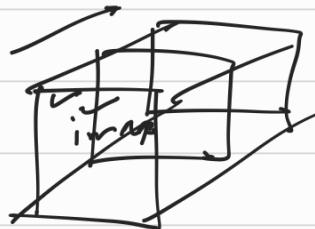


Natural Language Processing

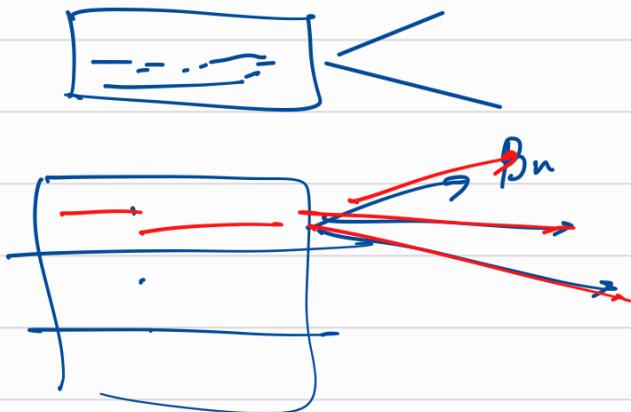
to make machines → human language.



→ don't have a universal representation in numbers.

→ modelling.

→ Classification

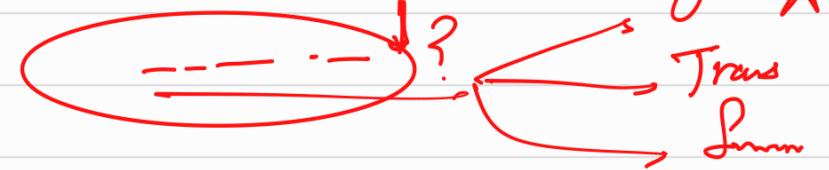


→ Named Entity

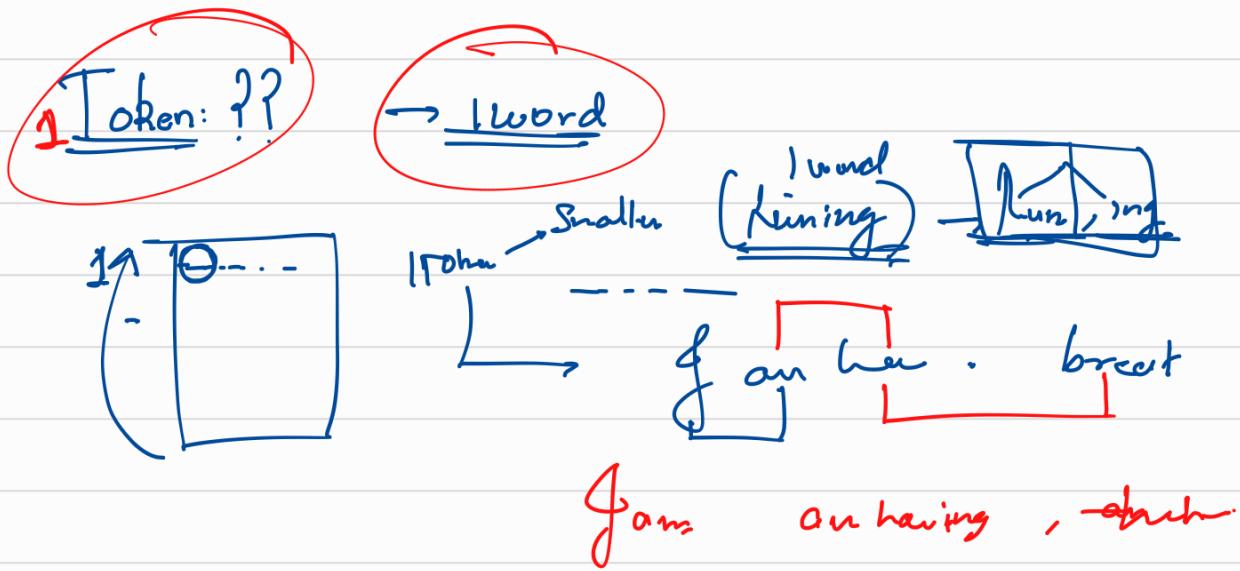
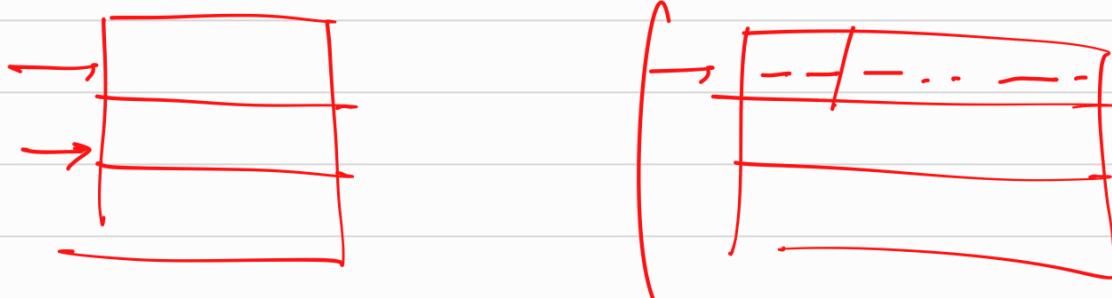
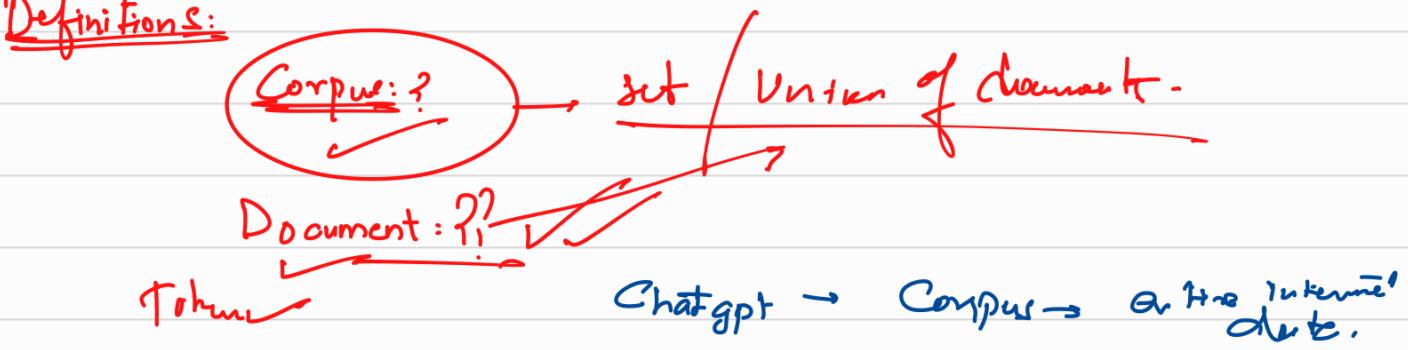
Timing, Topic

→ Have a ~~sitting~~ class at 11:00 am on Saturday about mp.

