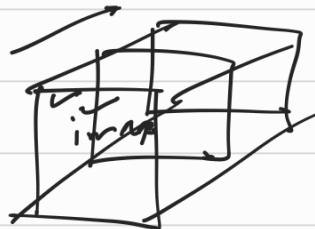


Natural Language Processing

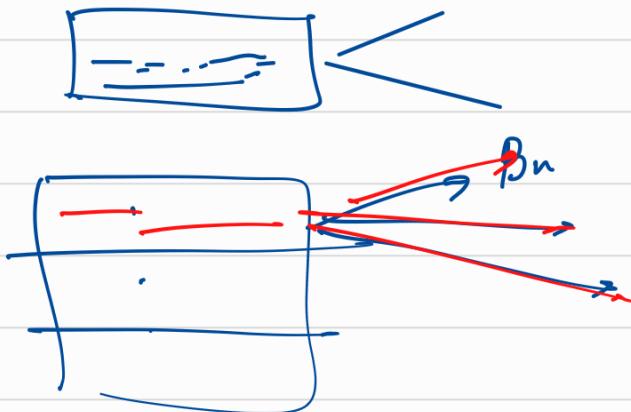
to make machines → human language.



→ don't have a universal representation in numbers.

→ modelling.

→ Classification

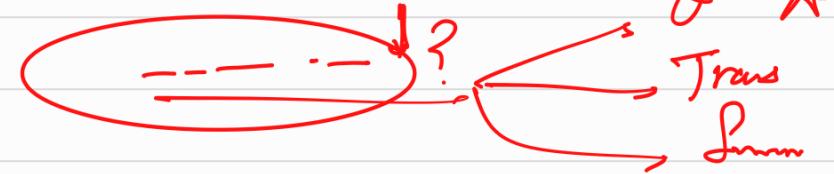


→ Named Entity

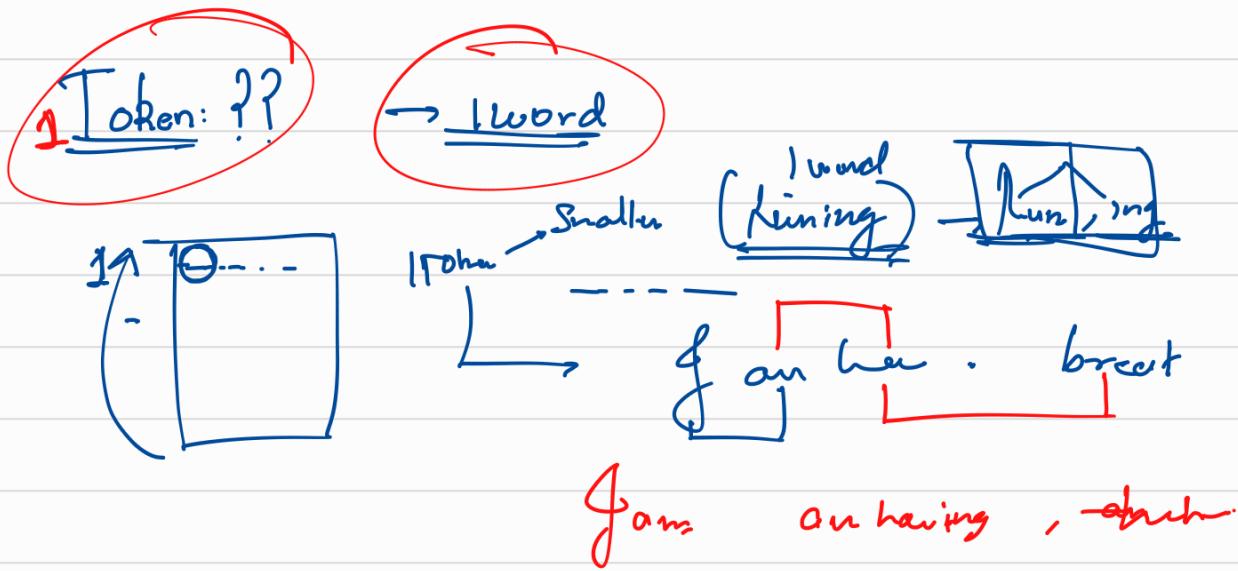
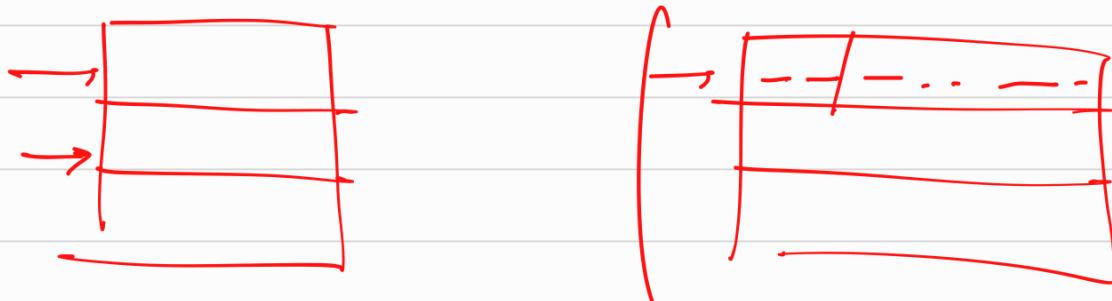
Timing, Topic

→ Have a ~~saturday~~ class at 11:00 am on Saturday about NLP.

(Large) Language modeling



Definitions:



Vocabulary:

Union
Collection of Tokens

- NLP
- ① text → numbers
 - ② Model

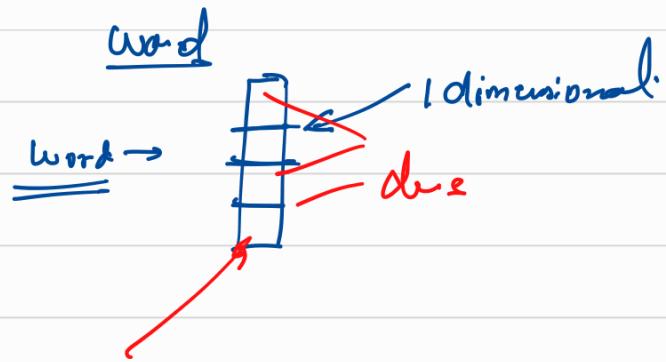
text → numbers (Embedding).

Prediction

→ (Frequency)

text → numbers

Frequency

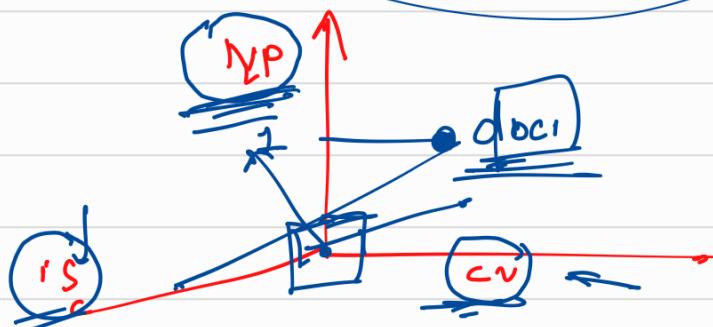
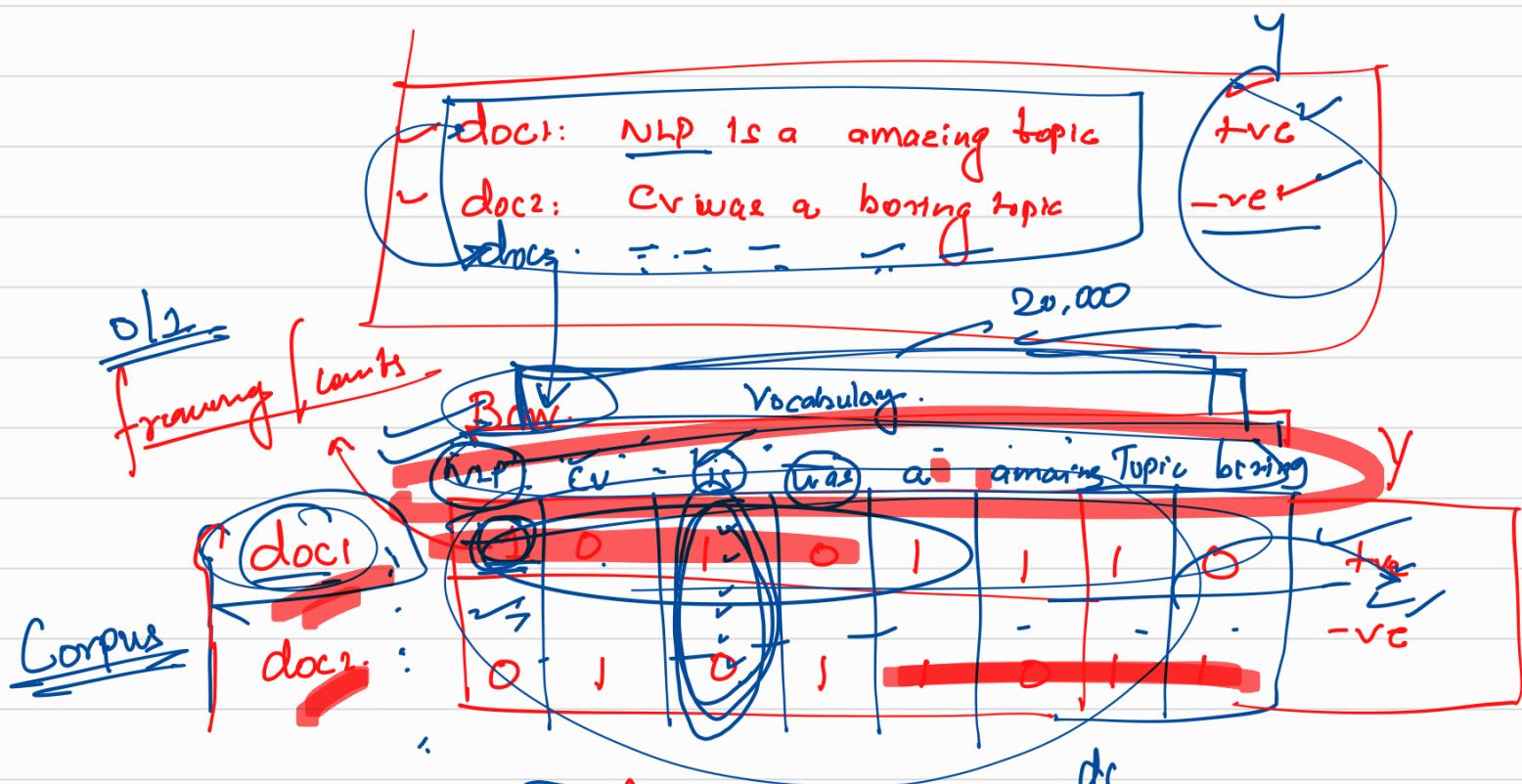


