

---

---

---

---

---



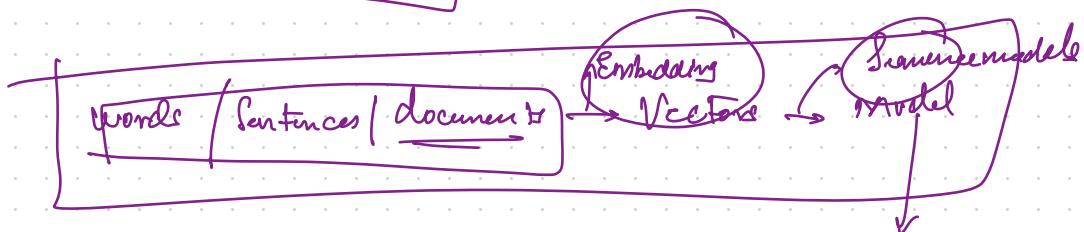
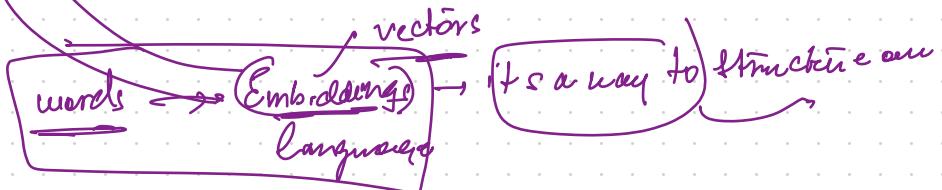
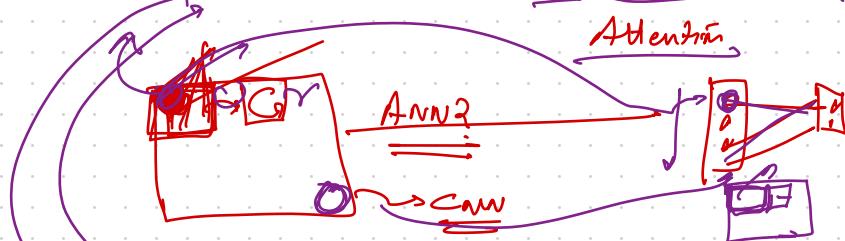
# Natural language Processing

① NLP problems are easel.

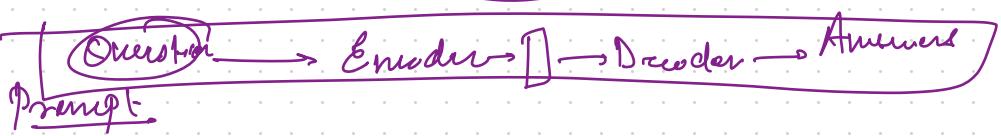
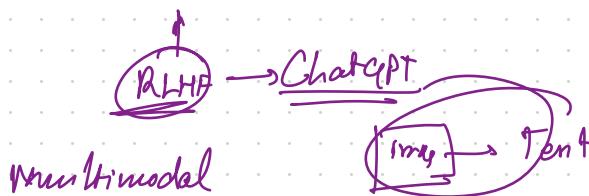
② How do we structure words / sentence

Feature engineering

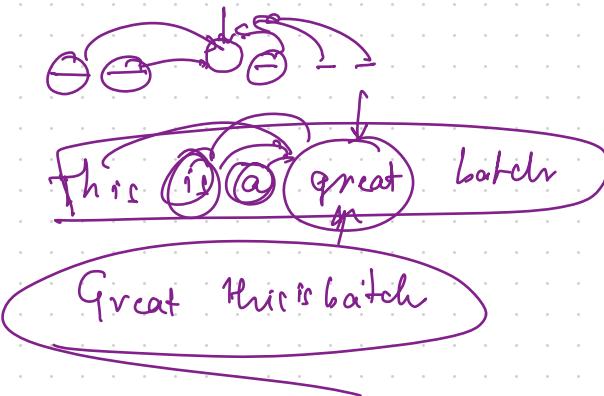
③ Same models → RNN, LSTM, GRU, Transformer



large language models



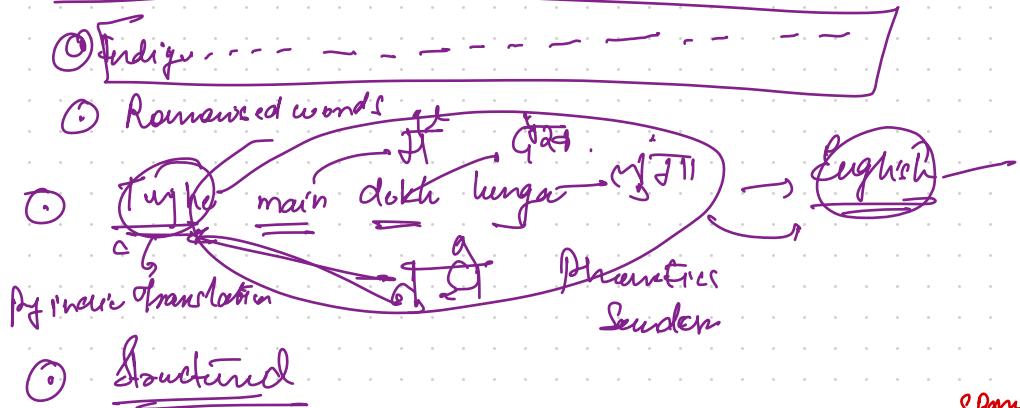
Natural Language - is a Sequence



## Natural language Processing

Understanding information from text

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



- ③ Structured

Use Cases:

Translation

Gmail auto type → auto fill

Search engines

Sentiment analysis

Small spam → Ham  
spam

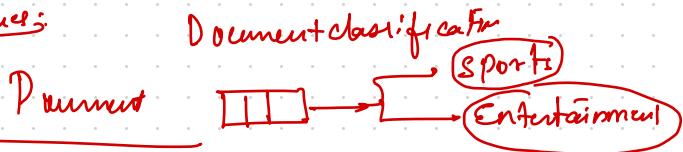
Chatbots

Text classification, Sum  
NLP, 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

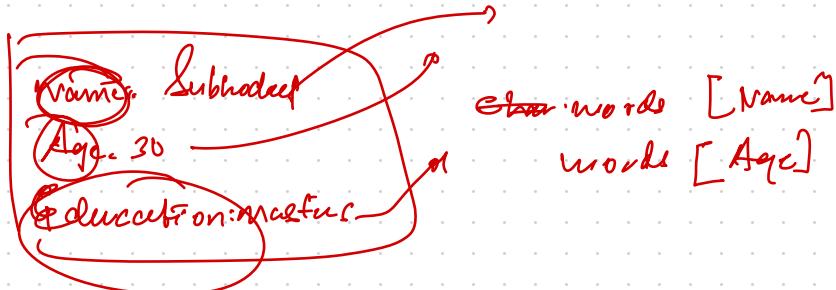
Ticket

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

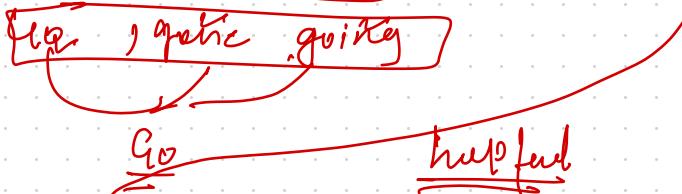
✓ Frequency based  
✓ Count of the words  
in the sentence

Context?