(Large) Language modeling



Definitions:



Vocabulary:

Union
Collection of Tokens

- NLP
- ① text → numbers
 - ② Model

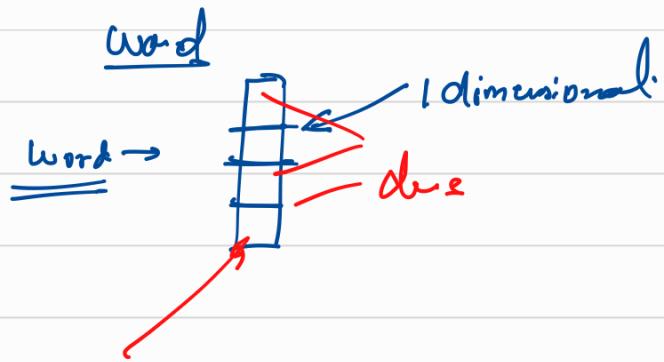
text → numbers (Embedding).

Prediction

→ (Frequency)

text → numbers

Frequency

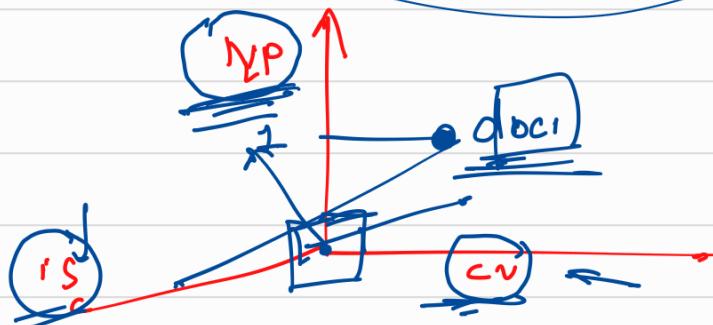
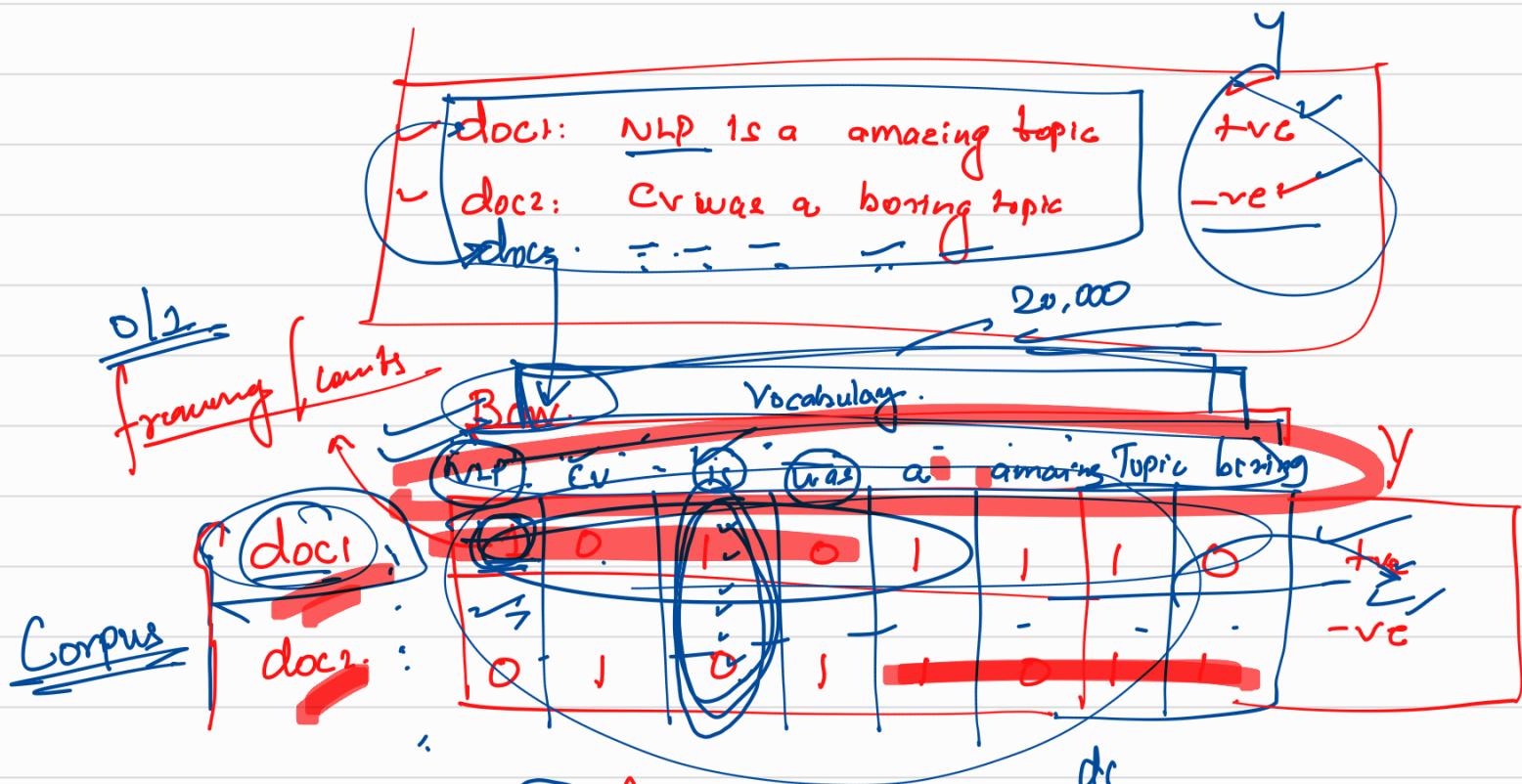