embedding



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

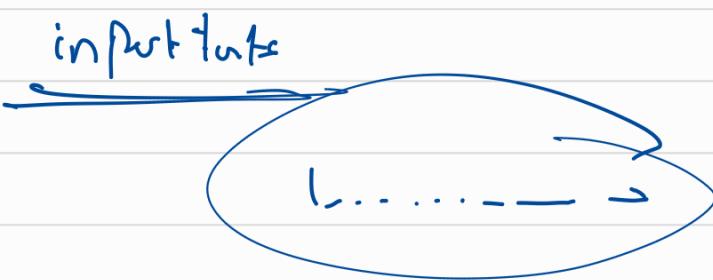
vectors



High dimensions → Sparcity

Ignores → order

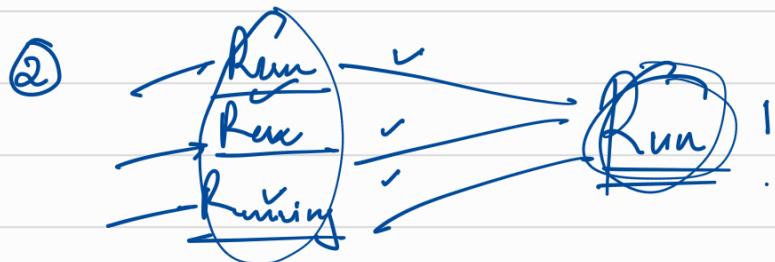
~~No Contextual Information~~

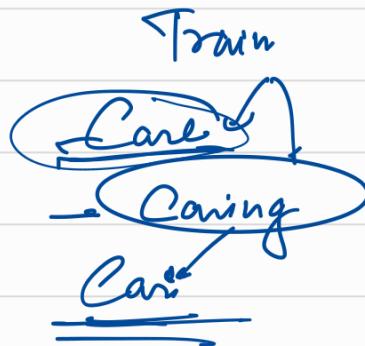
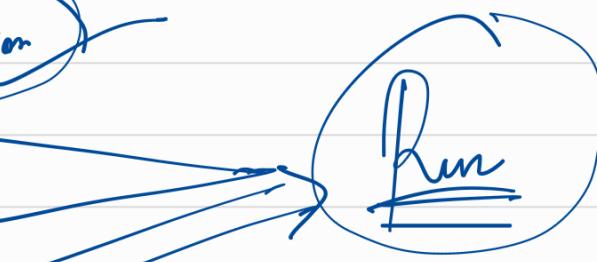
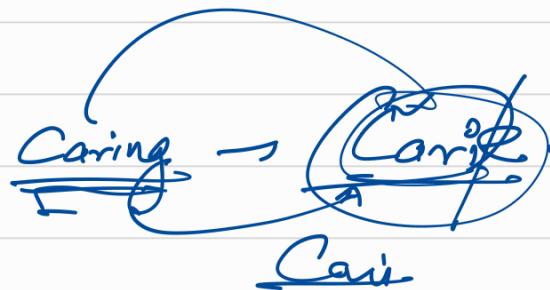
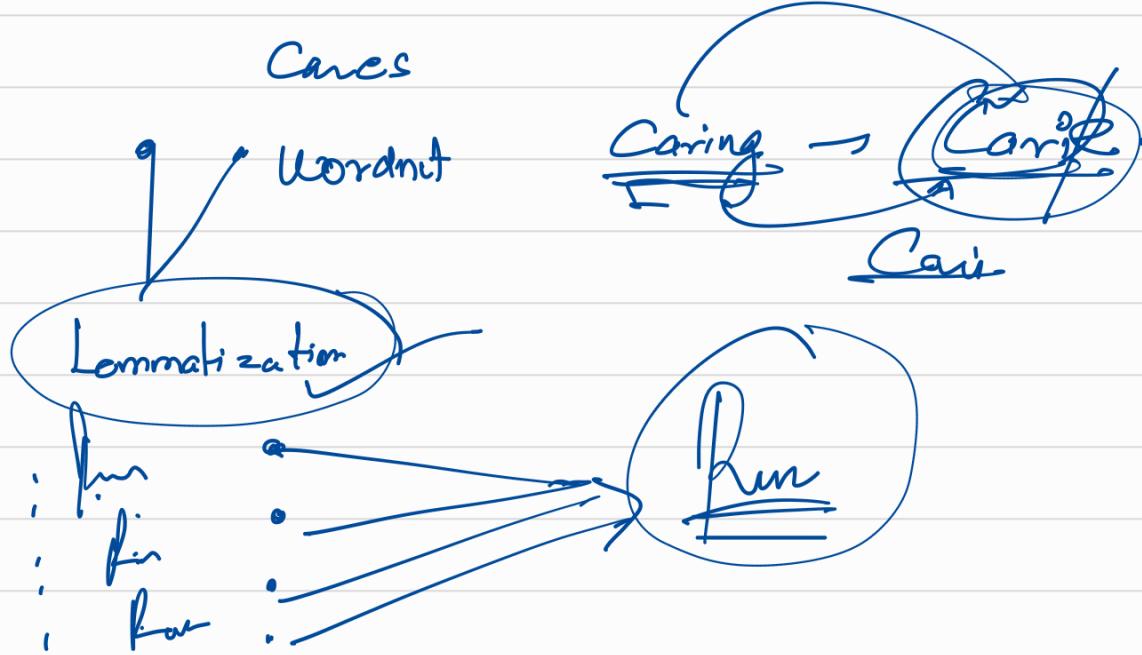
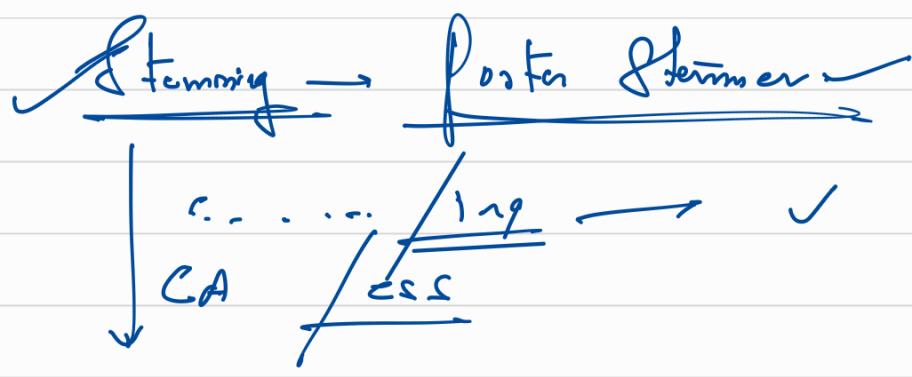


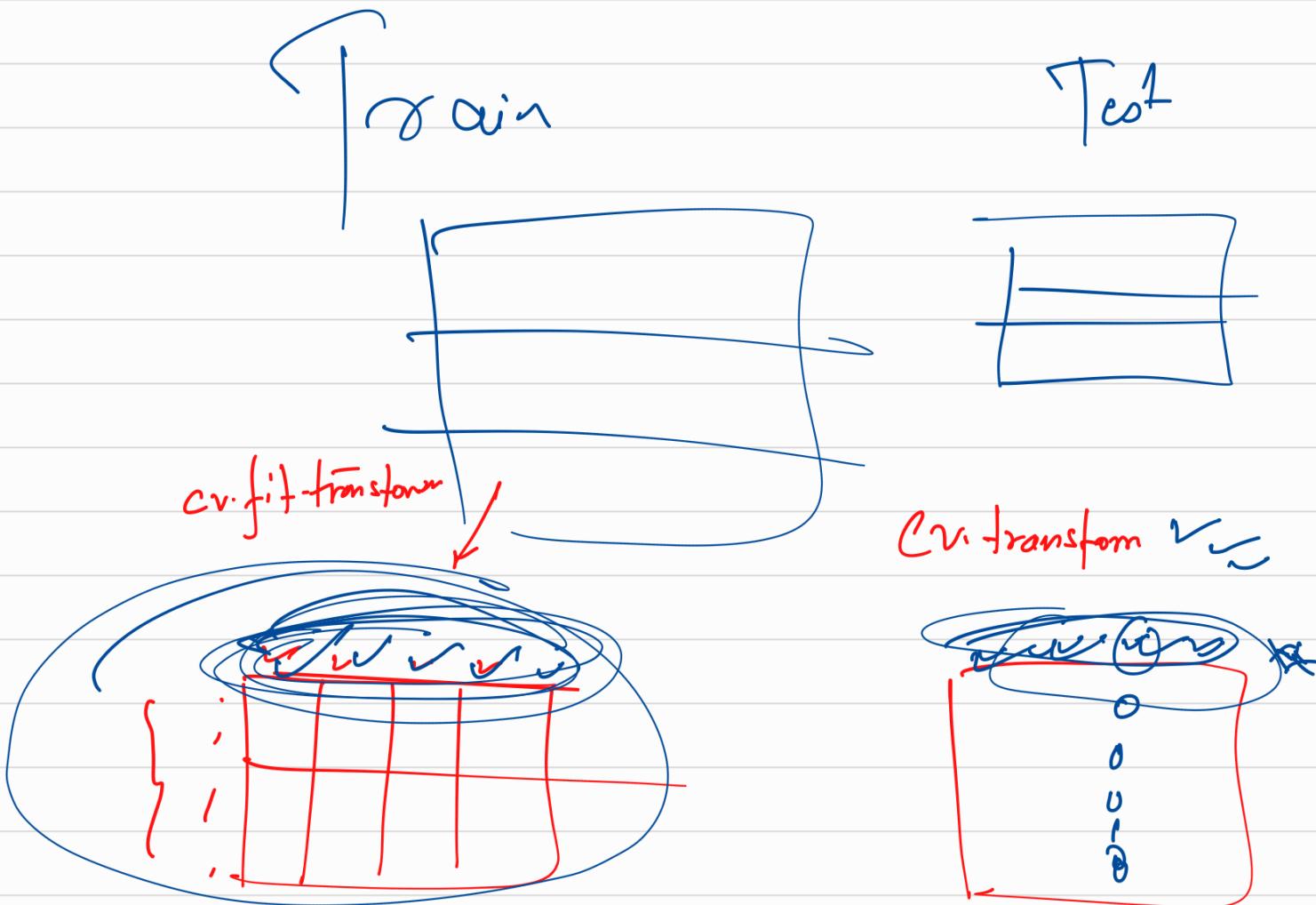
12: 18 AM

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	
doc1	1	1	1	0	0	0	0	train
doc2	0	0	0	1	1	0	0	test
log	0	0	0	0	0	1	1	

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

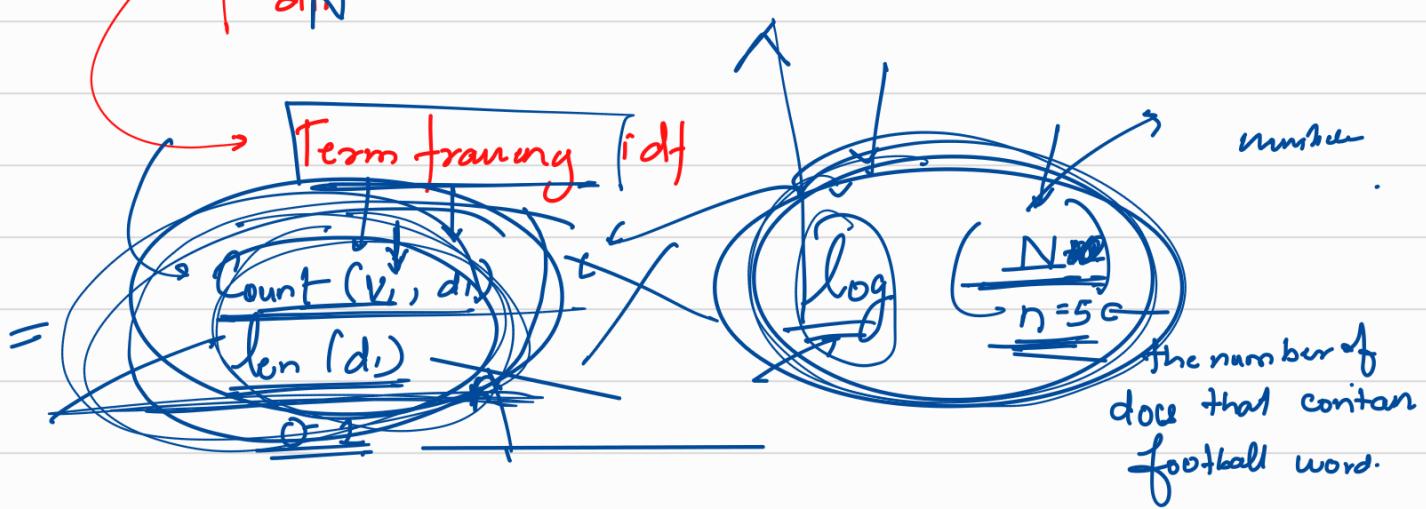
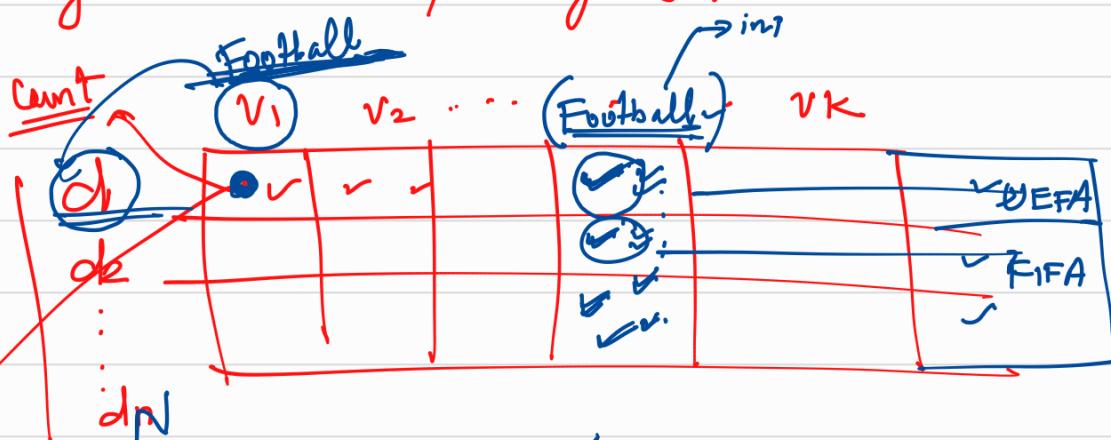
Trigram

