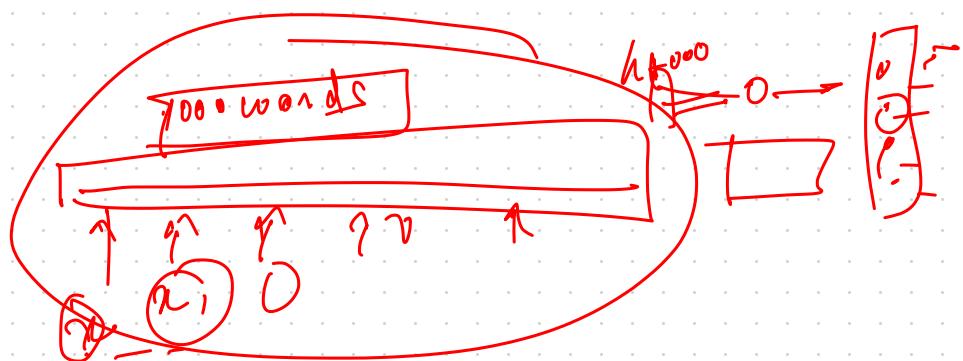
Non-linear learnt Exponential moving average



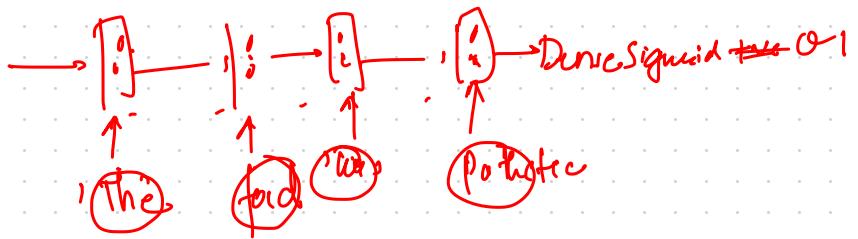


RNN





NLP



A hand-drawn diagram of a trapezoidal element with vertices labeled x , y , and z . The top edge has a length of 10. The left edge is labeled 1000×10 . The right edge is labeled 1000. The bottom edge is labeled 1000. The interior contains two sets of red numbers: 123400000 and 153678900.

<u>1000</u>	<u>π</u>	year-length $\rightarrow 10$
①	tiny food was plentiful	-ve
②	The audience was enthusiastic, definitely a <u>plus</u>	
3		