embedding



$$\text{an} \sim \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

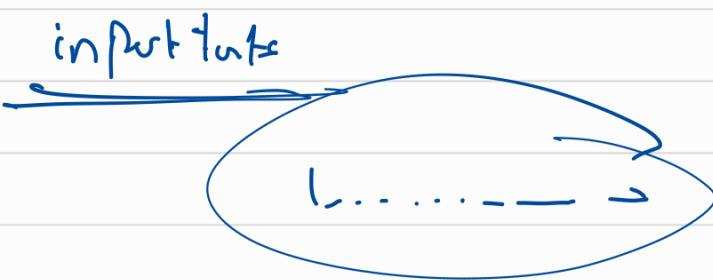
vectors



High dimensions → Sparcity

Ignores → order

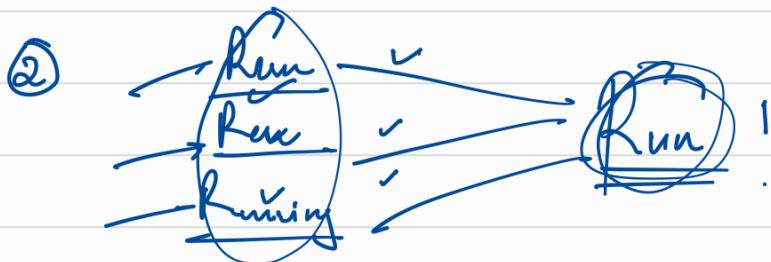
No Contextual Information

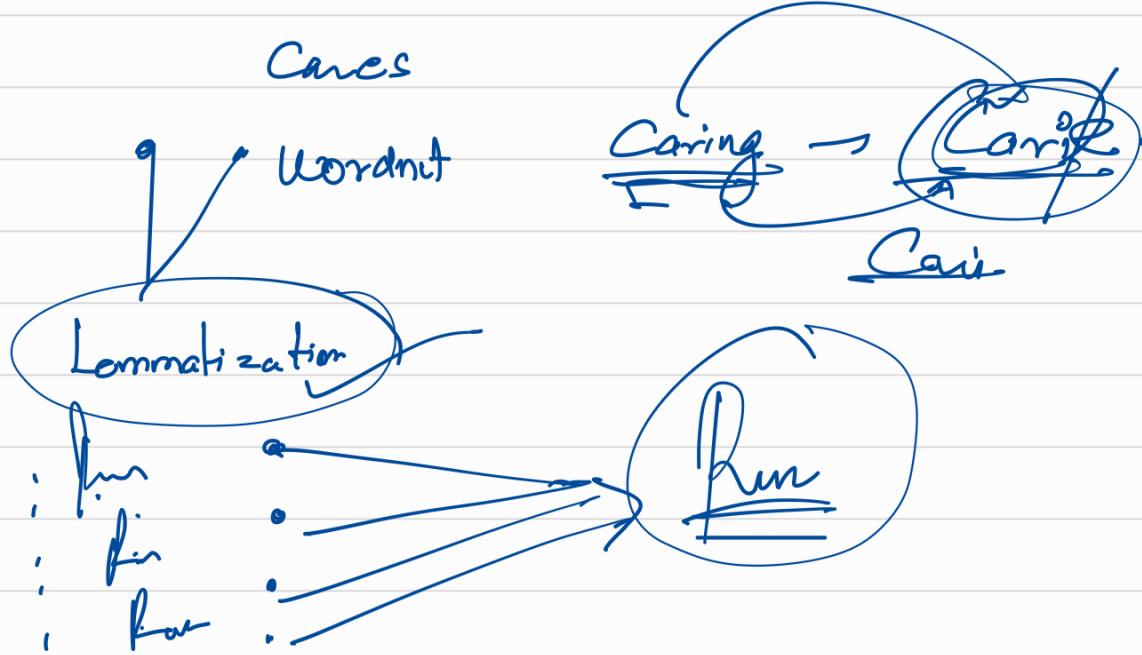
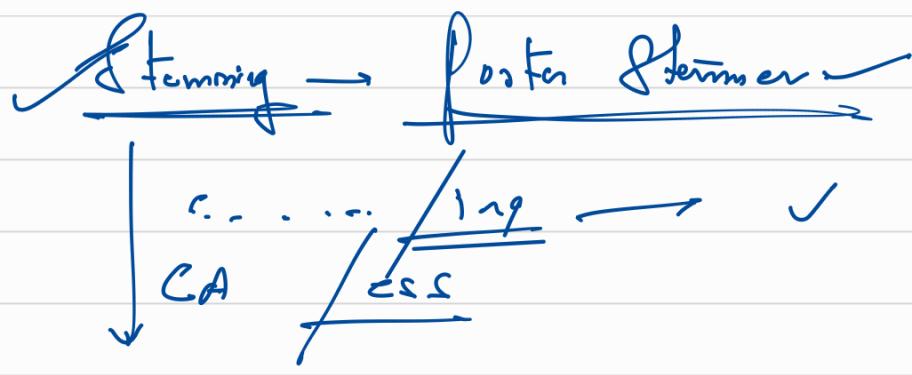


