

- ① Skip gram
- ② CBOW
- ③ Custom BOW
- ④ Sentence embedding (langchain, hugging, Sentence transform, transform)
- ⑤ RNN | LSTM | Transformer
- ⑥ LLM
- ⑦ API, huggingface

Word2Vec (NN)

=

- ① word → vector (Dense vector)
(Non-zero)
- ② Semantic info dim dimension

vector

=

(happy, joy)

(happy, sad)

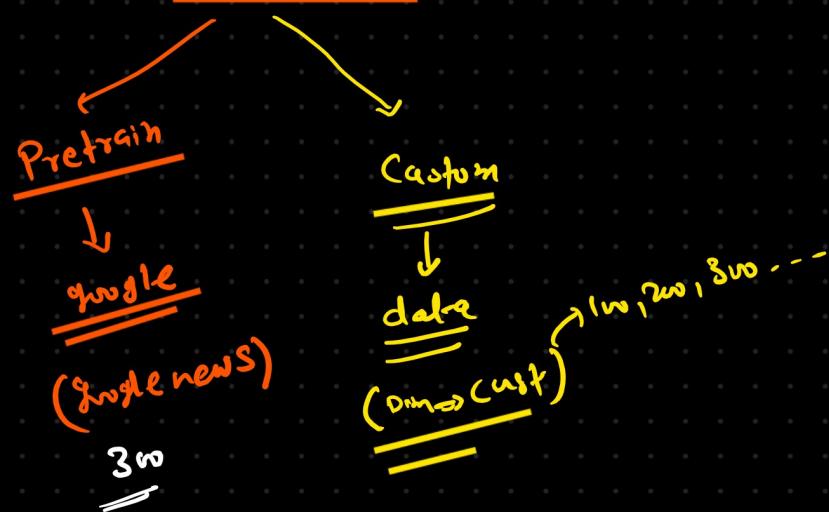
- ③ Low-Dim (300, 200, 1000, 1500)

⇒ computation

=

- ④ Dense (info)

Word2vec



[0-1]

	F ₁ gender	F ₂ Wealth	F ₃ Power	F ₄ Weight-	F ₅ speak
King	-1	1	1	0.8	1
Queen	-1	0.9	0.8	0.9	1
MAN	-1	0.5	0.2	0.8	1
WOMAN	-1	0.3	0.1	0.7	1
Monkey	-1	0	0	0.3	0

$$\begin{aligned}
 \text{Queen} &= \text{KING} - \text{MAN} + \text{WOMAN} \\
 \text{Queen} &= \begin{pmatrix} 1 \\ 1 \\ 0.5 \\ 1 \end{pmatrix} - \begin{pmatrix} 1 \\ 0.5 \\ 0.2 \\ 1 \end{pmatrix} + \begin{pmatrix} 0.3 \\ 0.1 \\ 0.7 \\ 1 \end{pmatrix}
 \end{aligned}$$

$$\begin{aligned}
 \text{King} &= \text{man} - (\text{MAN}) + \text{WOMAN}
 \end{aligned}$$

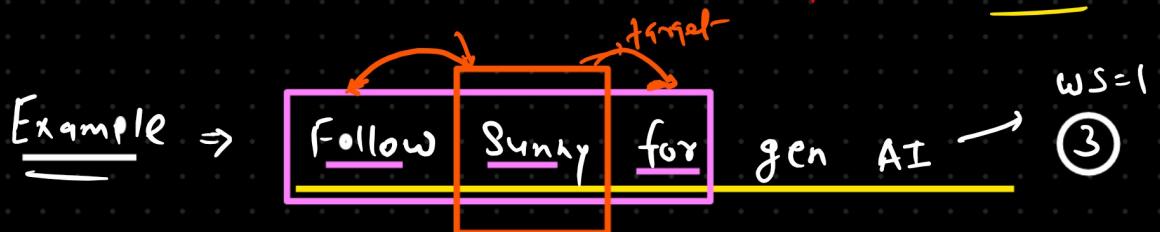
⑤ features

Type of wordvec

CBOW
(Continuous Bag of words)