Randomness based:  
vectors generated  
through LNN.

If am going home, home is where heart is

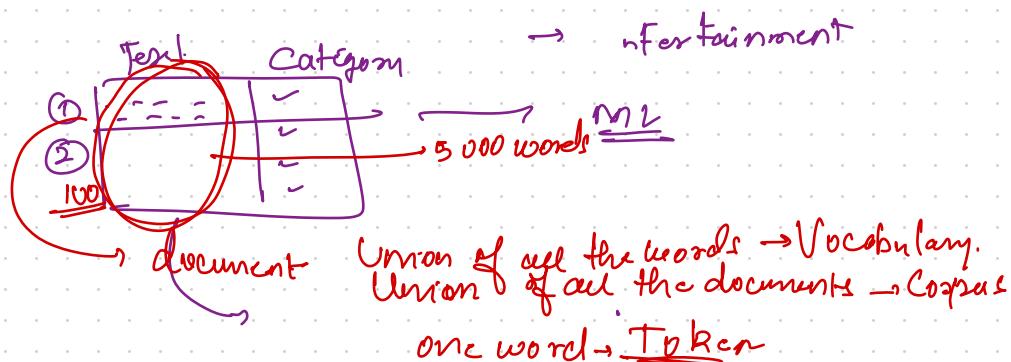
g:1, am:1, going:1, home:2, ts:2  
where:1, heart:1



10:32 → 10:40pm

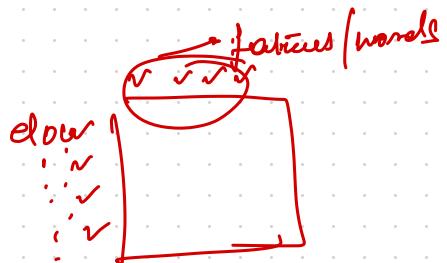
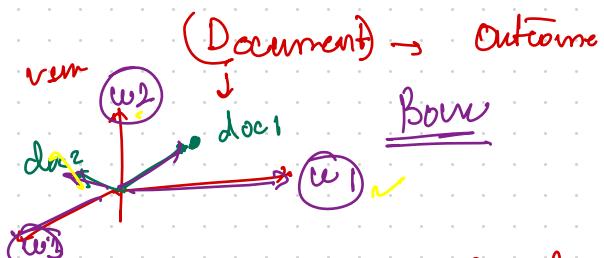


How to represent sentences / documents in a vectorised form using frequency embeddings.



## Bag of words model

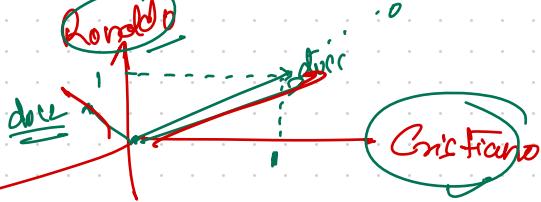
Document classification



- ① Cristiano Ronaldo got transferred Cristiano's great  
→ Sports

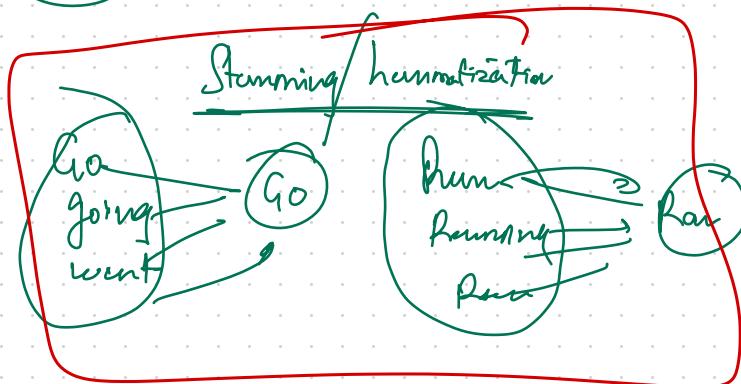
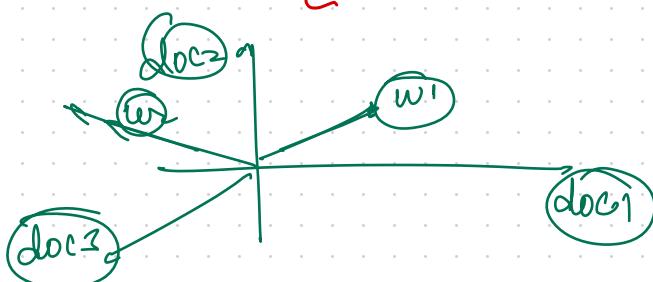
- ② Puerto Rican performs in Cannes.

	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





## Lernmaterialien

Verb

Running

Run

Going to go

Spreading ✓  
Specific ✓  
Concepts ✓  
Copy → Oral

Diagram illustrating Verb Phrases and Dependencies:

**Verb Phrase Analysis:**

- Root node: "will make" (verb phrase)
- Children: "will" and "make"
- "make" has children: "visit my home" and "and"
- "visit my home" has children: "visit" and "my home"
- "visit" has children: "I" and "go"
- "go" has children: "to" and "my home"
- "my home" has children: "my" and "home"
- "my" has children: "I" and "ac"
- "home" has children: "and" and "will make" (referred to as a recursive node)

**Dependency Parse:**

```

graph TD
    Root --- NP1[my own]
    Root --- NP2[is awesome]
    NP1 --- P1[is]
    NP1 --- NP3[my]
    NP1 --- NP4[home]
    NP3 --- P2[my]
    NP3 --- NP5[home]
    NP5 --- P3[and]
    NP5 --- VP[will make]
    NP4 --- P4[and]
    NP4 --- VP
    VP --- V[will]
    VP --- V2[make]
    V --- T1[to]
    V --- NP6[visit my home]
    NP6 --- V3[visit]
    NP6 --- NP7[my home]
    V3 --- T2[I]
    V3 --- NP8[go]
    NP8 --- T3[to]
    NP8 --- NP9[my home]
    NP9 --- T4[visit]
    NP9 --- NP10[my home]
    T4 --- T5[visit]
    T5 --- T6[my]
    T6 --- T7[home]
    T7 --- T8[and]
    T8 --- VP2[will make]
    VP2 --- V4[will]
    VP2 --- V5[make]
    V4 --- T9[visit]
    V5 --- T10[my]
    V5 --- T11[home]
  
```

The diagram shows the hierarchical structure of the sentence "I am going to visit my home and will make" and its dependencies. The root node is "will make". It branches into "will" and "make". "make" branches into "visit my home" and "and". "visit my home" branches into "visit" and "my home". "visit" branches into "I" and "go". "go" branches into "to" and "my home". "my home" branches into "my" and "home". "my" branches into "I" and "ac". "home" branches into "and" and "will make" (a recursive node). The dependency parse provides a detailed breakdown of these relationships.