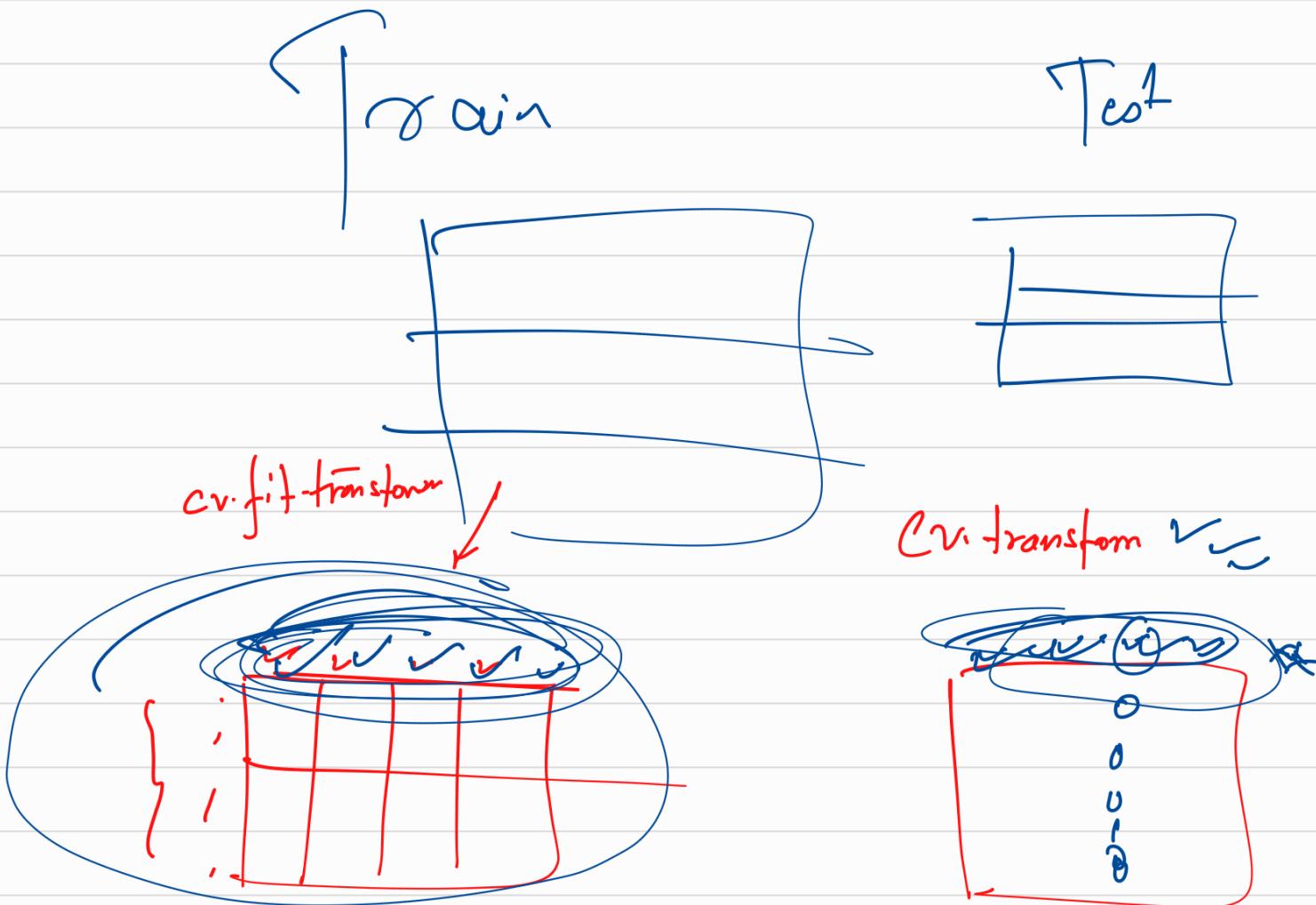
12: 18 AM

Higher dimensionality →

① can remove stop words, remove punctuation,
remove #, Remove numbers ✓







doc1: NLP is nice
 doc2: CV is great.
 doc3: ML awesome

~~train~~

	NLP	is	nice	CV	great	ML	awesome	total
doc1	1	1	1	0	0	0	0	5
doc2	0	0	0	1	1	0	0	4
doc3	0	0	0	0	0	1	1	2
Total	1	1	1	1	1	1	1	9

~~test~~

Cant

order of the words

doc1 : food great ambience food
 doc2 : ambience great food bad.

food great ambience bad

d ₁	1	1	1	1	1	1
d ₂	1	1	1	1	1	1

Bigram

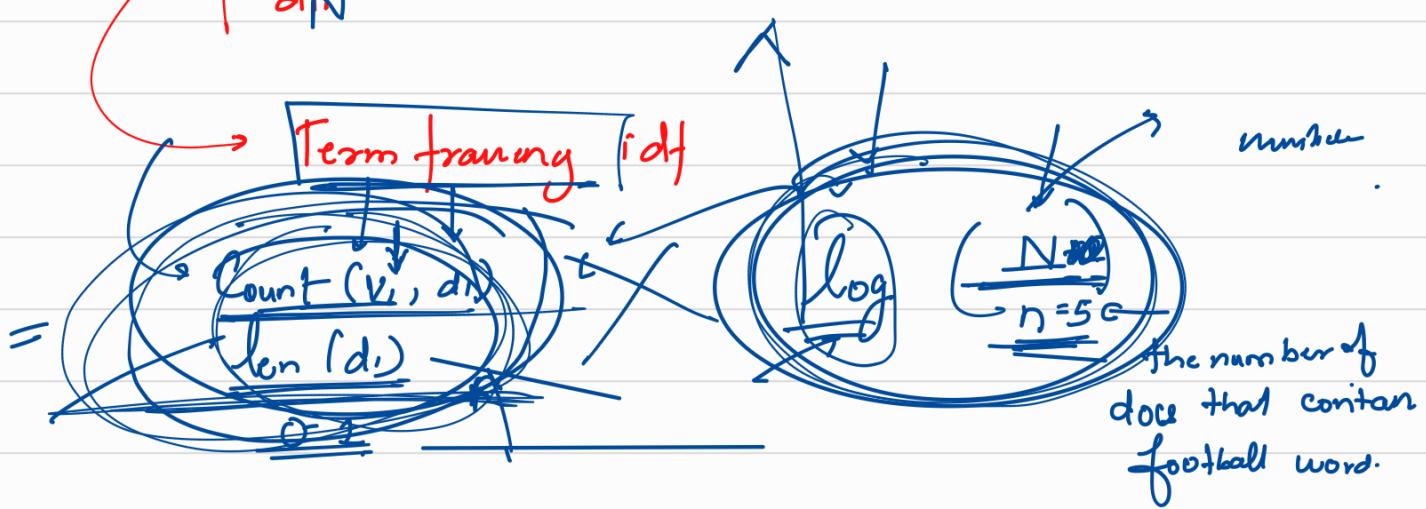
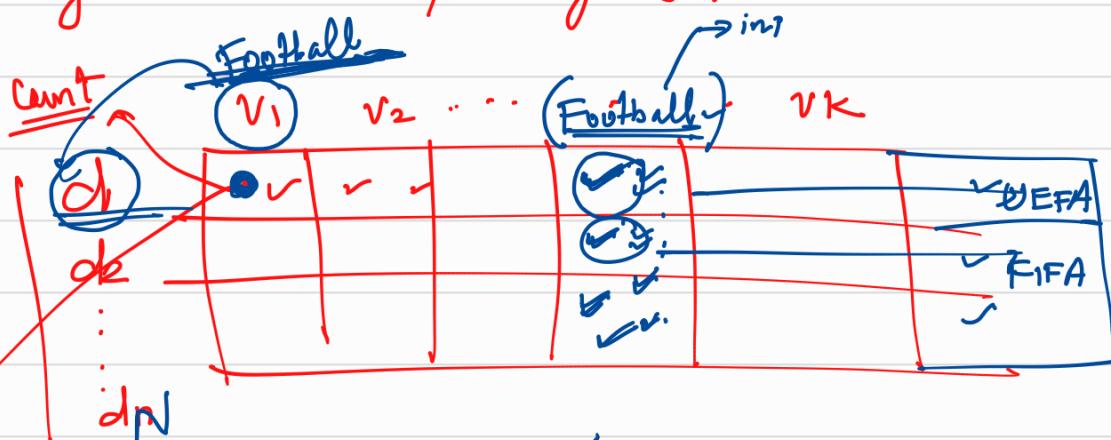
Trigram