Skip-gram

CBOW (Neural Network Knowledge is required) (Prerequisite)



Window Size $\Rightarrow 1, 2, 3, 4, 5, \dots, n$

\Rightarrow Dimension of vector $\Rightarrow 2, 3, 4, 5, \dots, n$

Window Size \Rightarrow 3 words continuously

WindowSize = 1

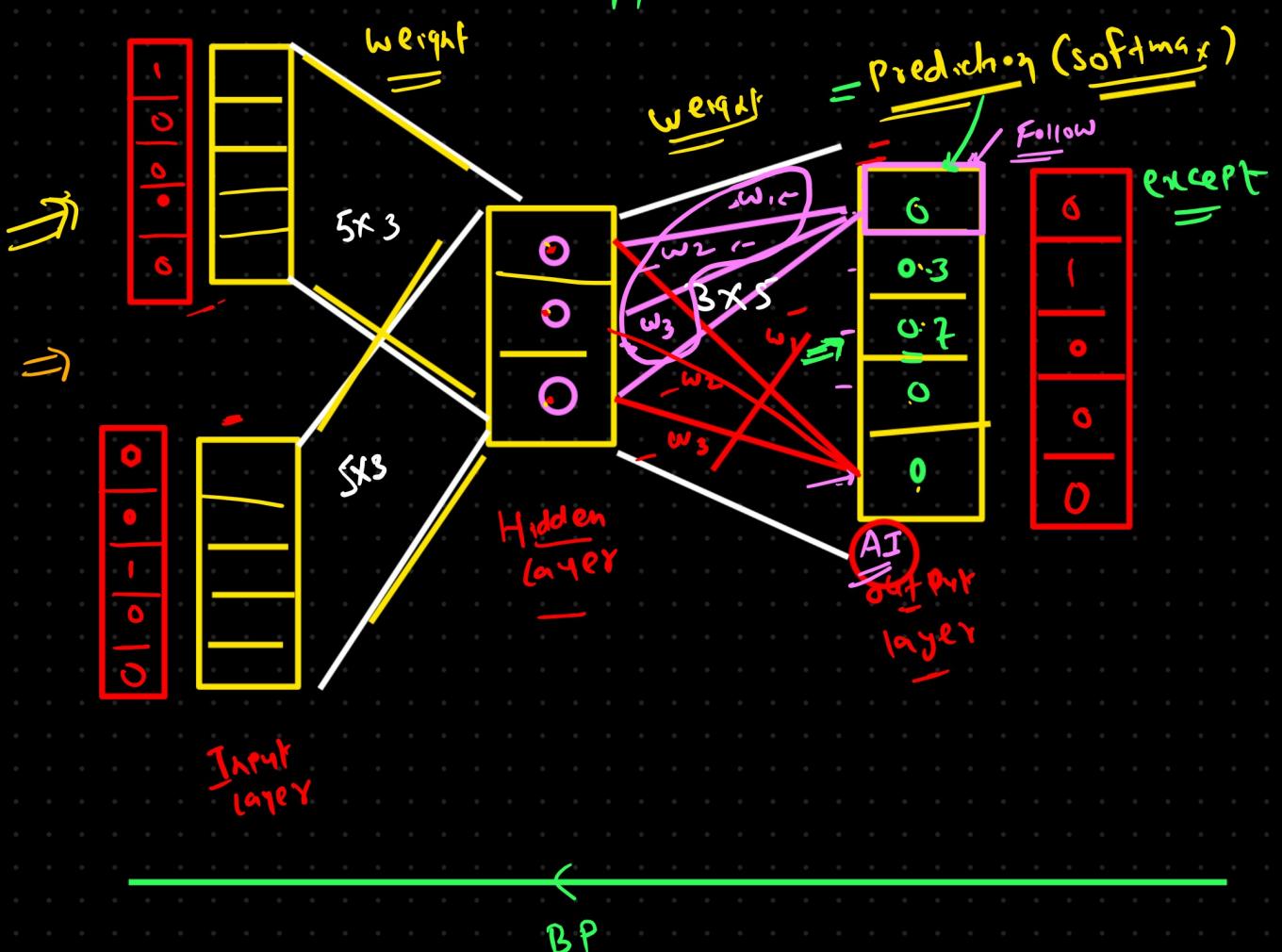
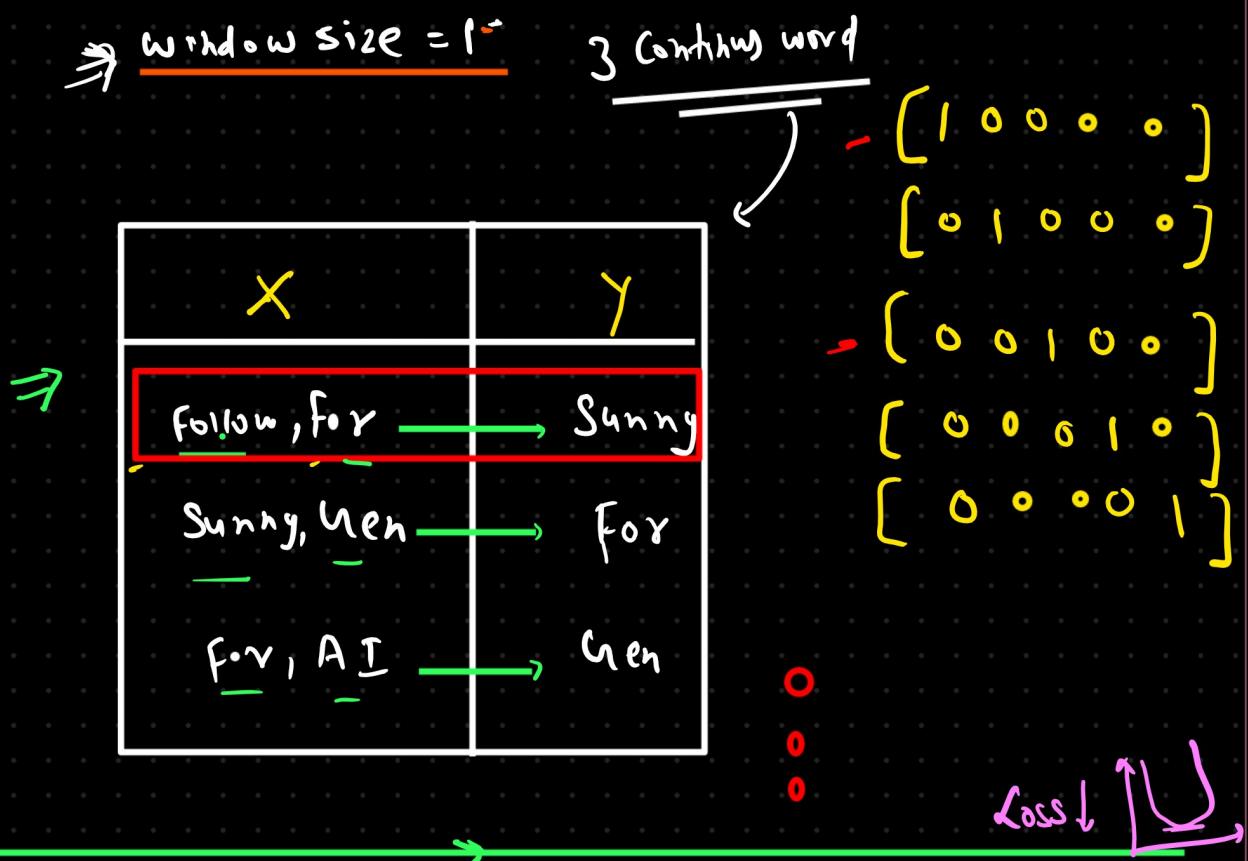
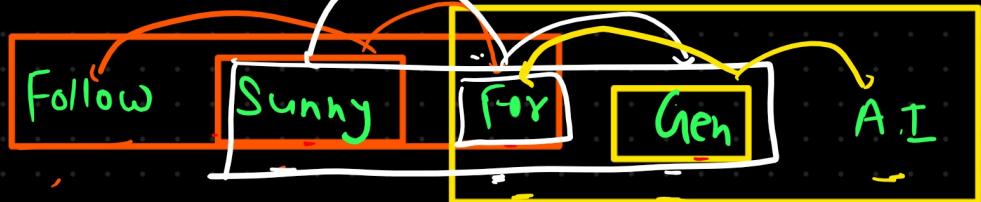
\Rightarrow 5 word cont.