"our great at public Speaking"

Stemming (word)

lemmatization

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

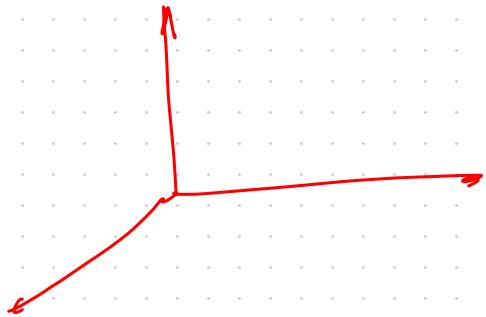
game awesome

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

$$\text{pos} = 130$$

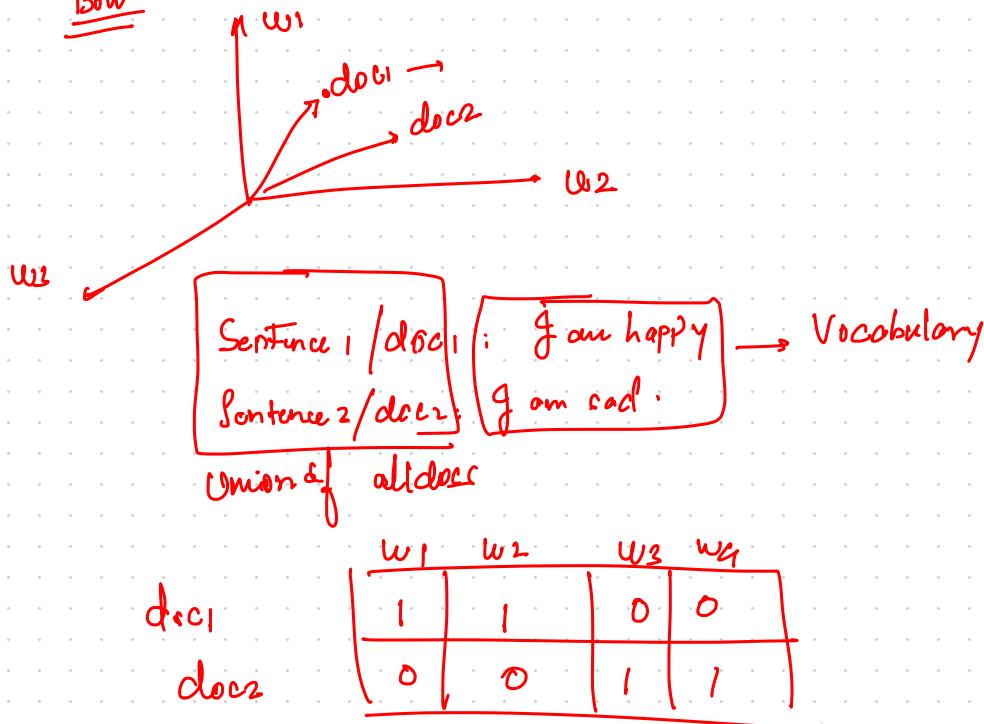
$$\text{neg} = 70$$

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

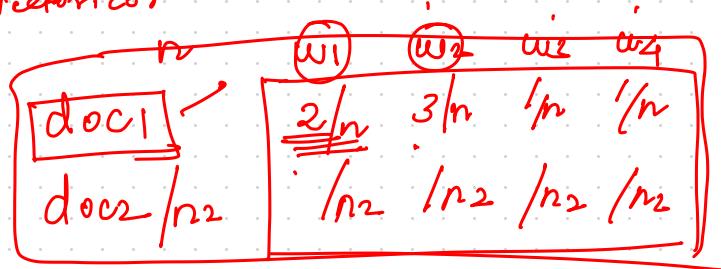


NLP  $\rightarrow$  words / docs into vector space models

BOW



Count vectorizer



- ① do not maintain any order

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

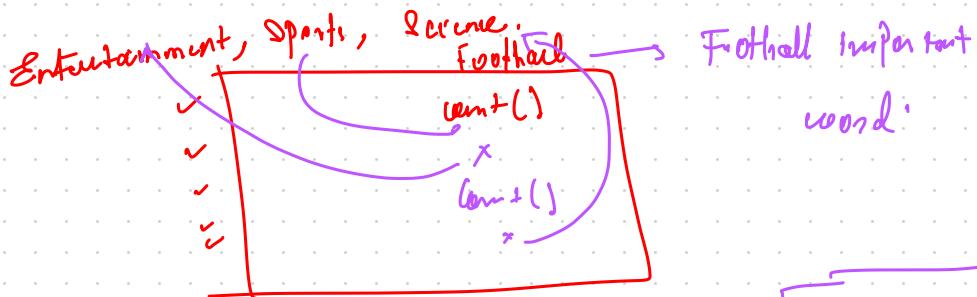
doc1 

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

frequency based embeddings don't capture the context

# Term Frequency vs Document Frequency

Football



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

$$\text{tf}(\text{Term frequency}) \times \text{idf} = \frac{f(w_i, \text{doc}_i)}{\text{Total number of tokens in doc}_i} \times \log\left(\frac{N}{n}\right)$$

where  $w_1, w_2, w_3, \dots, w_m$  are words.