N-gram

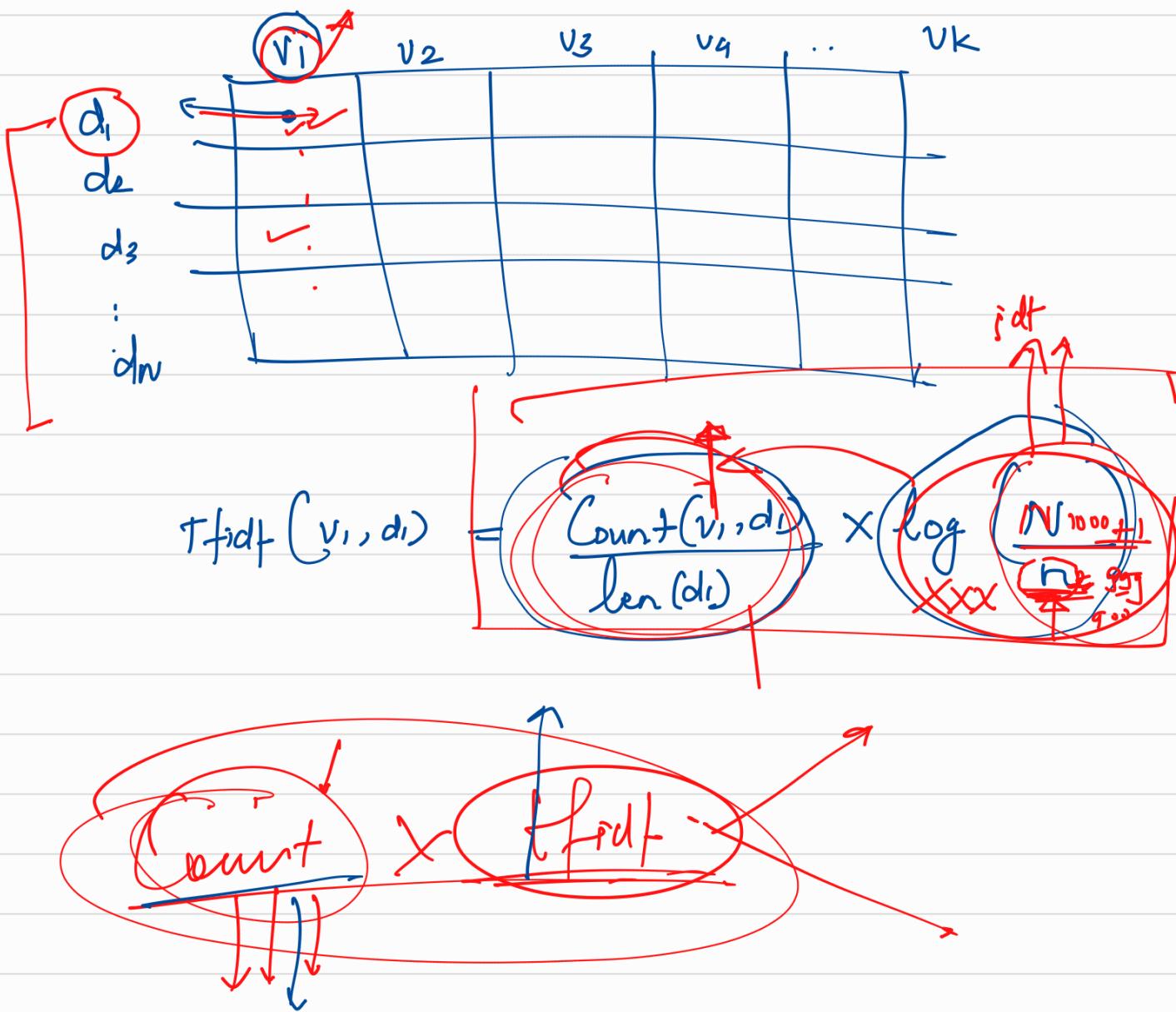
fg ga ab ag gt fb

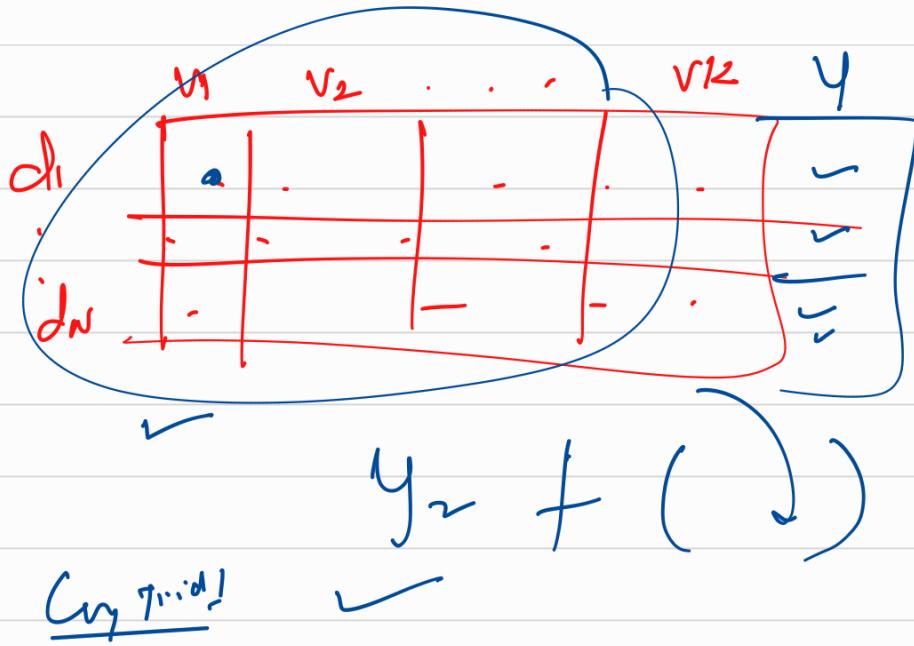
d ₁	1	1	1	0	0	0
d ₂	0	0	0	1	1	1

Term frequency inverse document frequency (tfidf) →

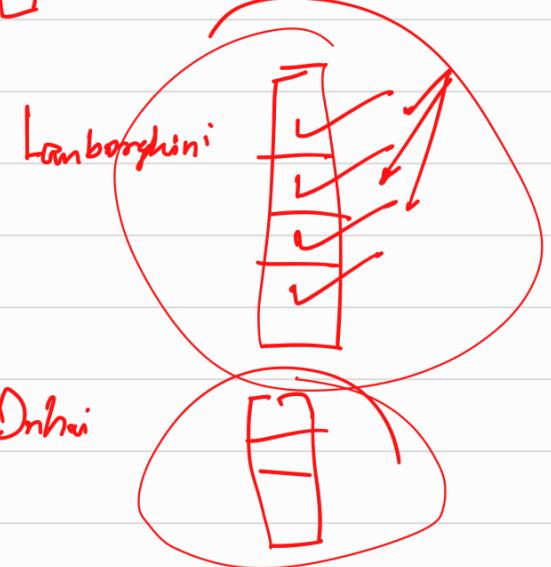
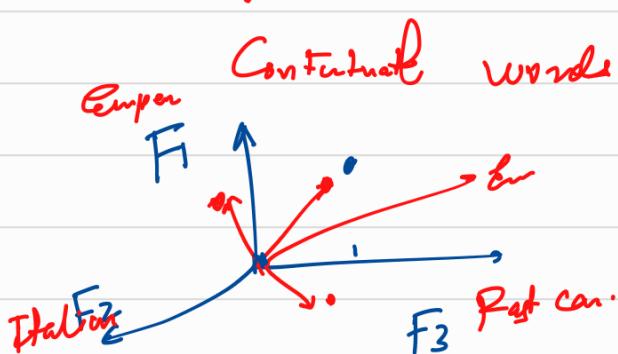
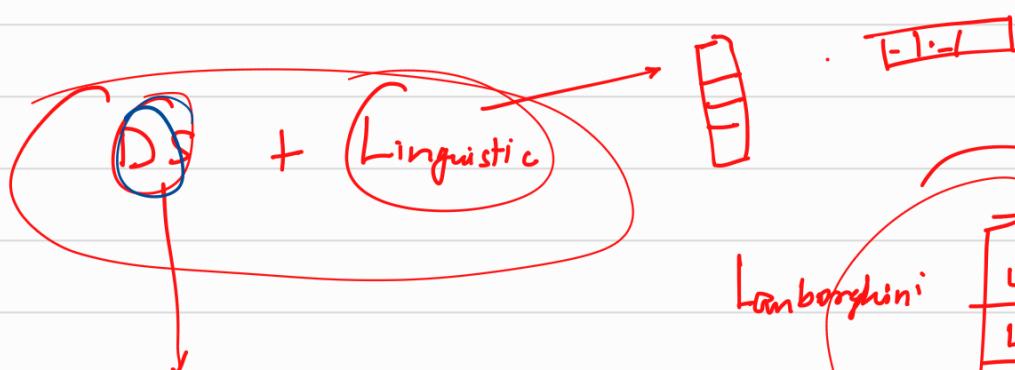
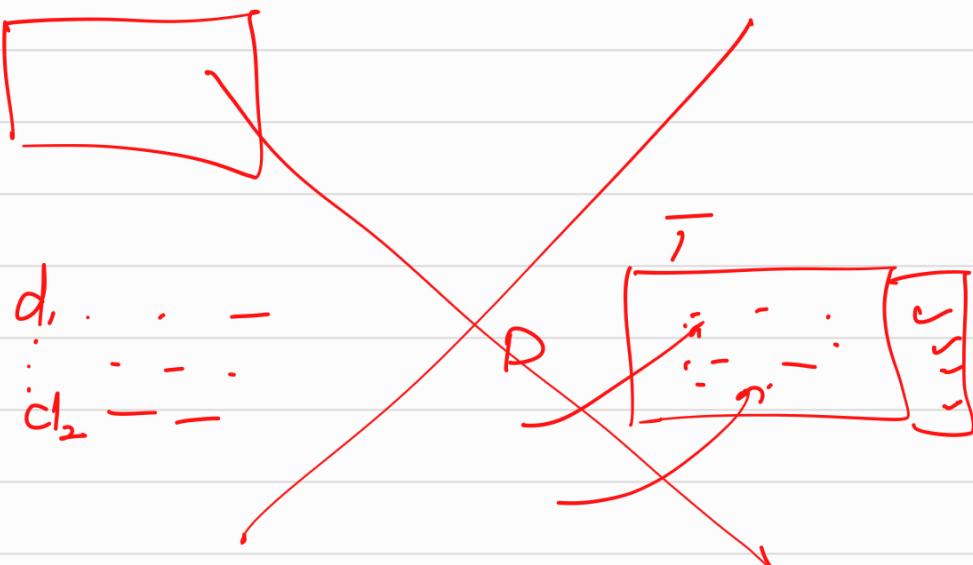