N-gram

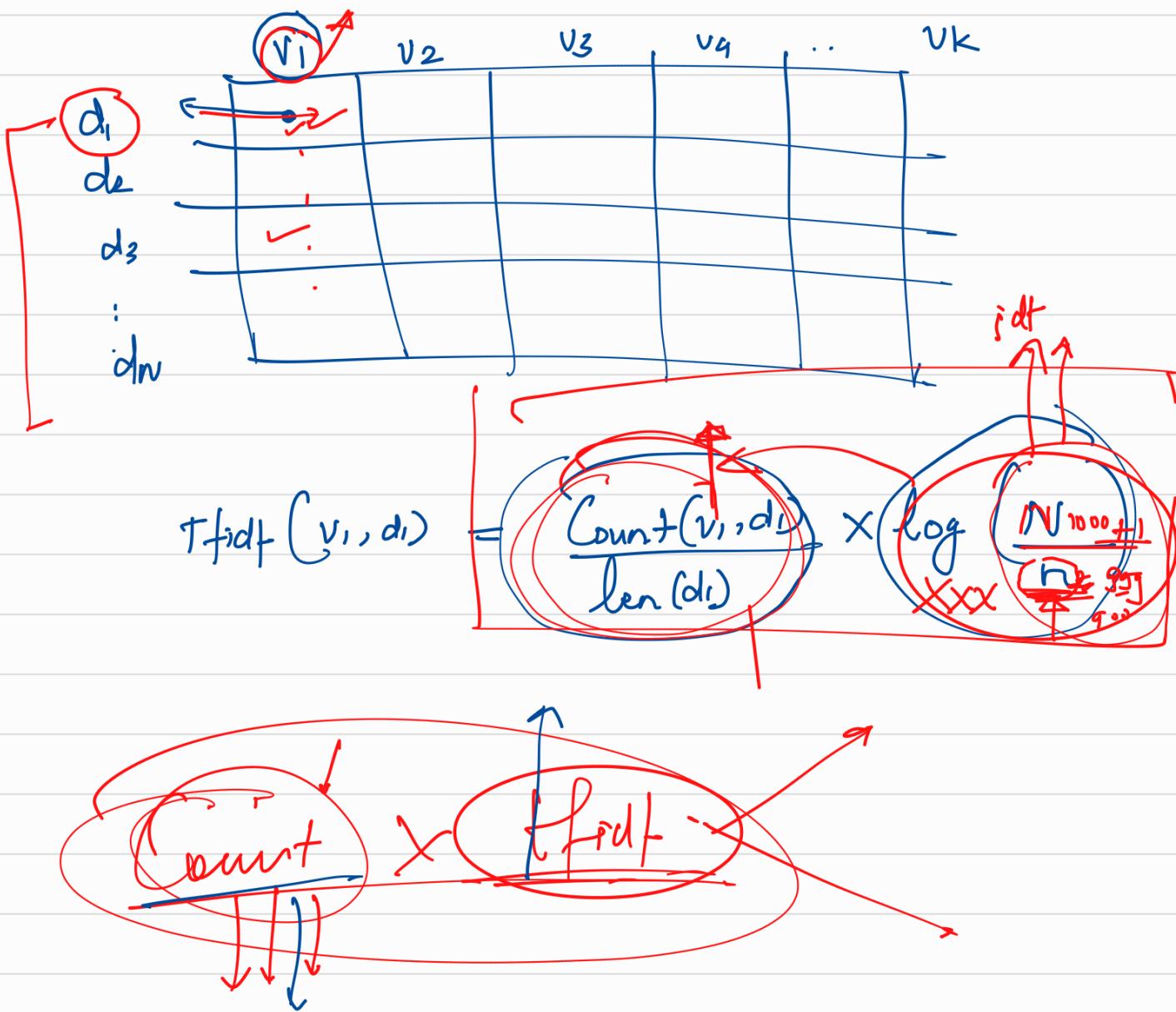
fg ga ab ag gt fb

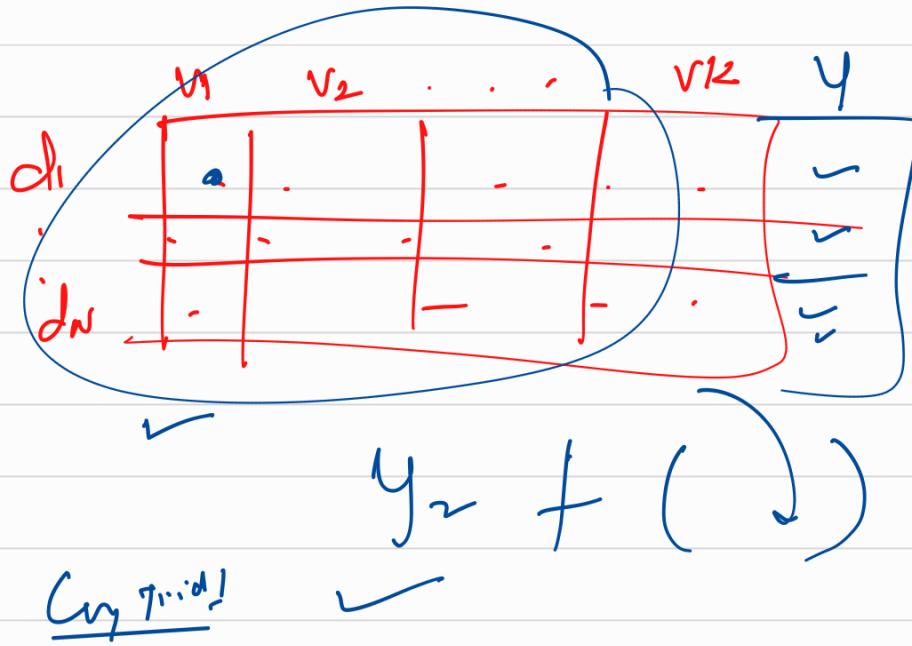
d ₁	1	1	1	0	0	0
d ₂	0	0	0	1	1	1

Term frequency inverse document frequency (tfidf) →

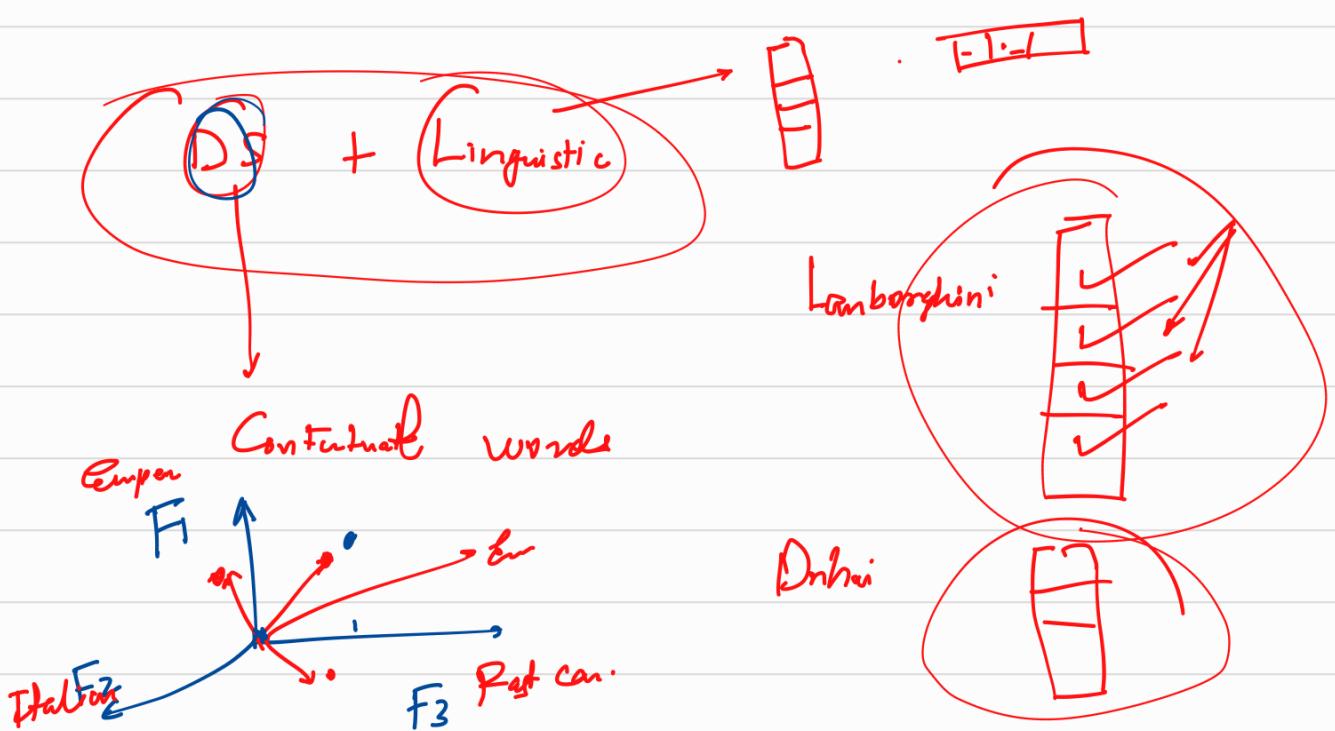