$N = 100$

$n$  is the 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

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

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

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

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

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

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

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

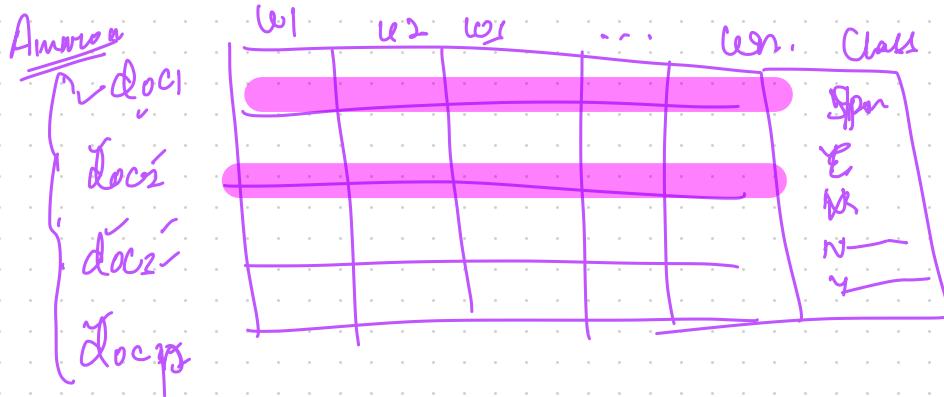
$\log \left( \frac{N+1}{n} \right) \approx 1$

TF-IDF ( $t_i$ )<sub>doc1</sub>

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

$\overbrace{N}^{\rightarrow}$  total number of doc in corpus.

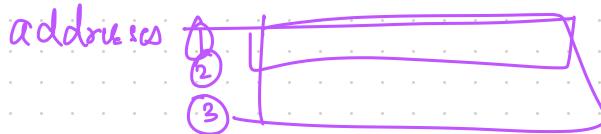
$\overbrace{n}^{\rightarrow}$  total number of doc in which term occurs.



Find out Similar Documents

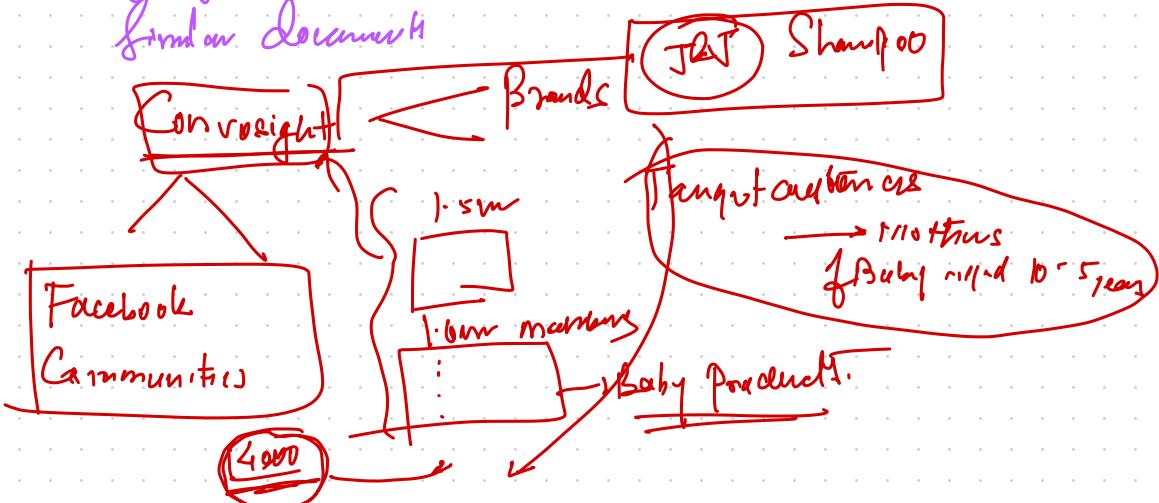
doc1 (Product)      doc3 (Product)

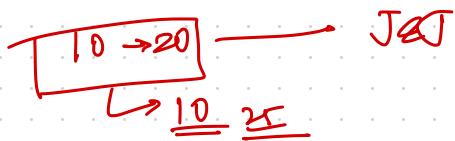
doc2 (Product) : :



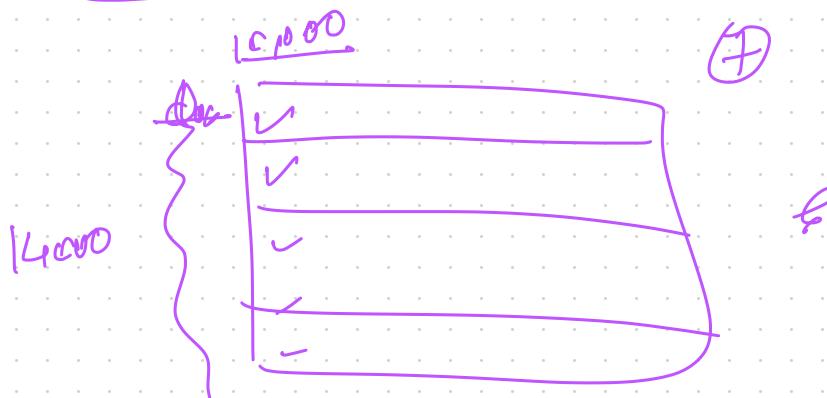
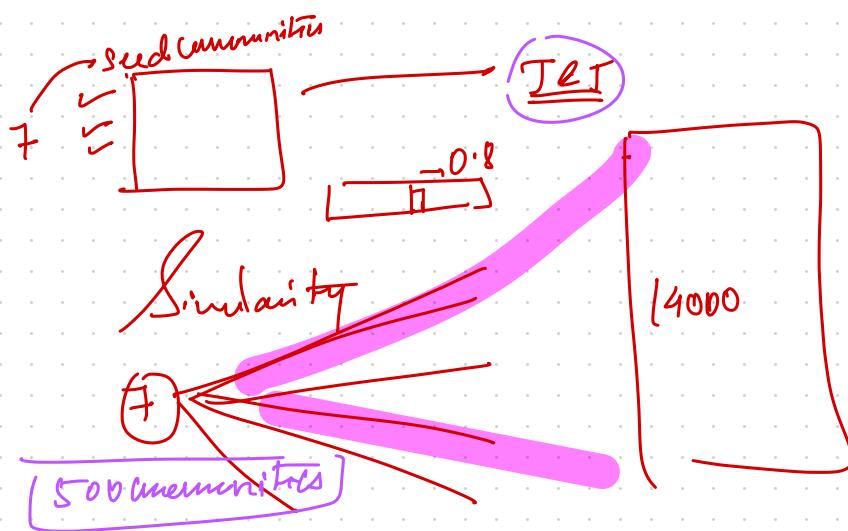
Using document based embedding of word2vec