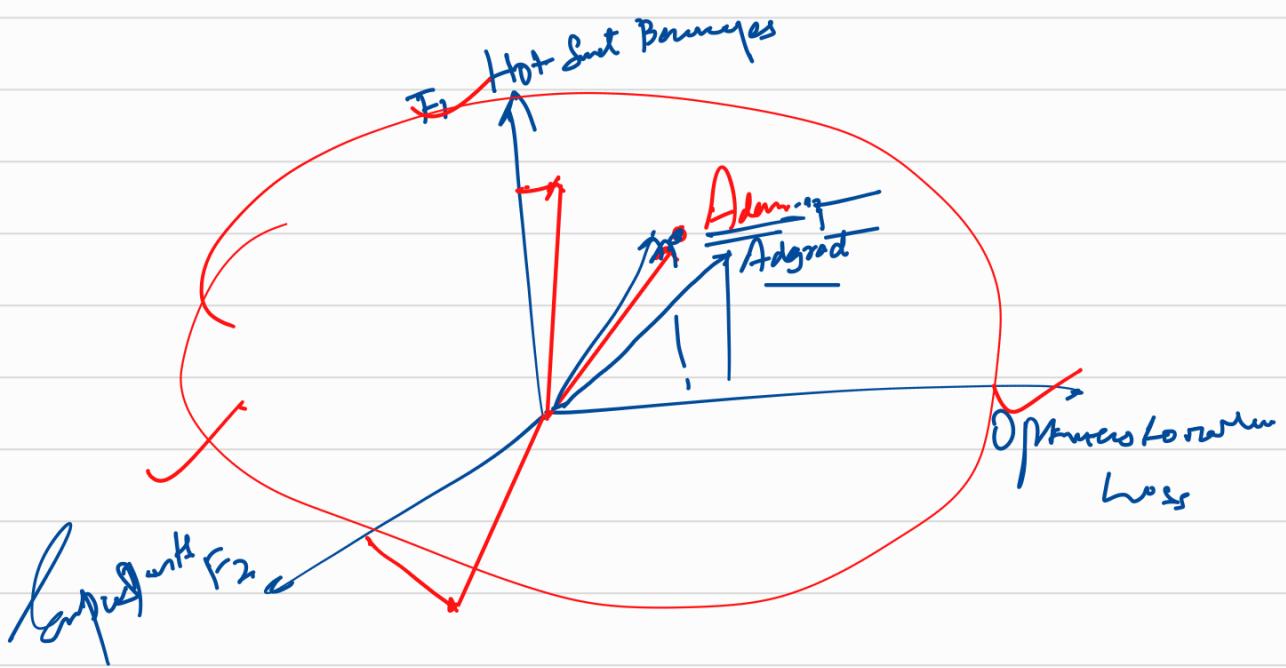
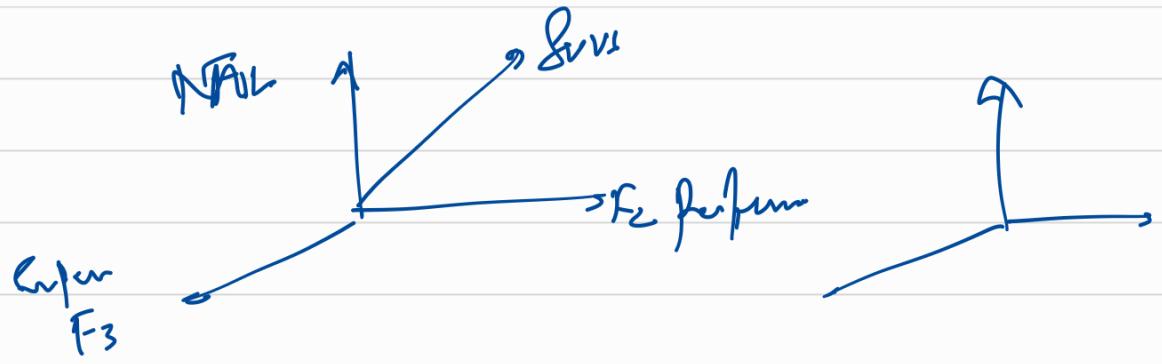
$$\underline{tfidf} = \text{imp}(w, d_i) \times \log \left(\frac{n}{n_i} \right).$$