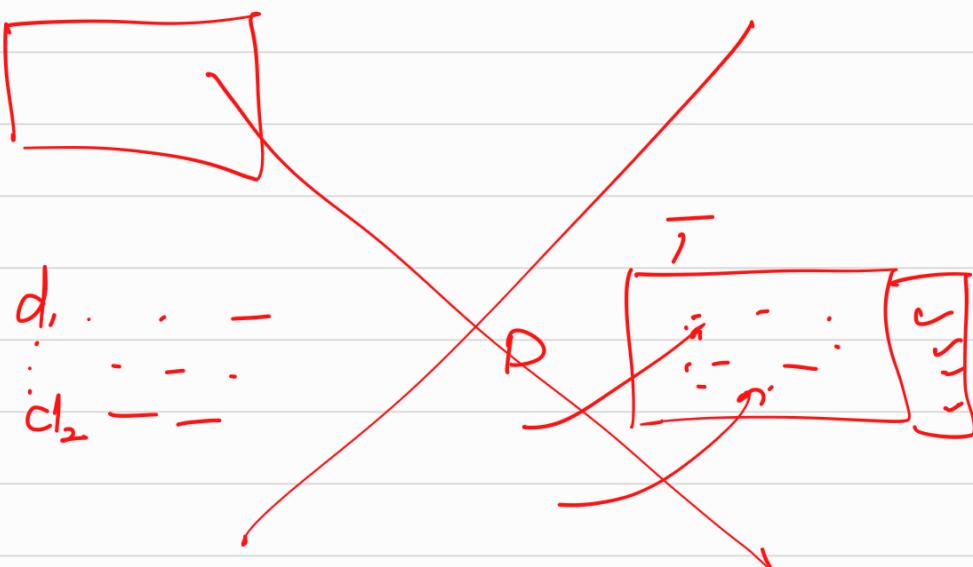


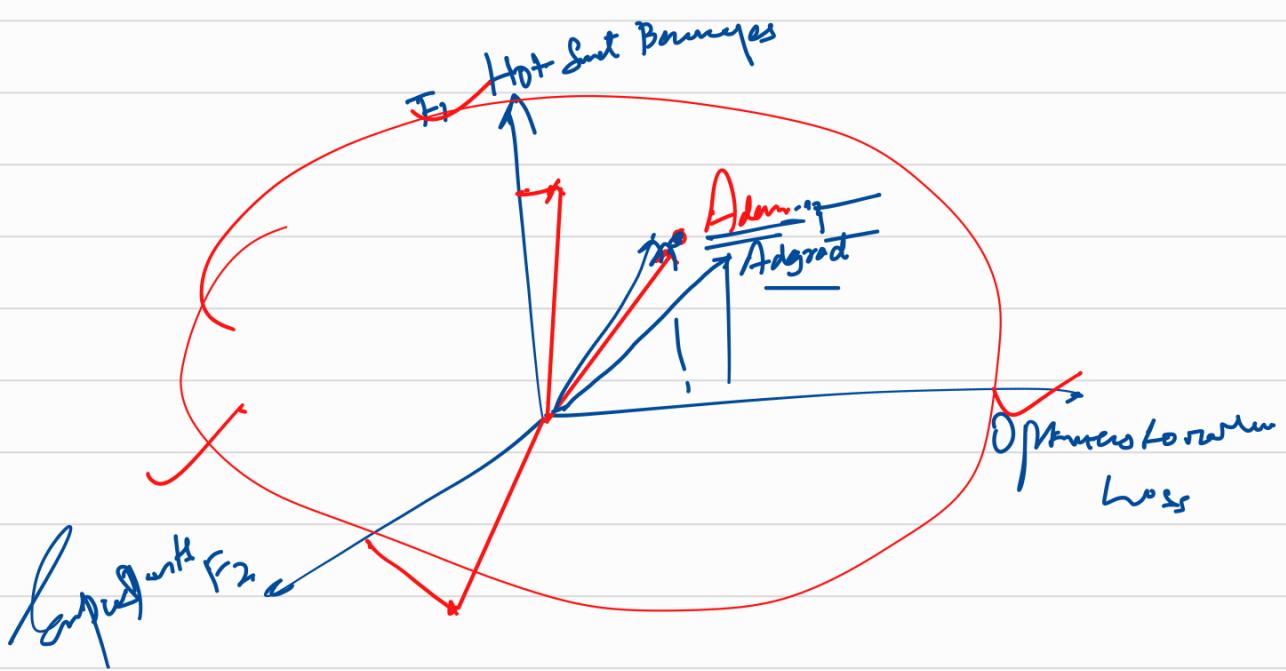
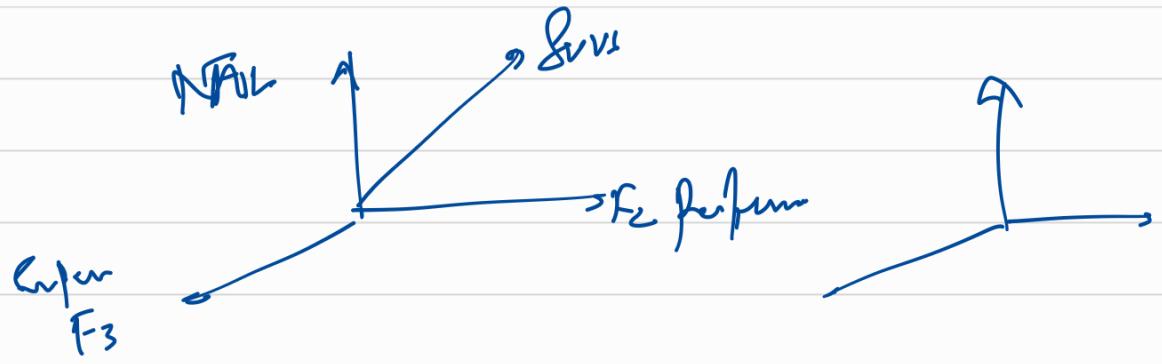
$$\underline{tfidf} = \text{imp}(w, d_i) \times \log \left(\frac{n}{n_i} \right).$$