Find our documents





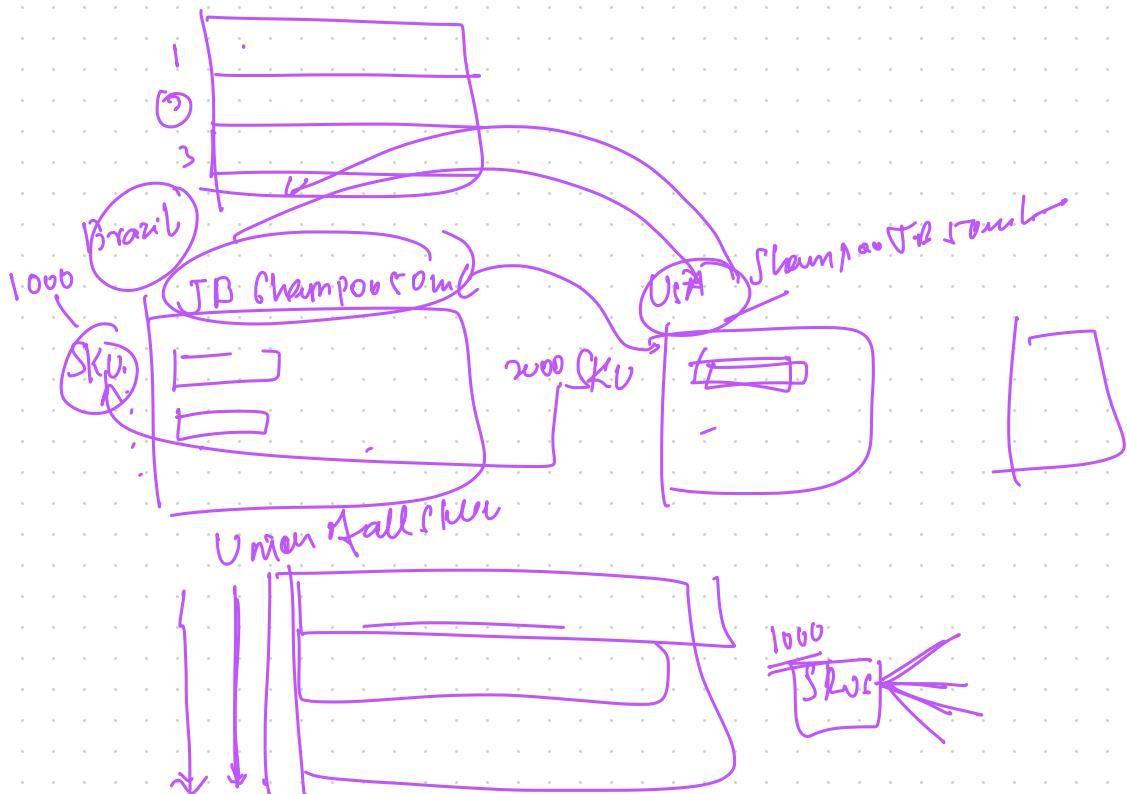
1000 users  
↳ 1.6 billion conversations

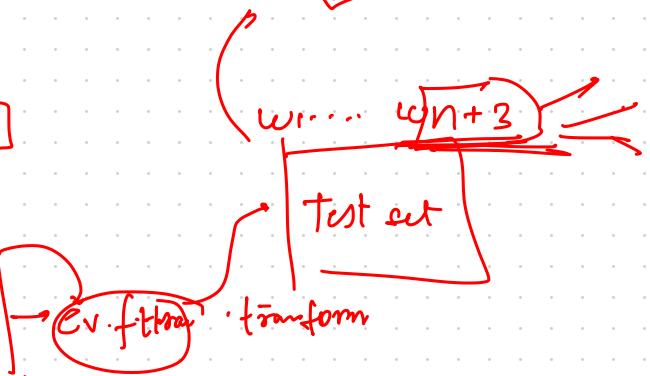
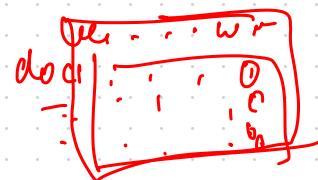
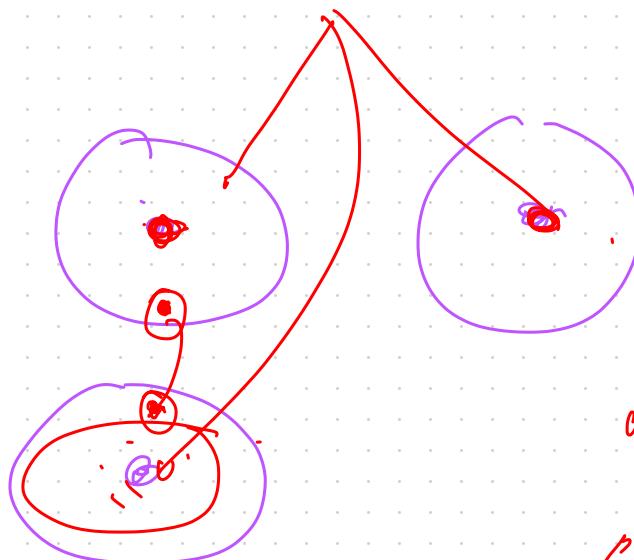


fan  happy to teach you

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

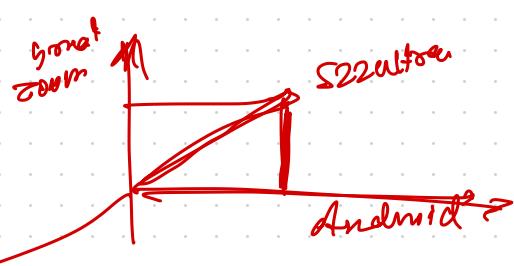
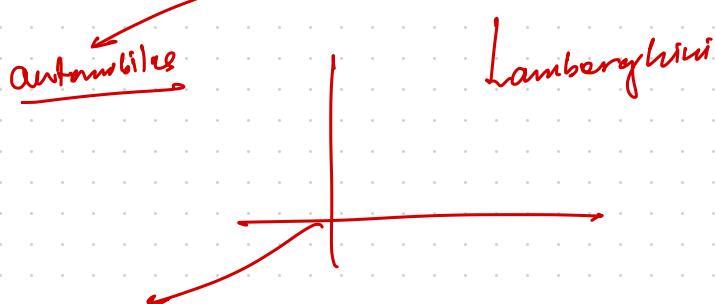
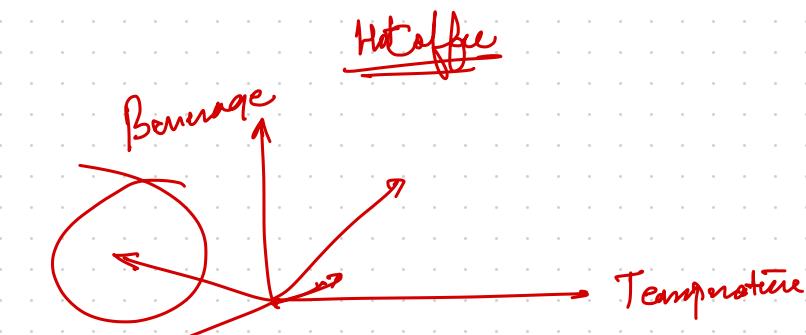
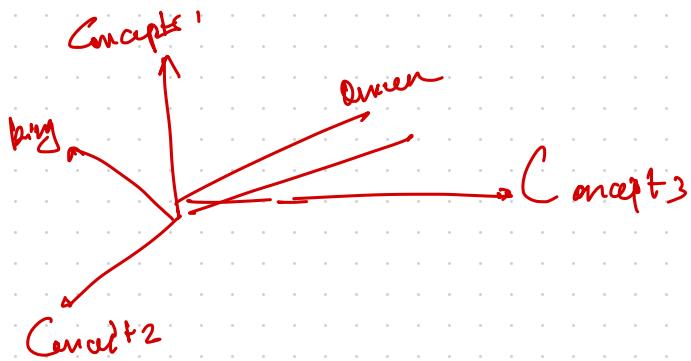
docs