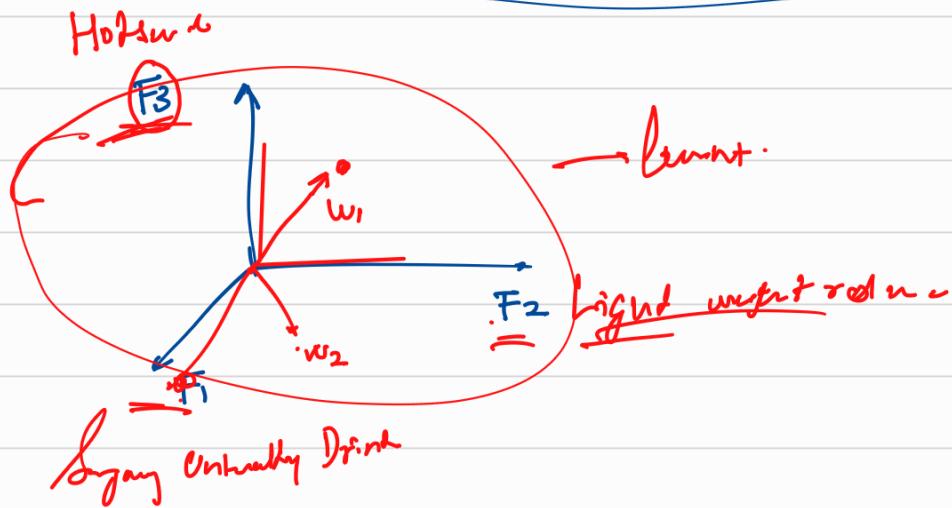


Euston

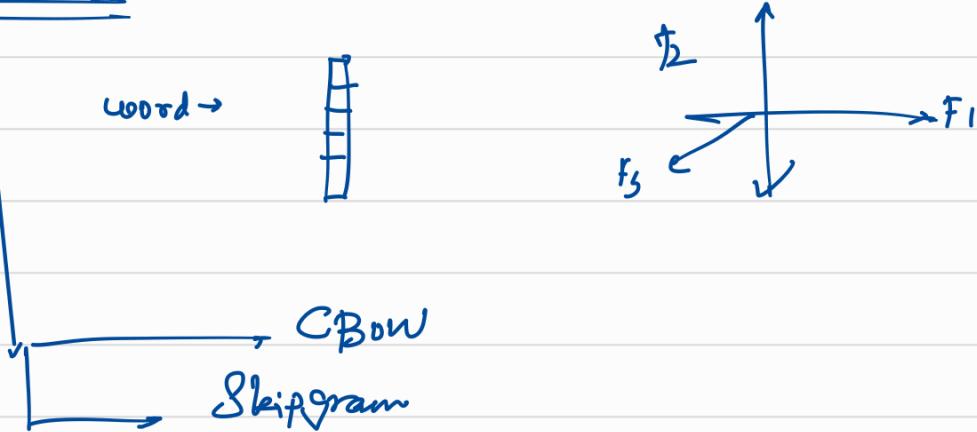




Word Similarity \rightarrow high cosine similarity



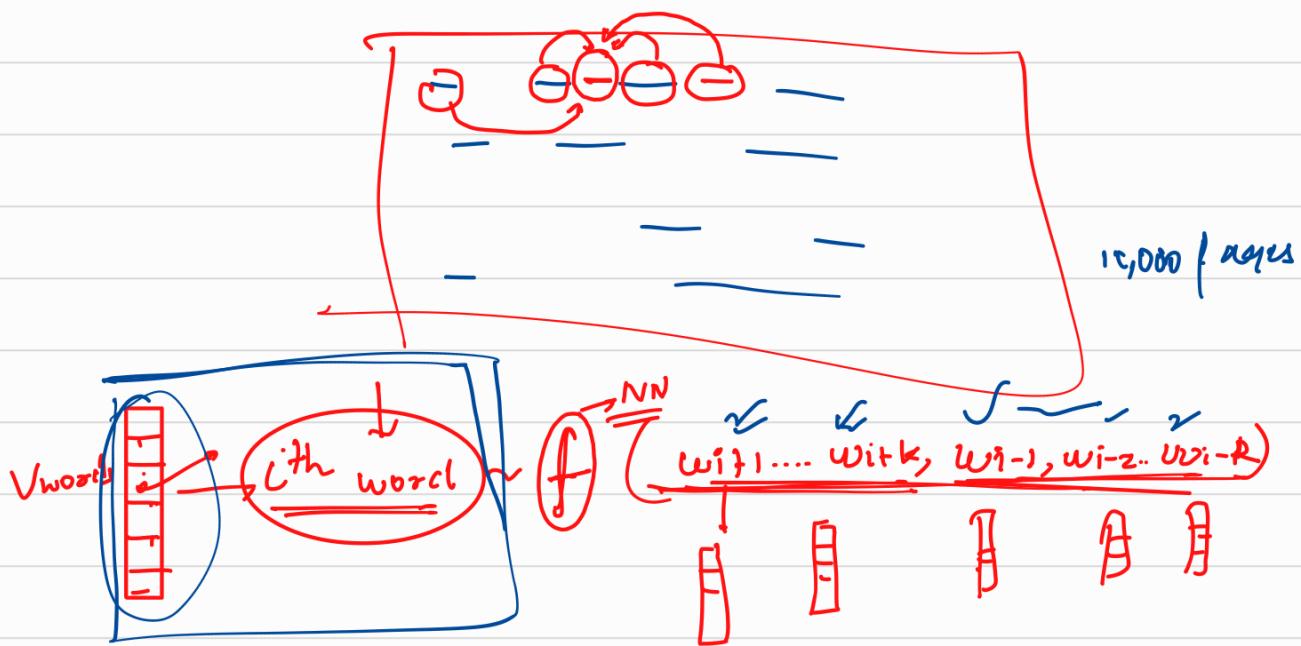
Word2vec

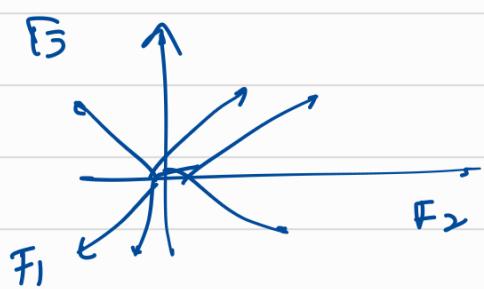
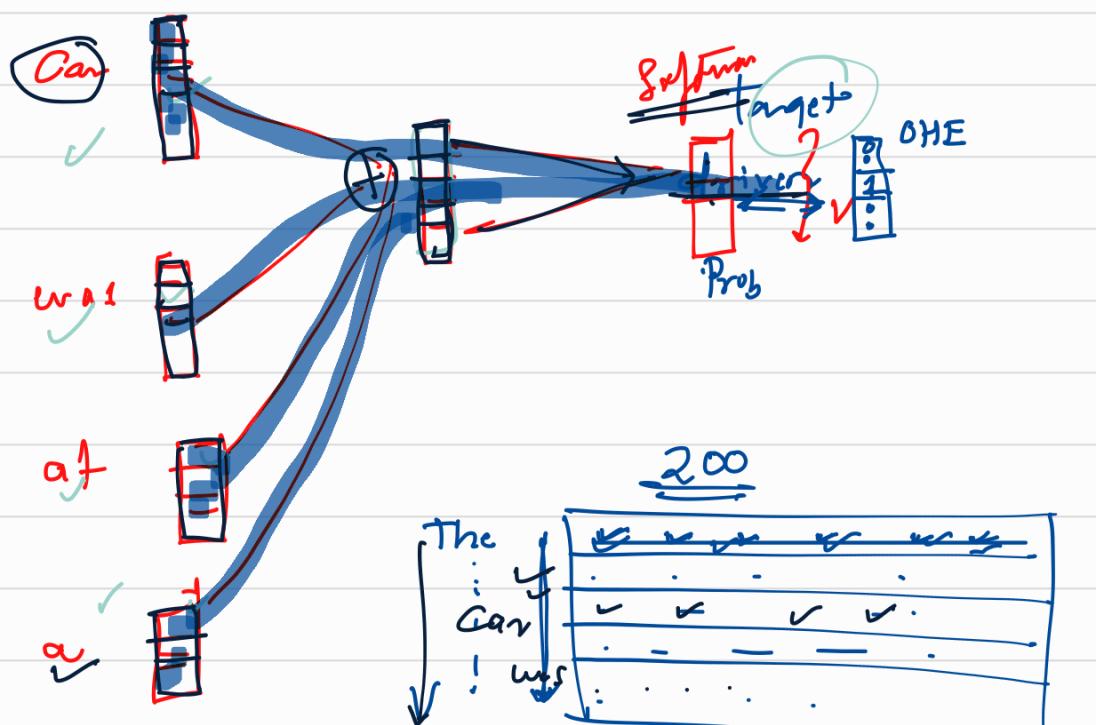
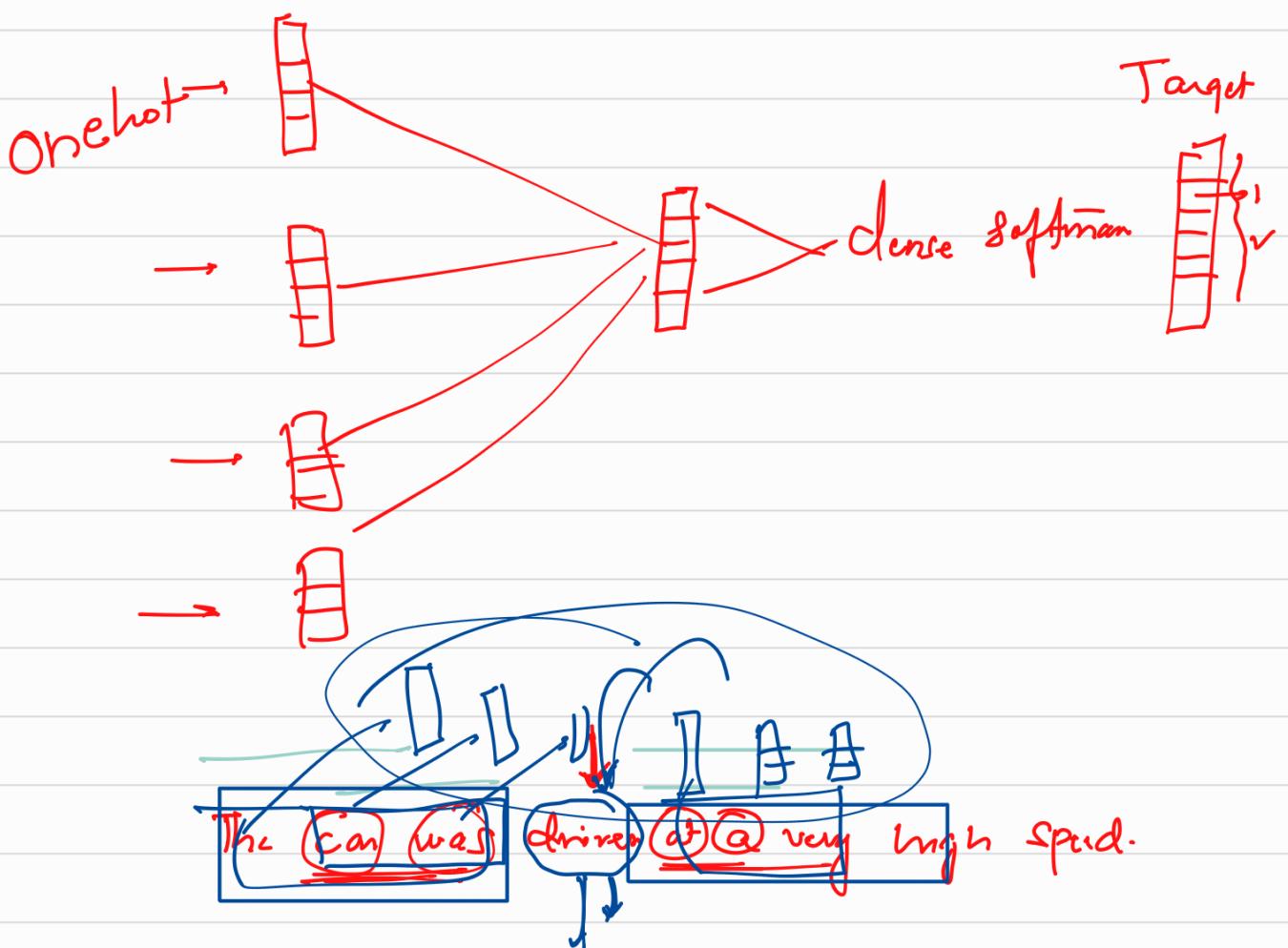


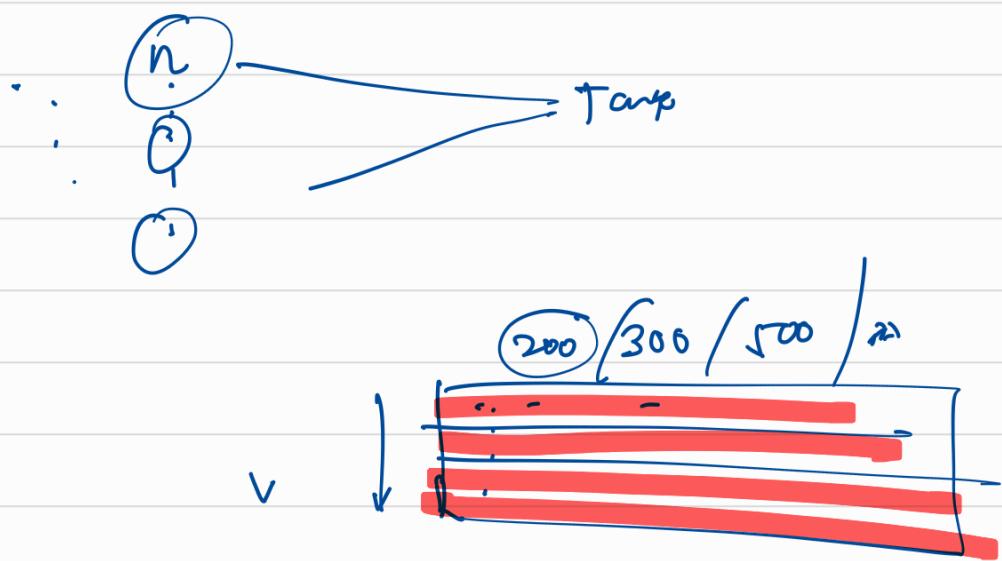
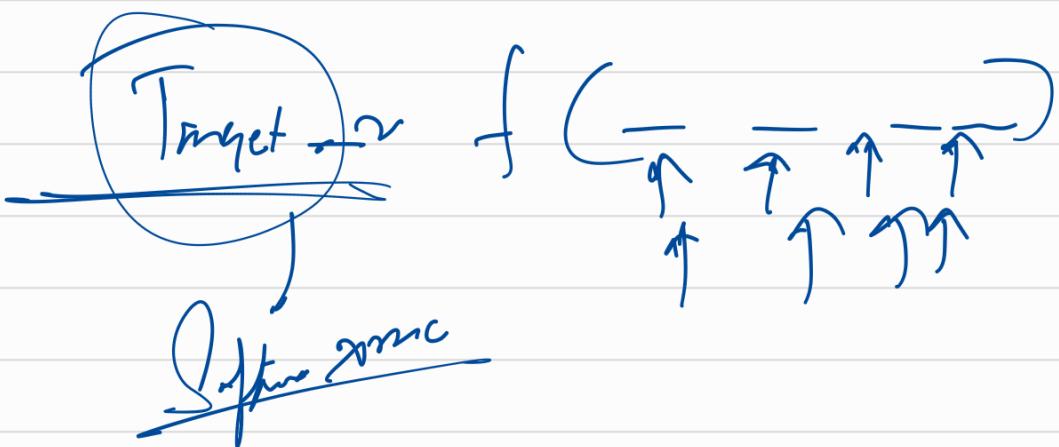
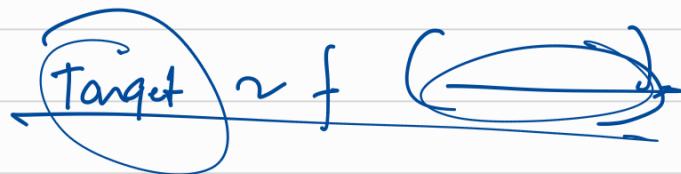
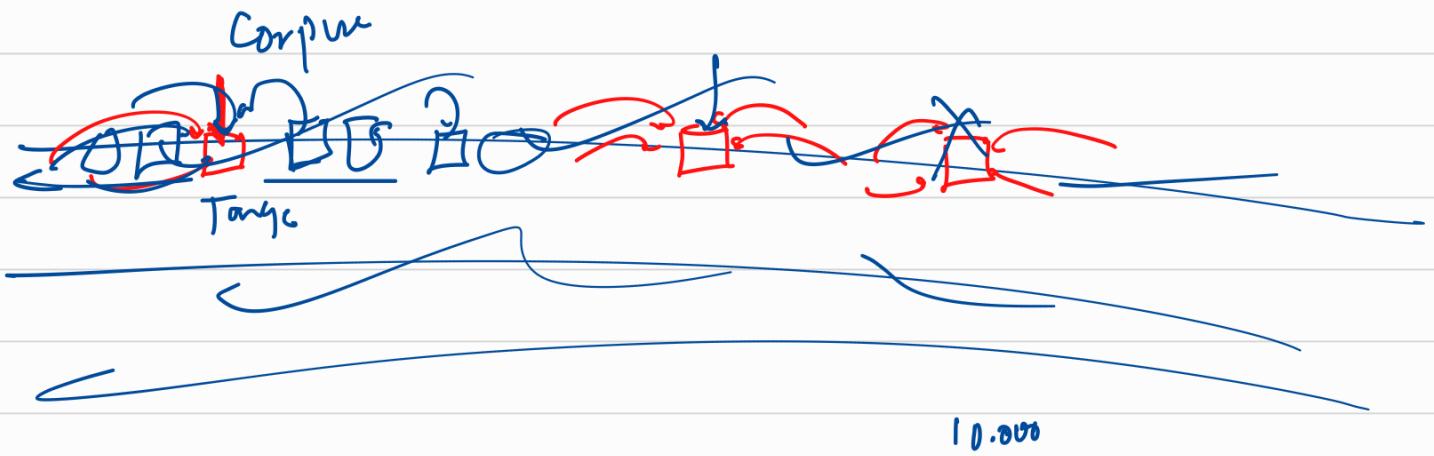
CBOW