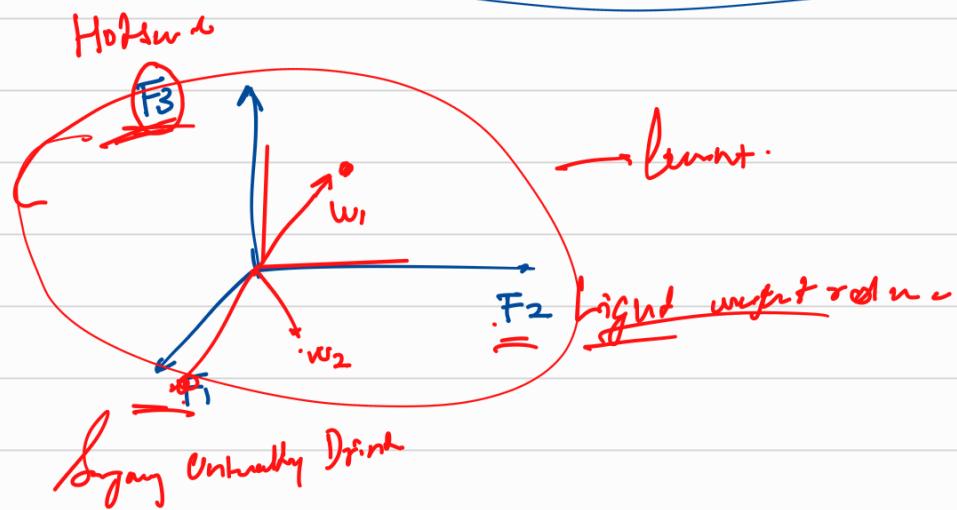


Endorse



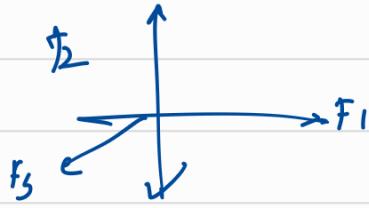


Word Similarity \rightarrow high cosine similarity



Word2vec

word \rightarrow



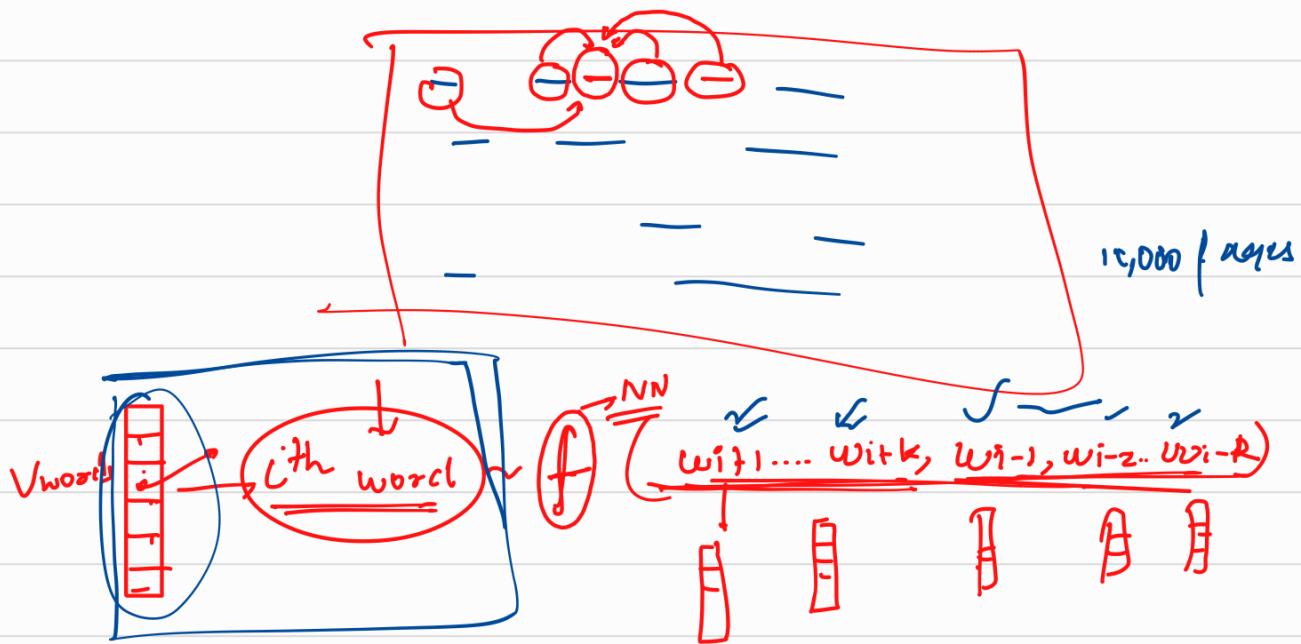