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  $\propto \propto x$

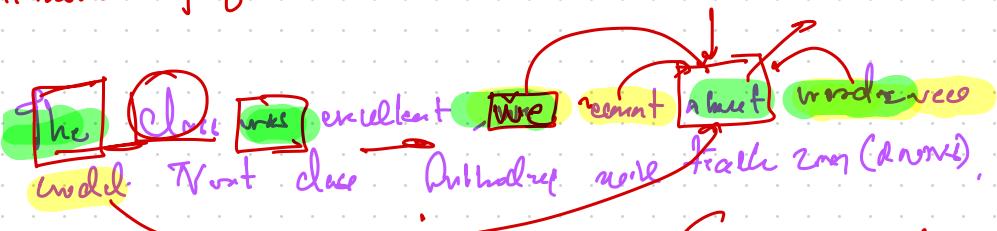
Embedding from sequences models  $\rightarrow$  (Supervised way)

Lm. Seqs ( $p_{im}, h_{im}$ )

Word2vec

cbow  
continuous Bag of words

Skip gram model



Context words  
The, was  
close, excellent  
was, we

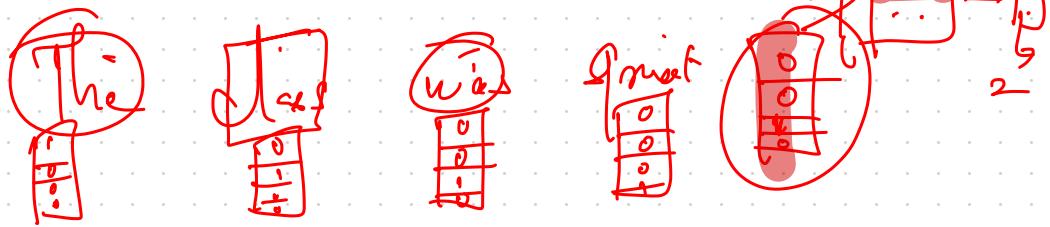
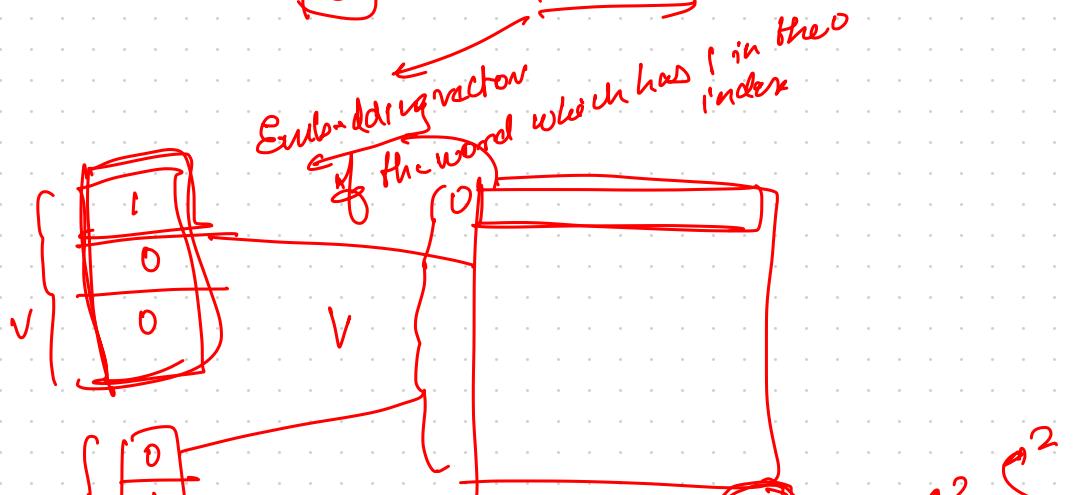
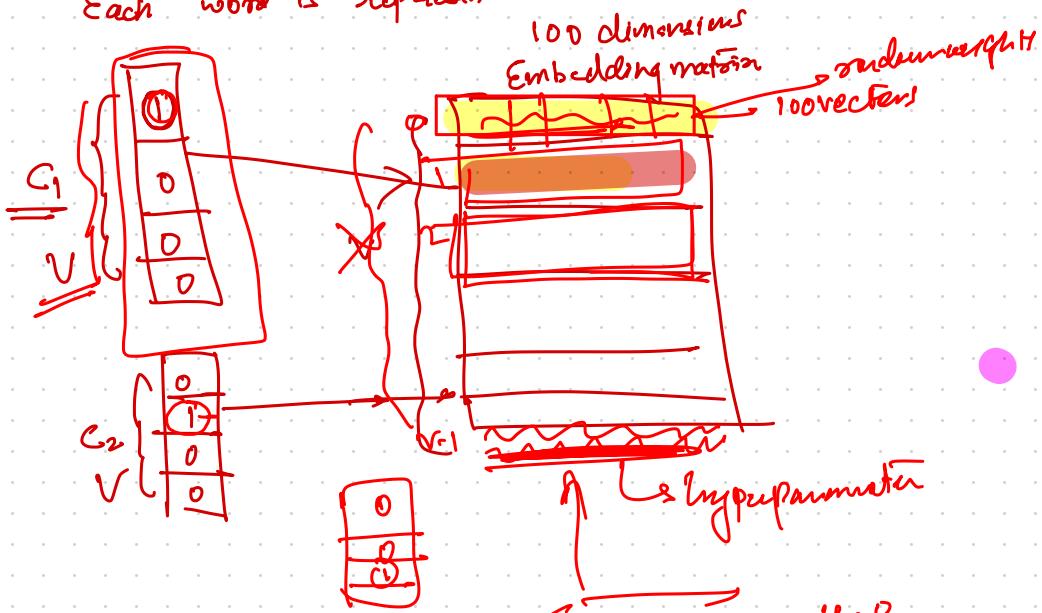
Target word  
close  
was  
excellent