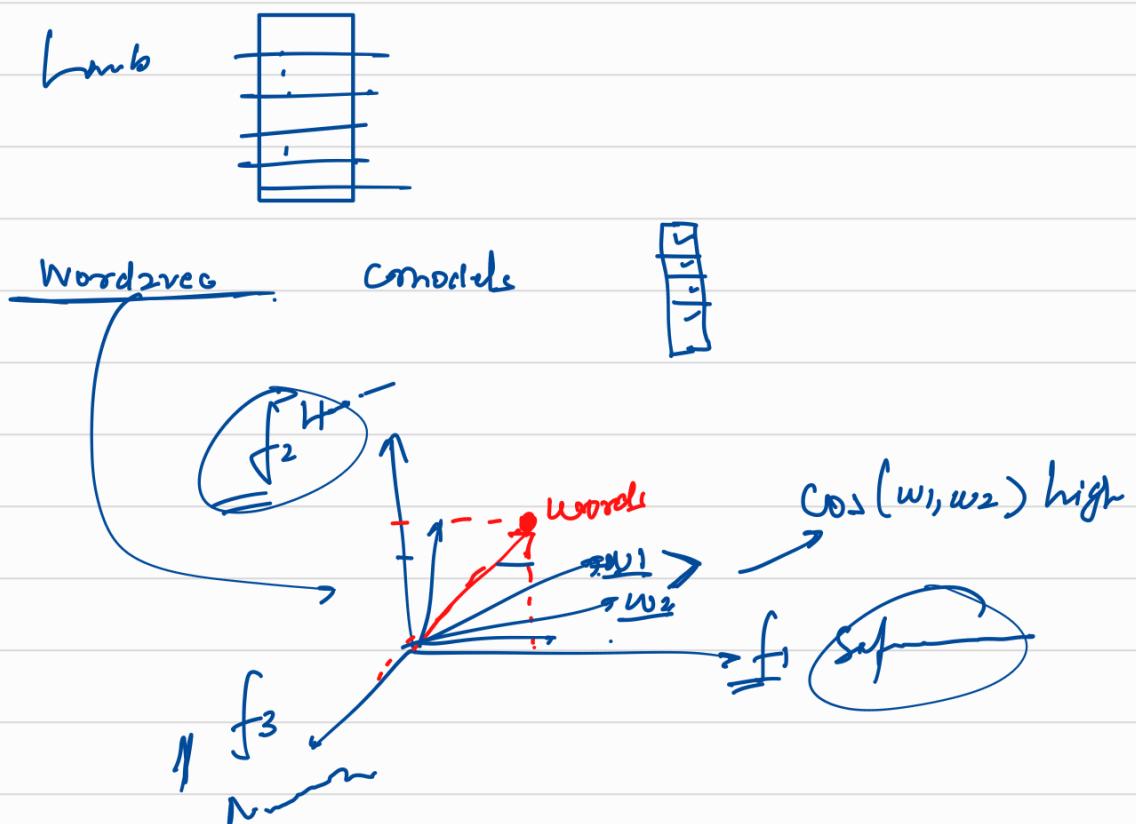
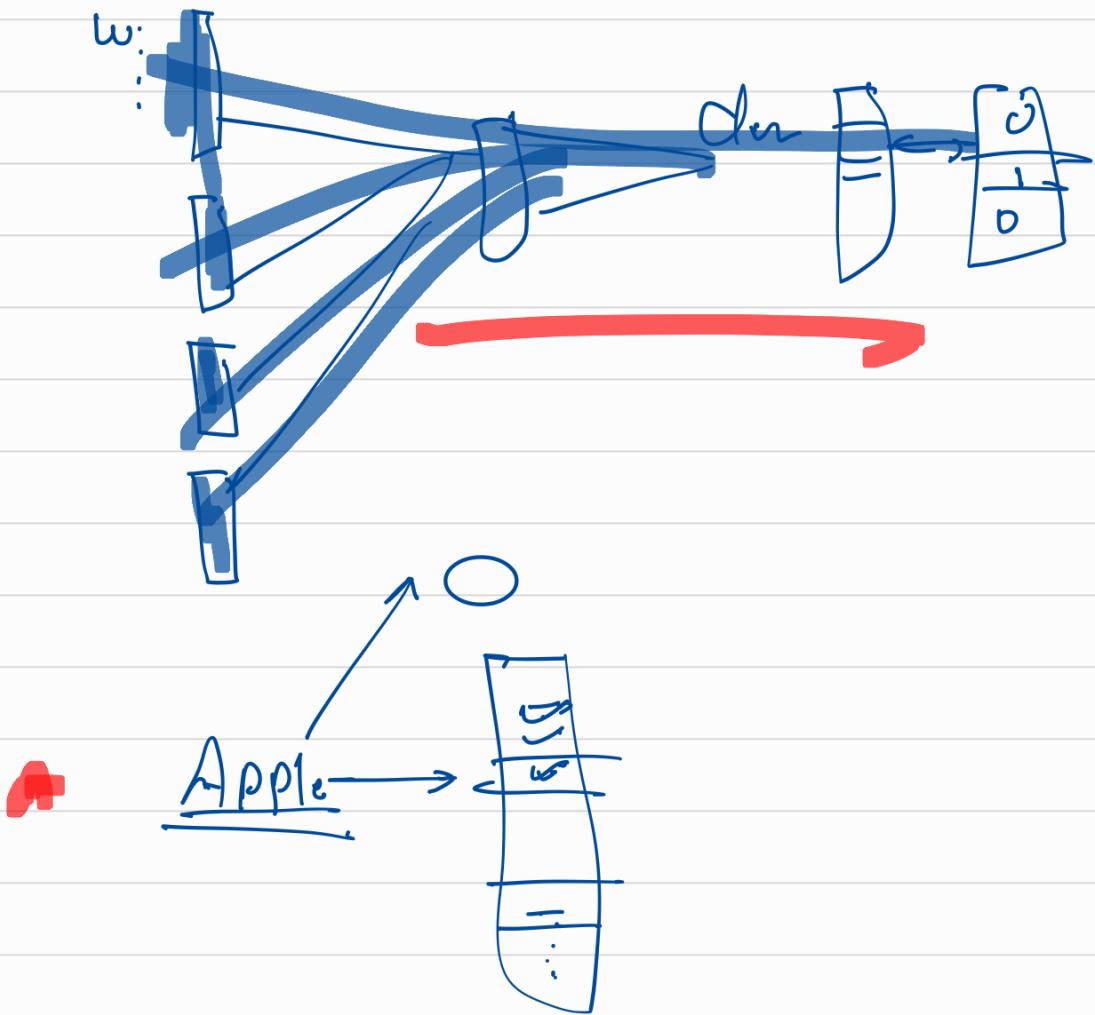
Corpus

Wikidict. | Google news



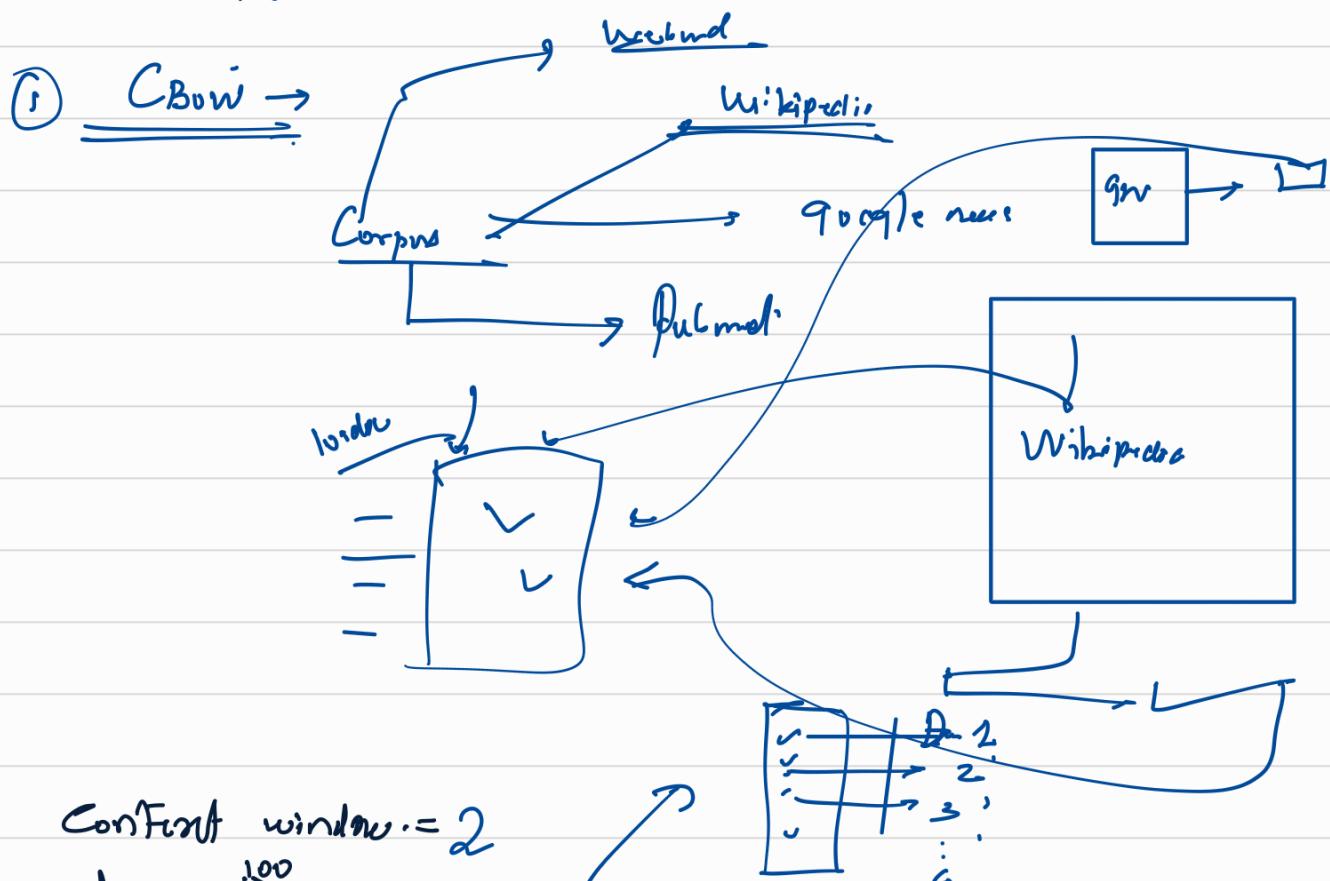






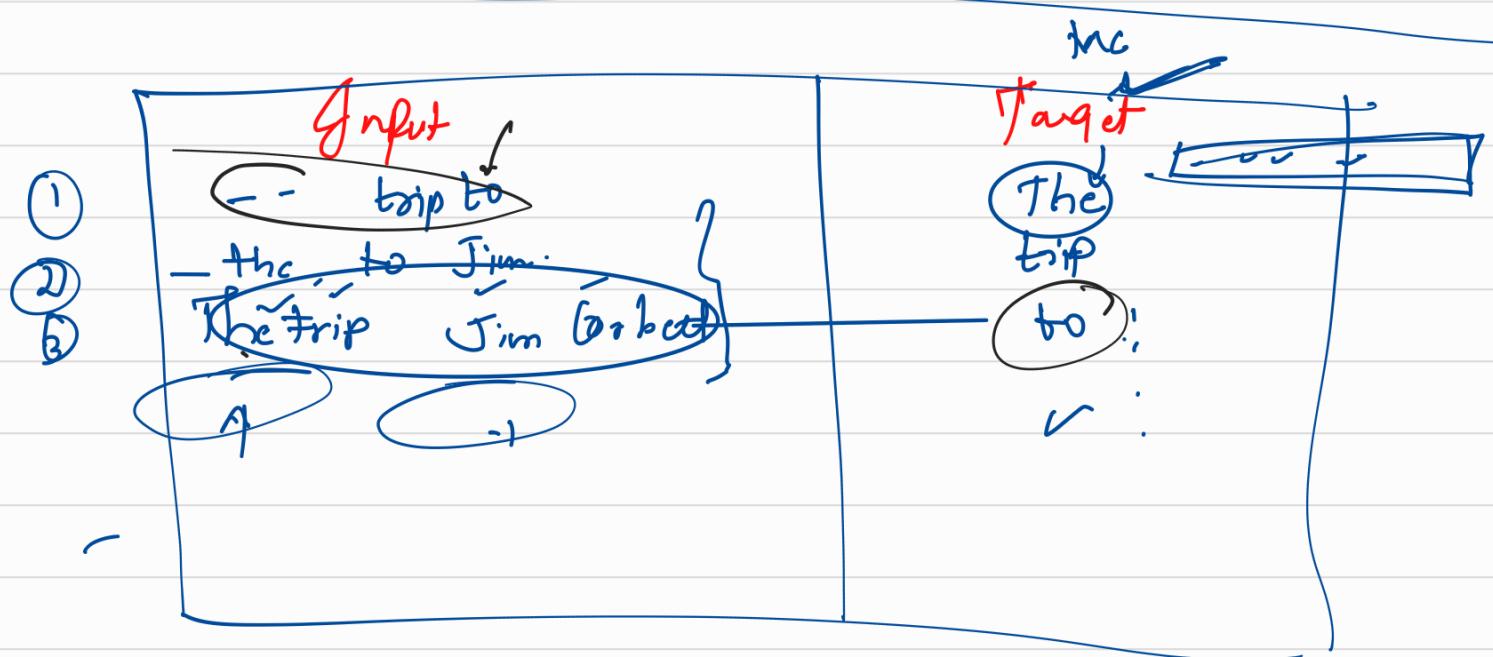
CBow (Continuous Bag of words model).