CBOW

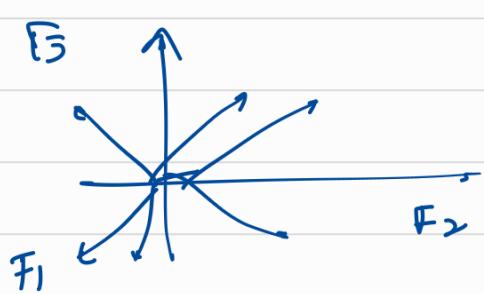
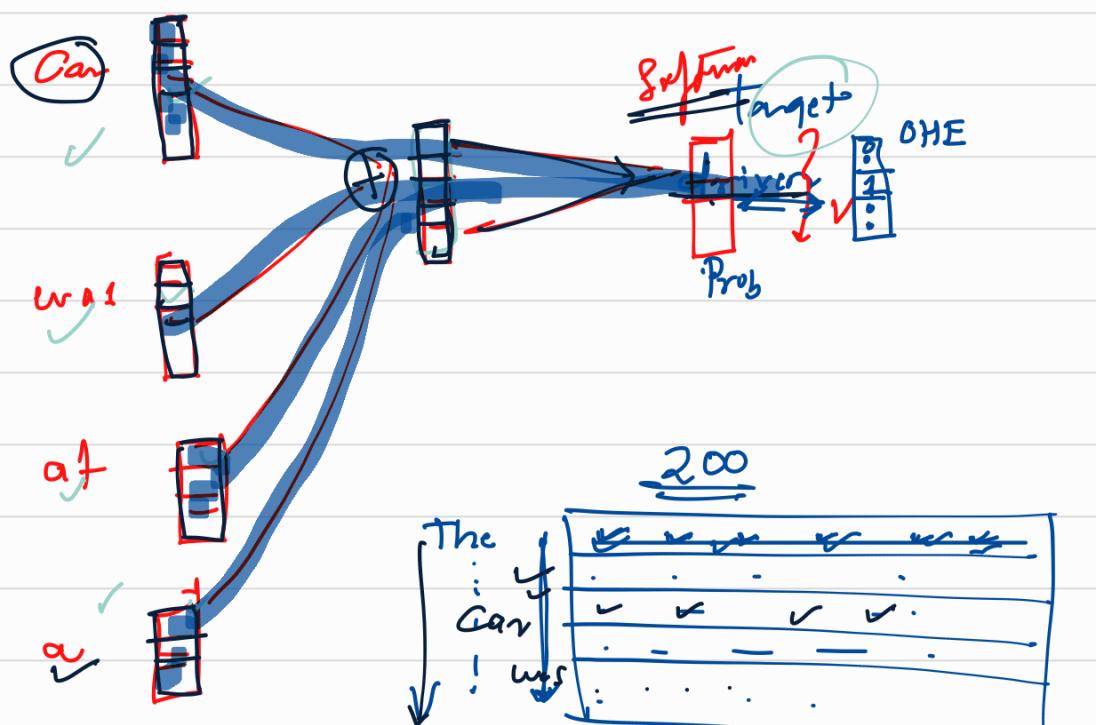
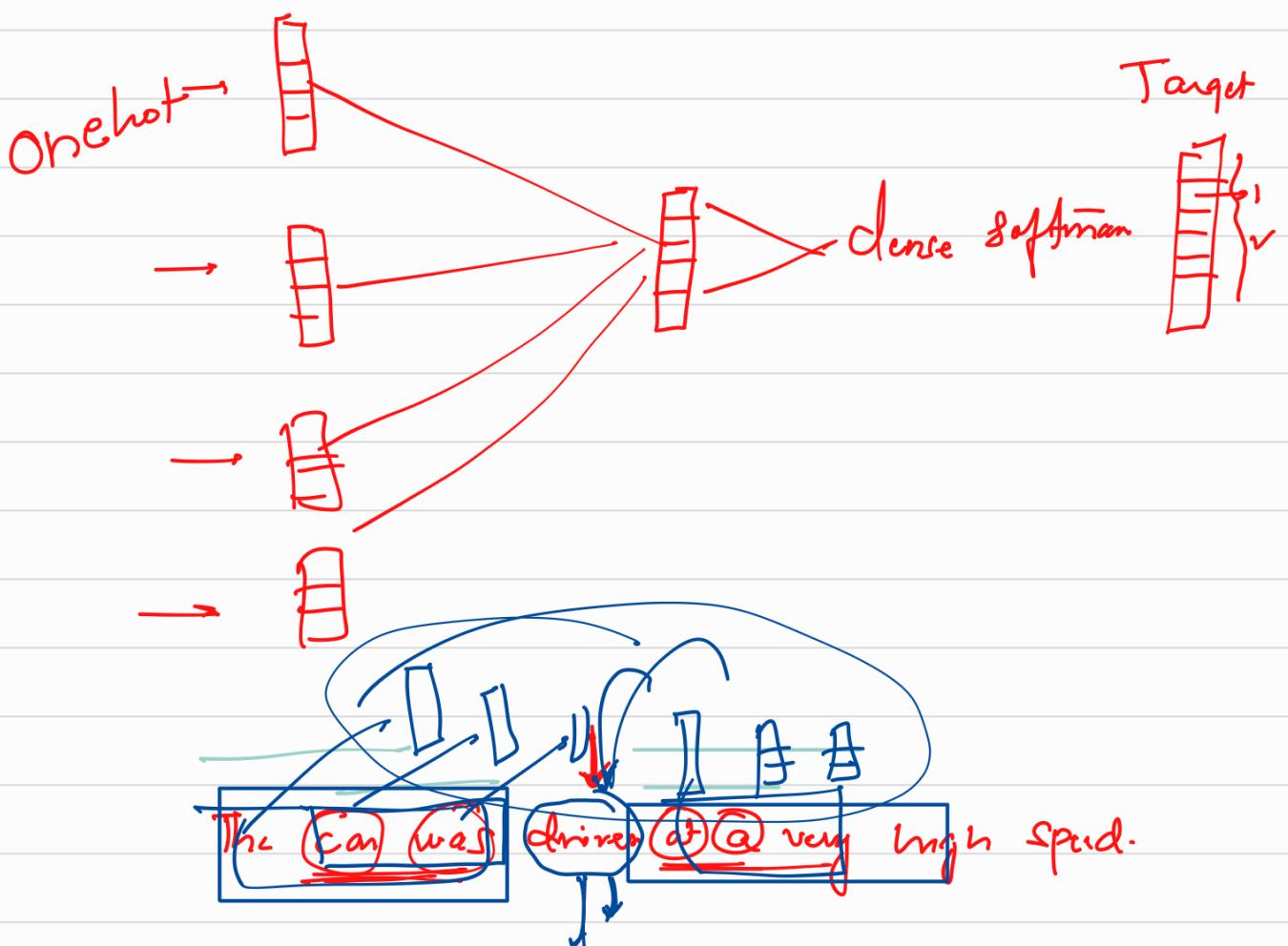
skipgram

CBOW

Corpus

Wikidict. | Google news





Embedding matrix α