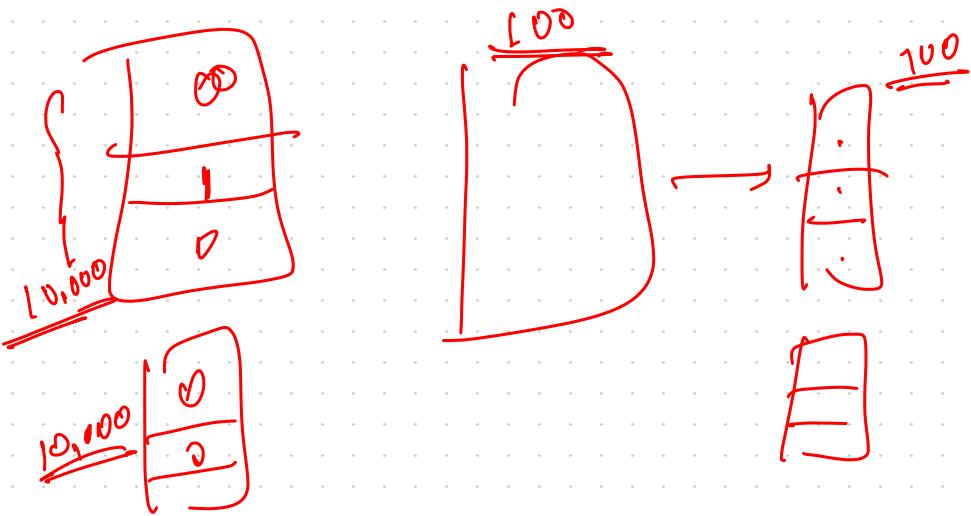
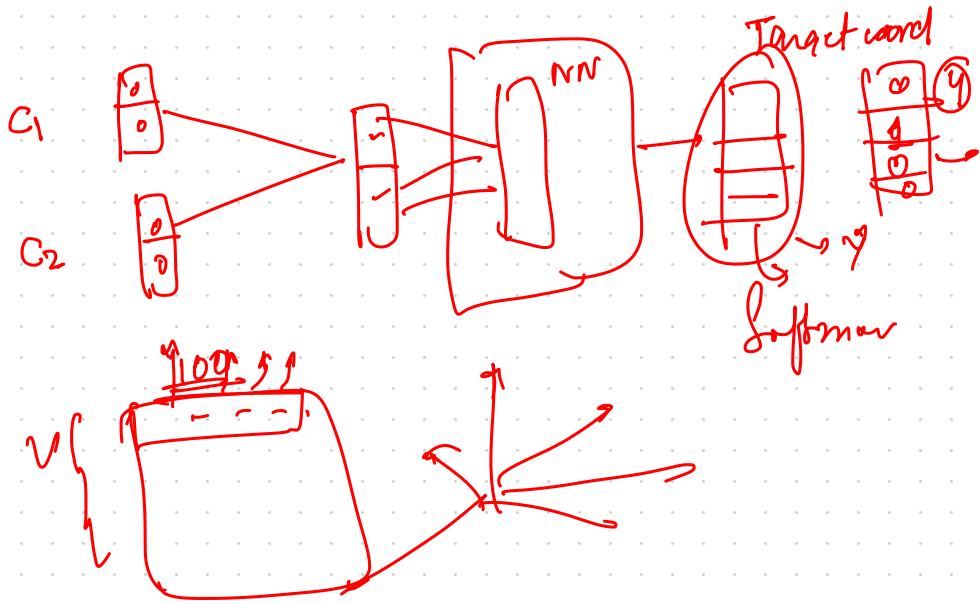
Context window  
 $= 1$   
Context window  
 $\leq 2$

close

The

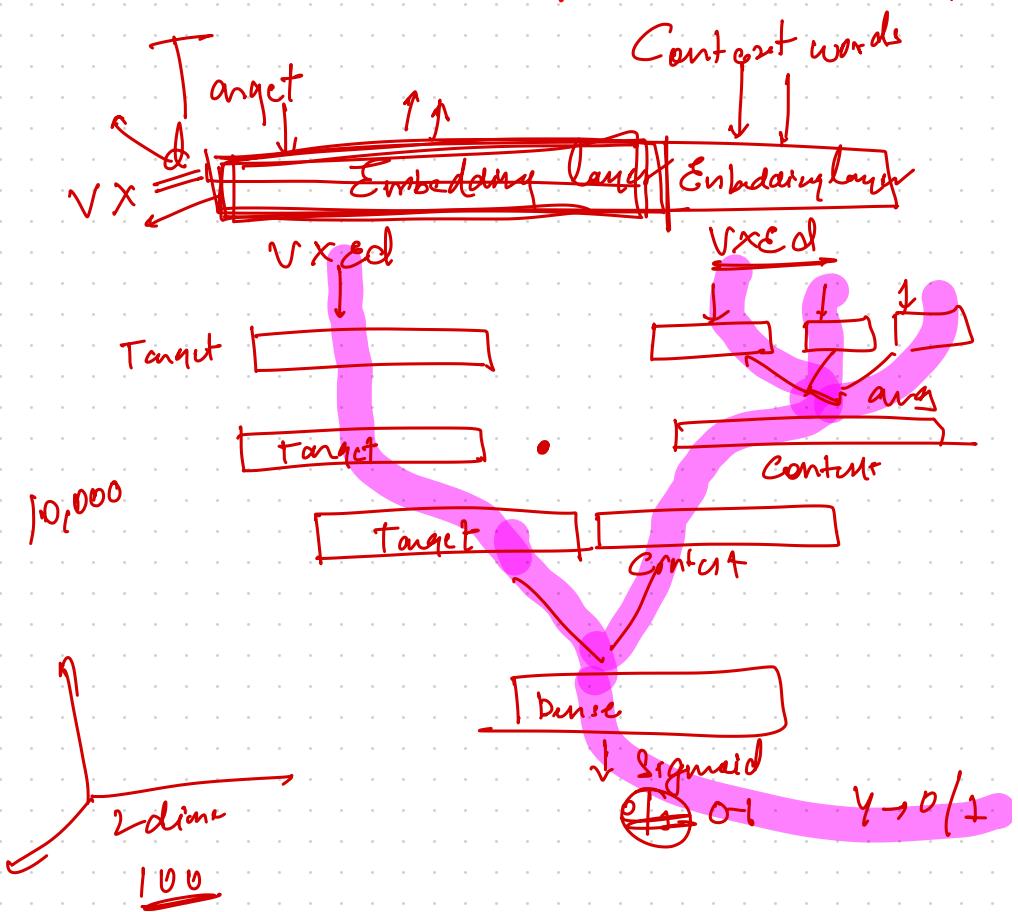
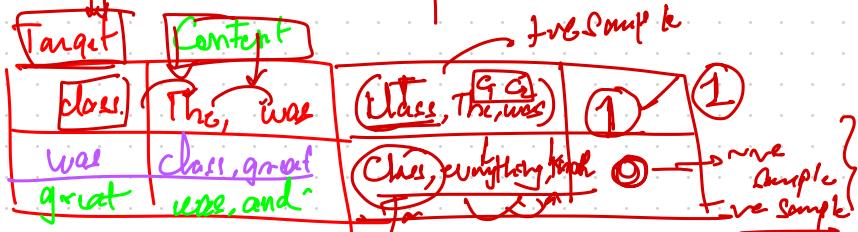
Each word is represented as a one-hot encoded vector