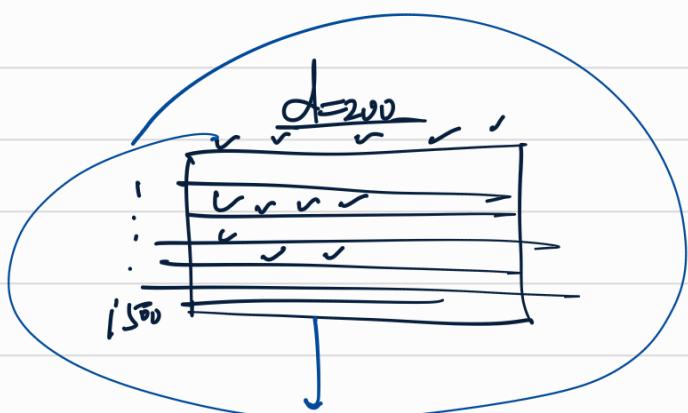
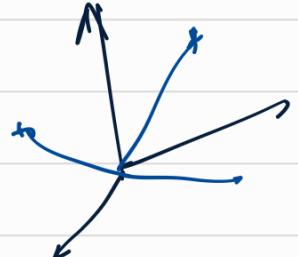
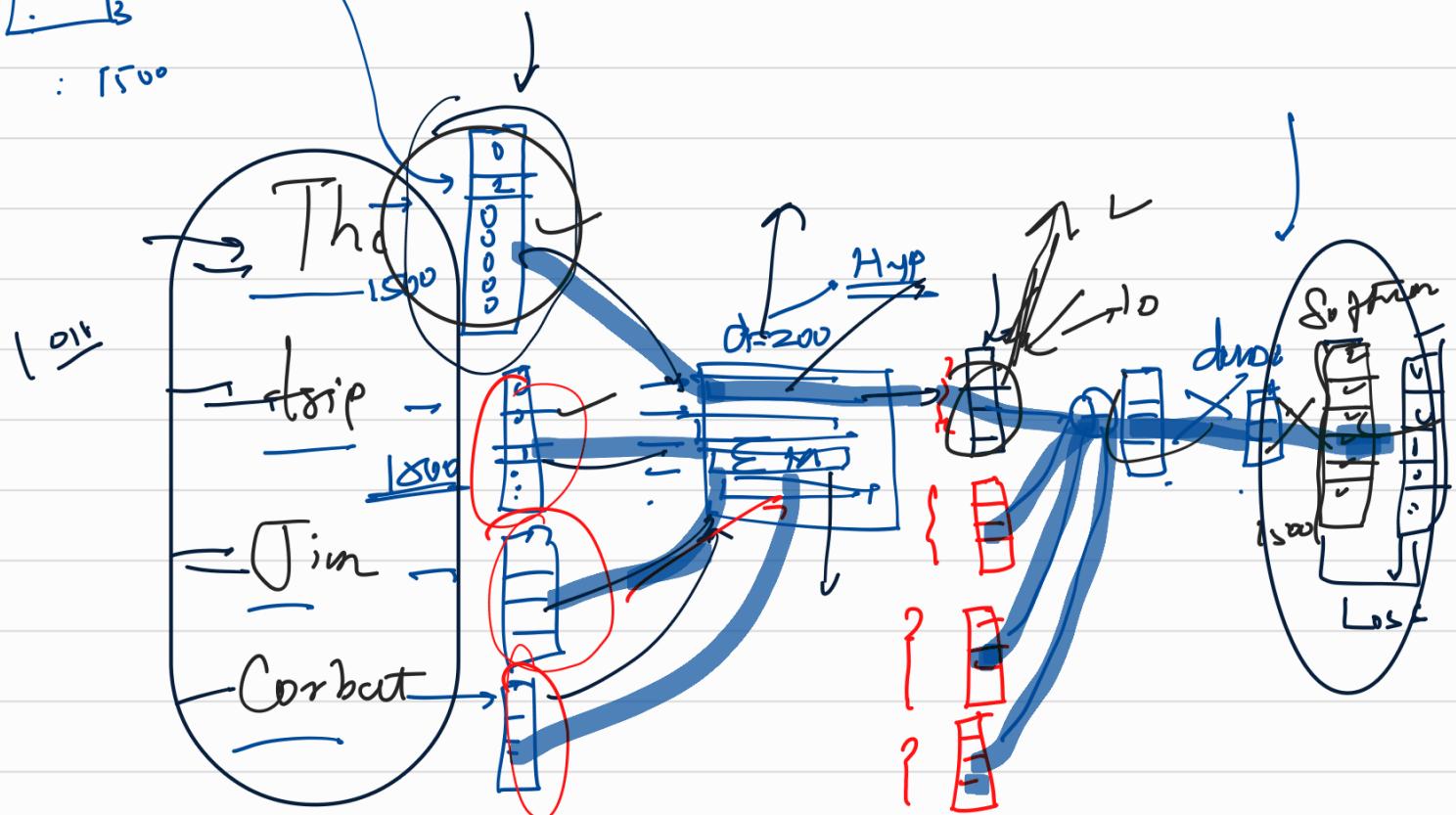
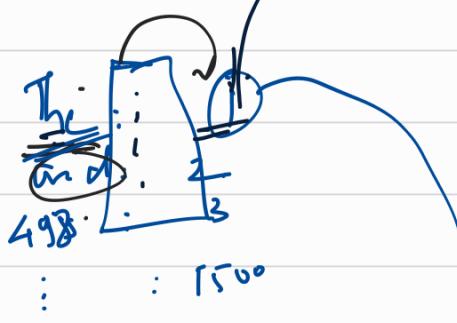
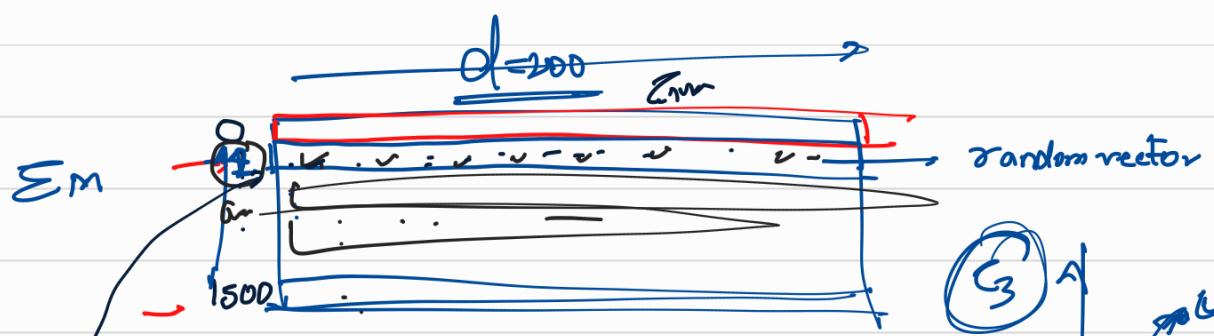
Skipgram model



ConFirst window = 2

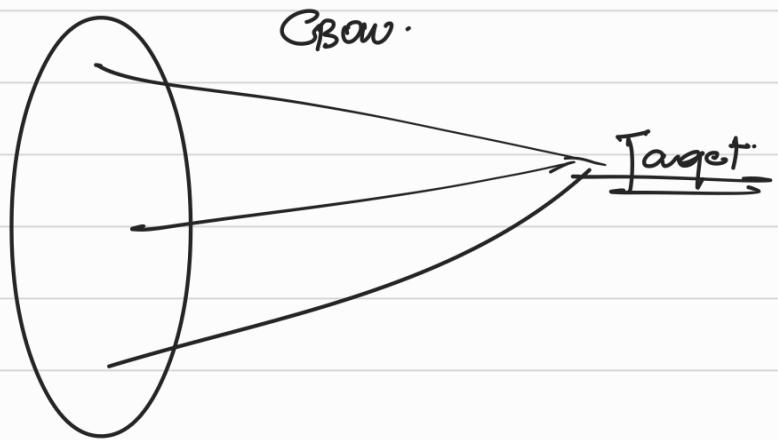
A hand-drawn diagram illustrating fluid dynamics around a ship's hull. The hull is shown as a blue line curving from the bottom left towards the top right. A horizontal line extends from the hull at the point where the curve begins, representing the free surface. Two arrows point away from the hull along this line, labeled '1000' above them. On the upper part of the hull, two circles are drawn, each with an arrow pointing towards the hull, labeled '100' above them. A large circle is drawn on the lower part of the hull, with an arrow pointing away from it, labeled '1500' above it. To the right of the hull, the text 'Sim looked was amazing. I saw 4 tigers' is written.



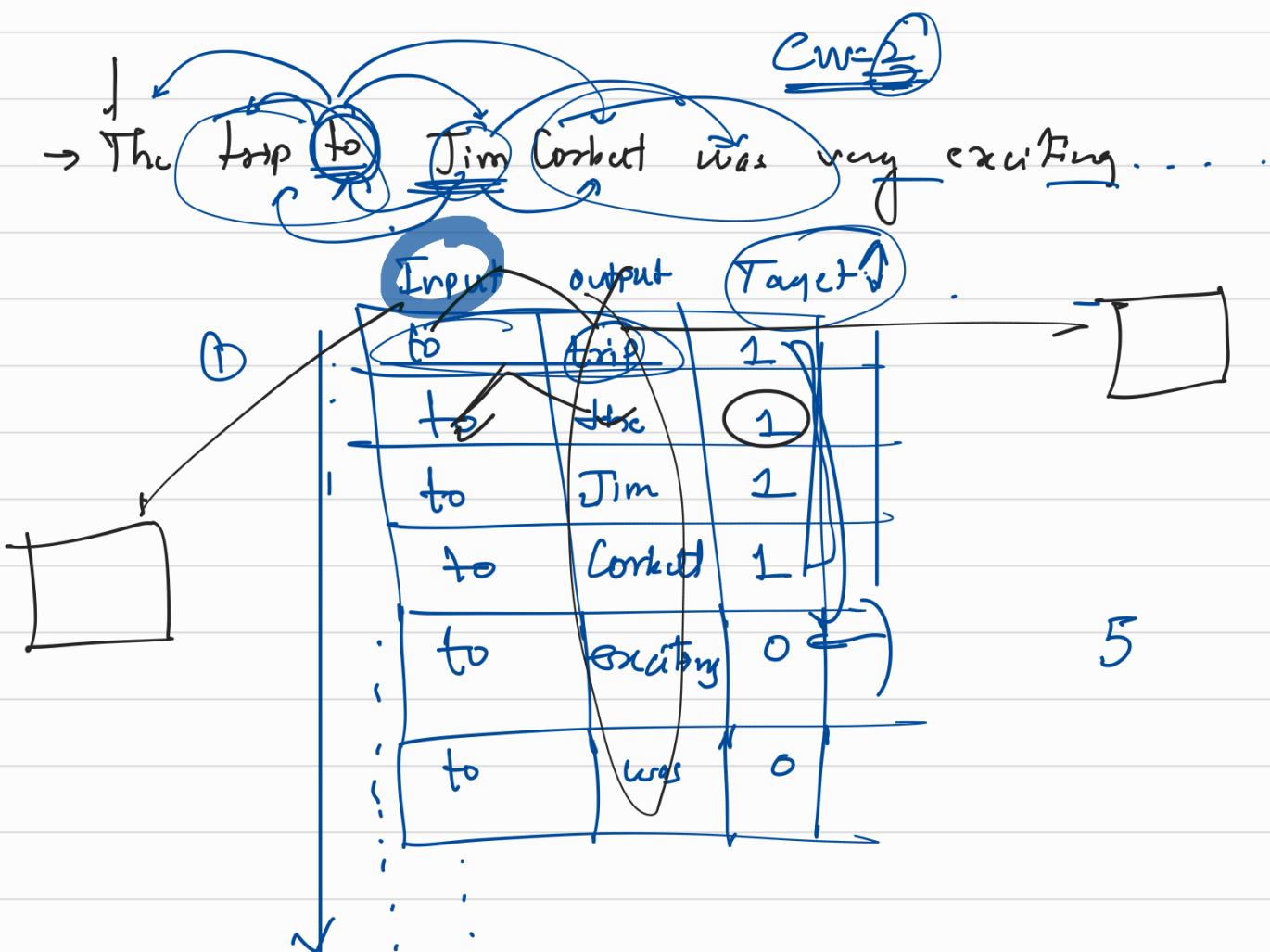
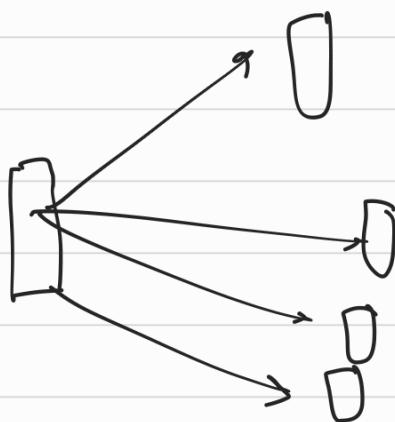