WS \Rightarrow 2

\Rightarrow 7 words

WS \Rightarrow 3

Dimension of vector = 3

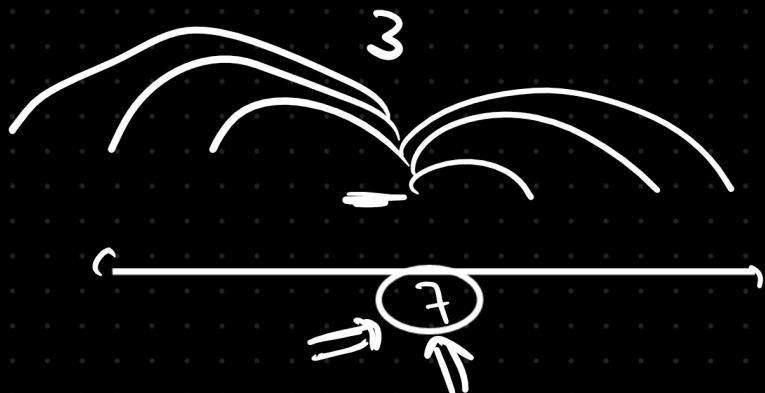


(Window Size)

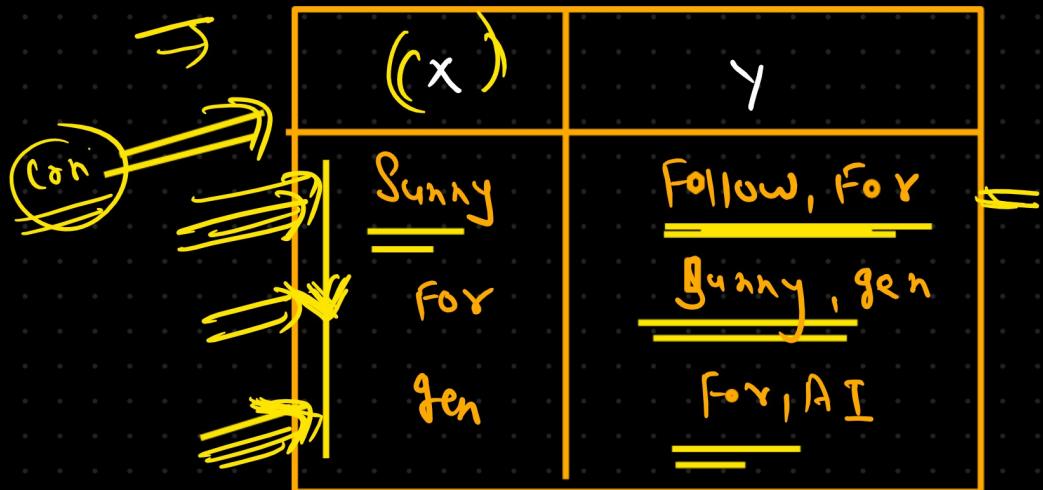
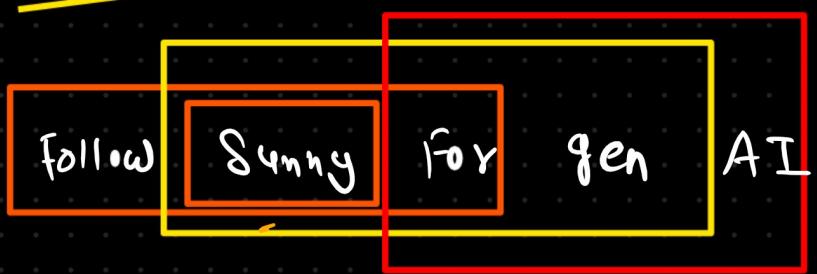
Follow Sunny Savita For Gen AI and ML
= - -