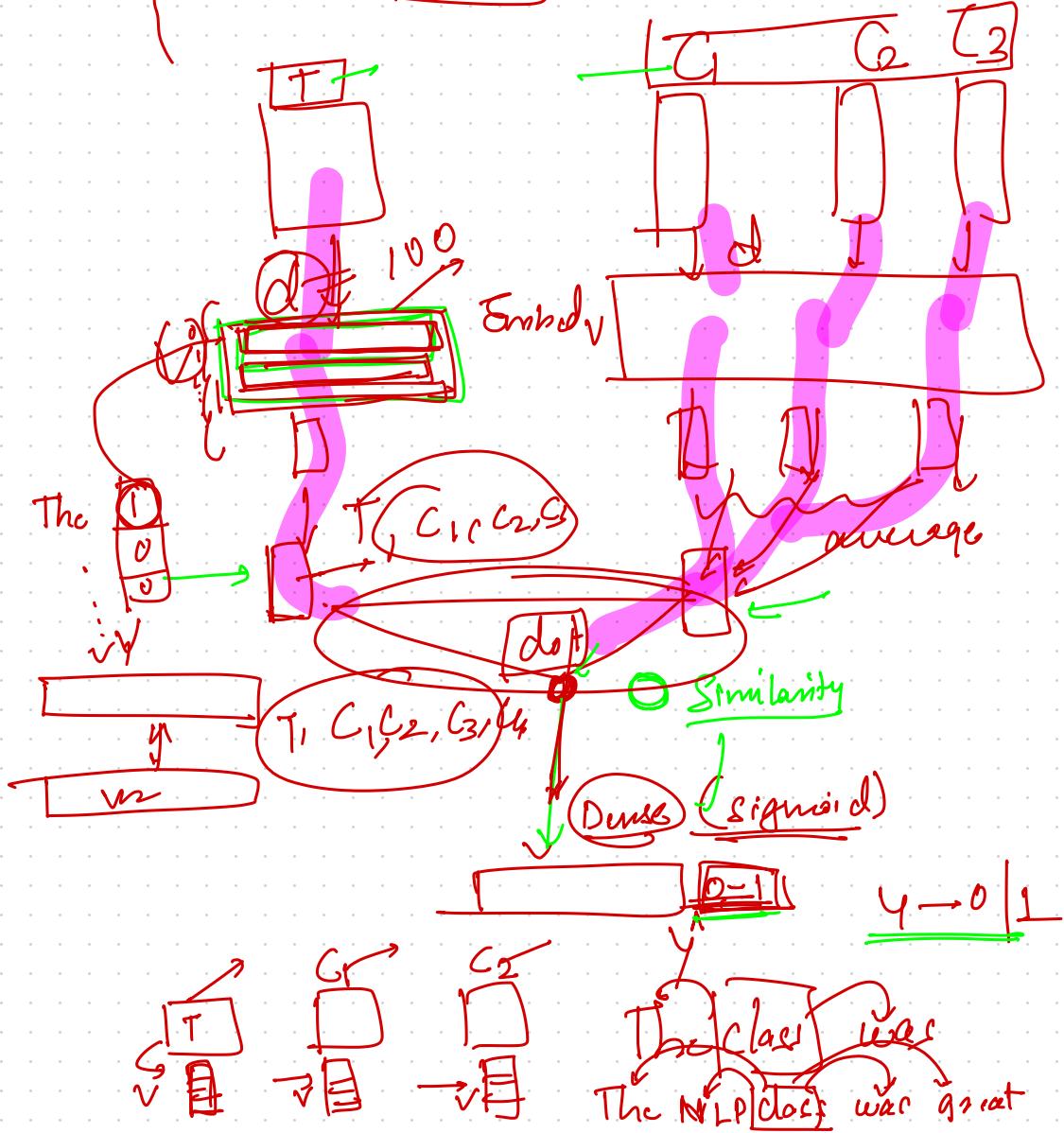
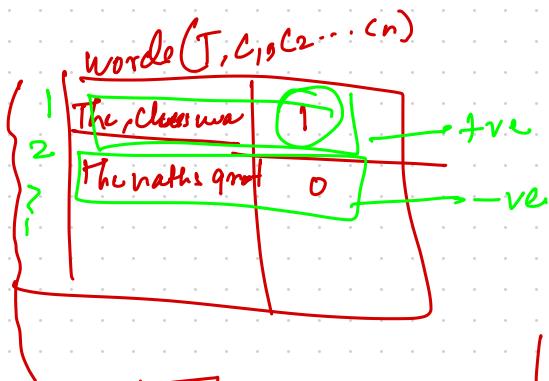




## Skip-gram model

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)

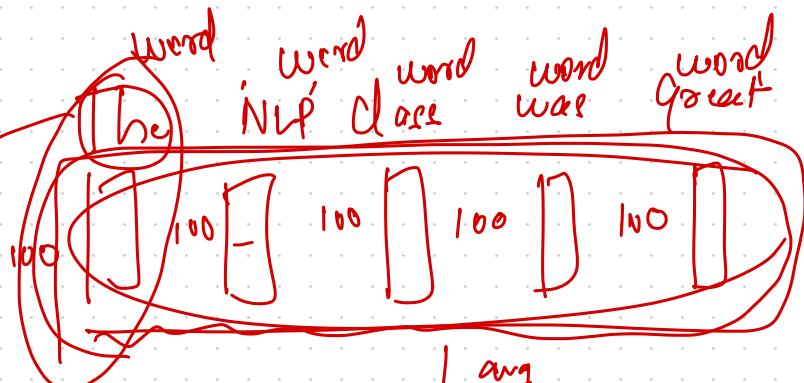


The class was great

The session was great

"The class was great."  
The NLP is interesting

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



model. wv ['the'] 100 →  $\sum$  avg  
 $\rightarrow$  100 → Representative vector for  
 the sentence