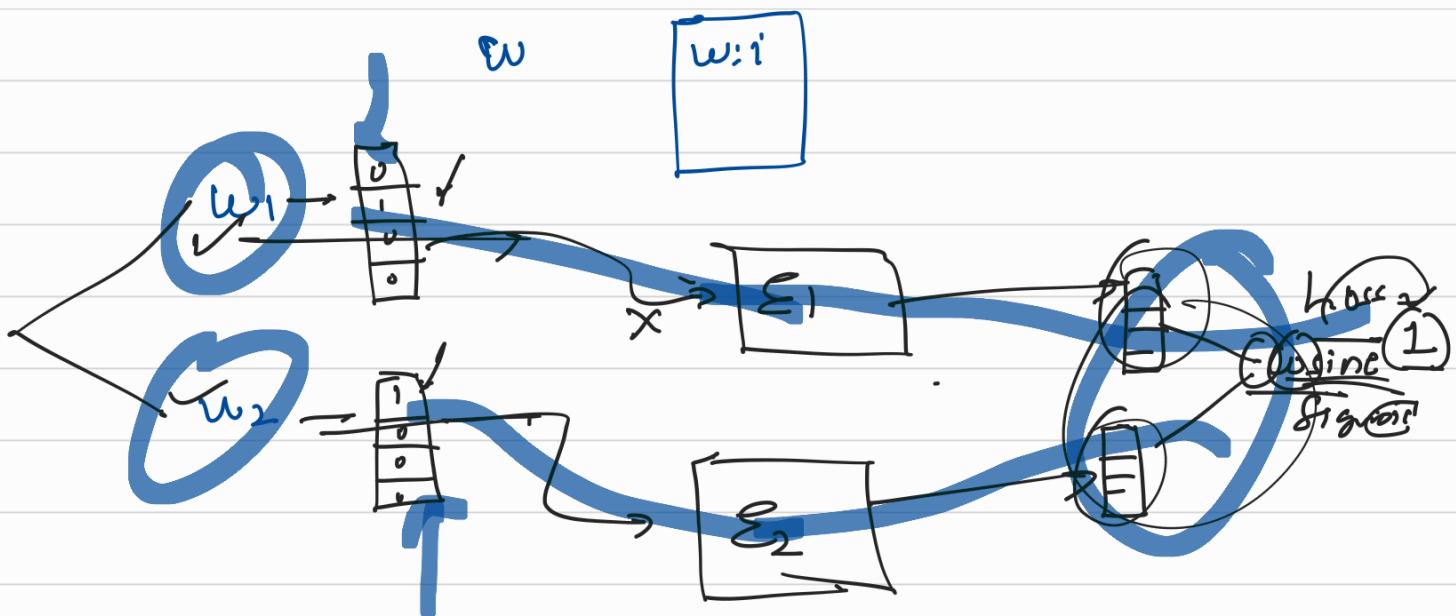
Sripgram

DM Jank



Skip-gram



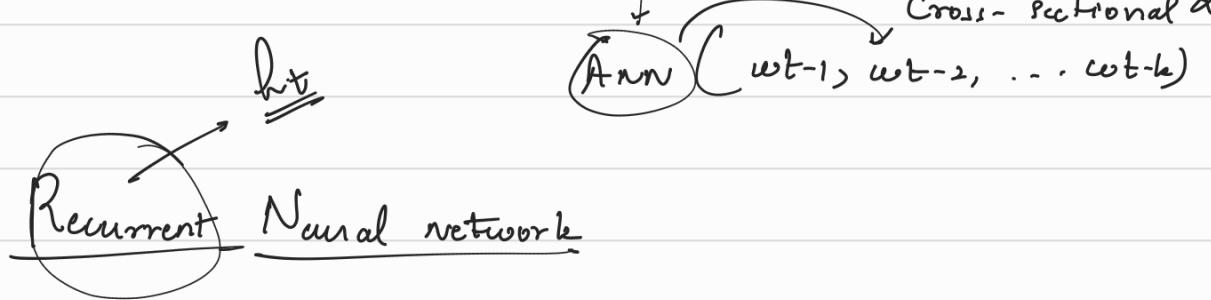


Recurrent Neural networks:

Tent \rightarrow Sequence

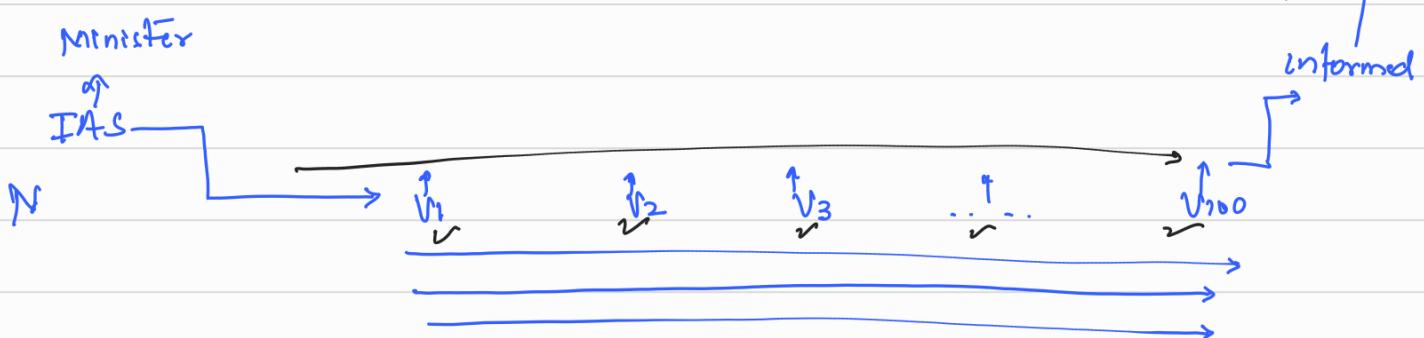
$$w_t \sim f(w_{t-1}, w_{t-2}, \dots, w_{t-k})$$

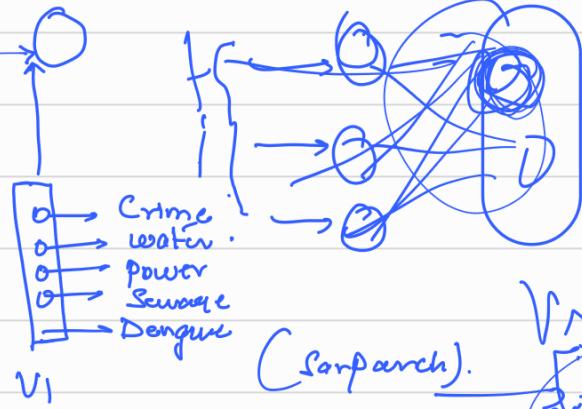
which is a neural network designed for 'sequence models'



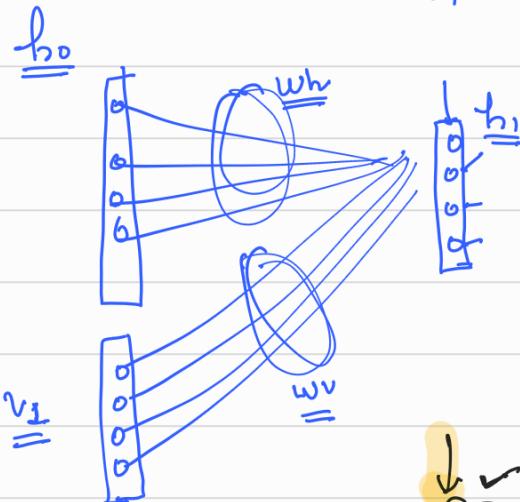
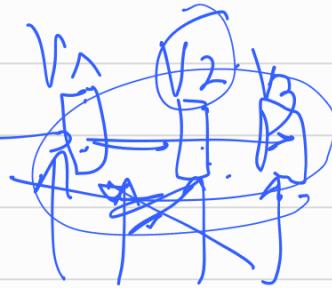
assume there are 100 villages in a town

Election - Yes
Informed



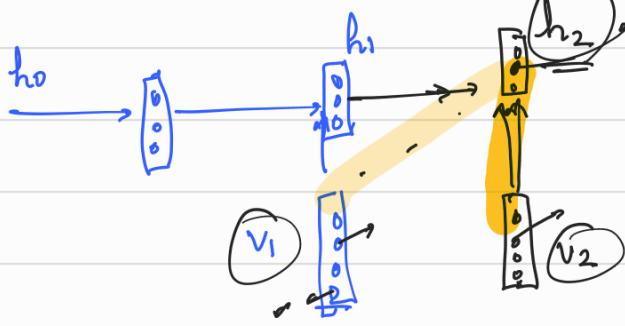


(Sarpanch).



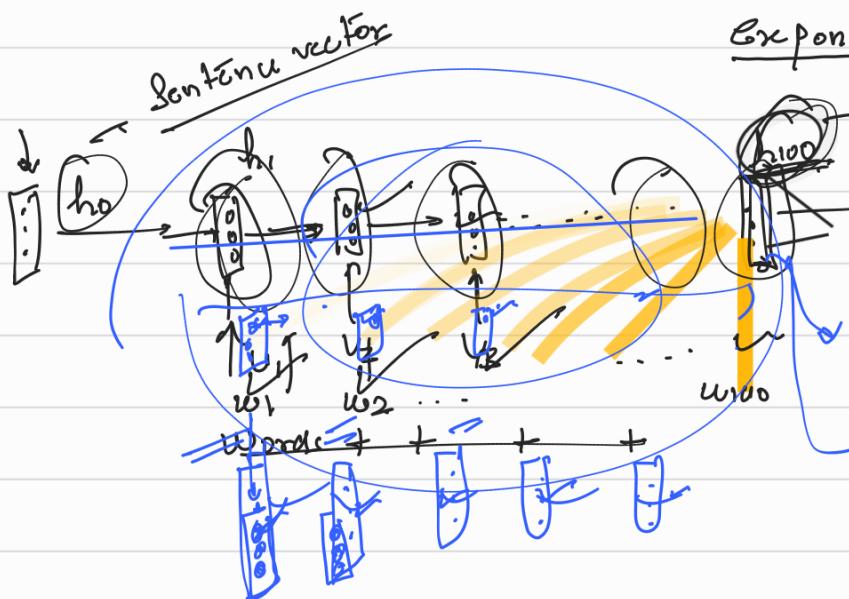
$$h_1 = \text{tanh}(\underline{w_h} \cdot h_0 + \underline{w_v} \cdot v_1 + b)$$

$$h_2 = \text{tanh}(\underline{w_h} \cdot h_1 + \underline{w_v} \cdot v_2 + b)$$



$$h_2 = \text{tanh}(w_h \cdot h_1 + w_v \cdot v_2 + b)$$

$$h_2 = \frac{w_h^2 h_0 + w_h w_v v_1 + b}{w_v^2} + \frac{w_h^2 h_0 + w_h w_v v_1 + b}{w_v^2} + b$$



Exponential moving average equation.

Sentence

It has information of all
denic signals
the villages.

have understanding
of all the important
aspects across 100 villages