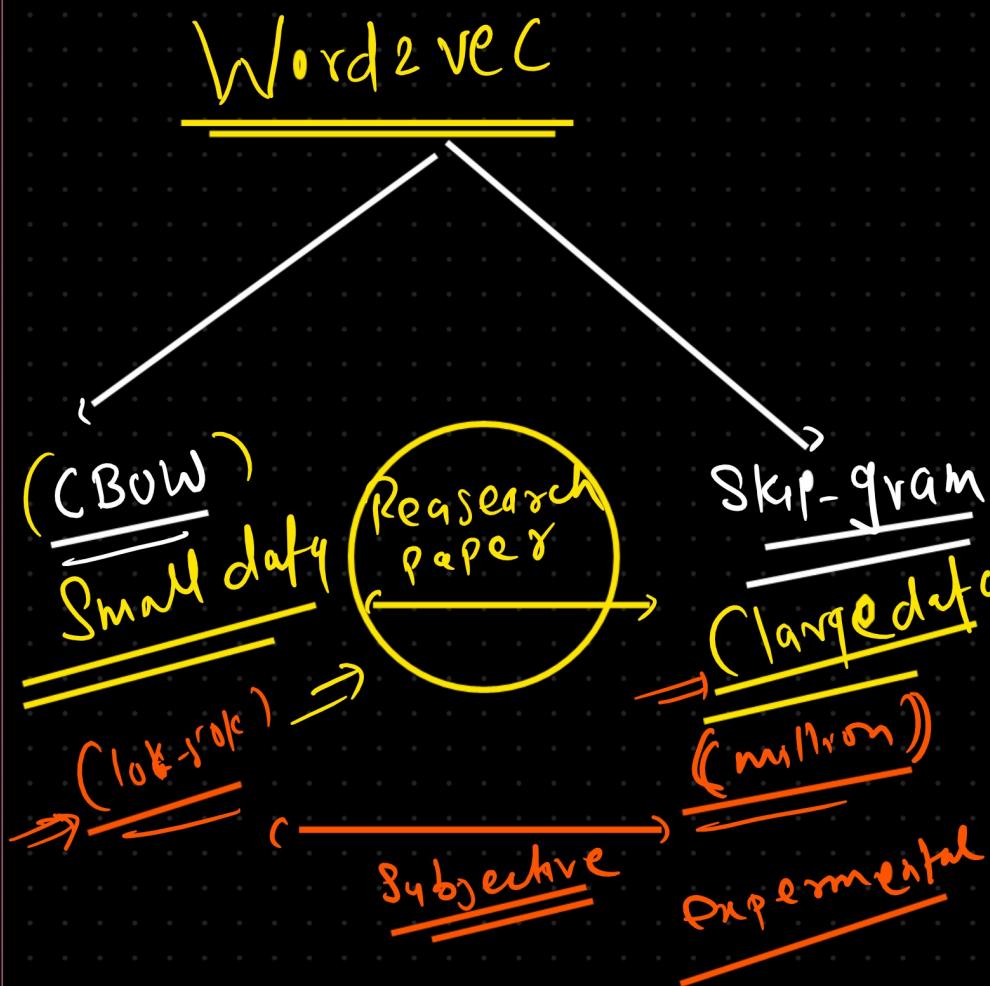
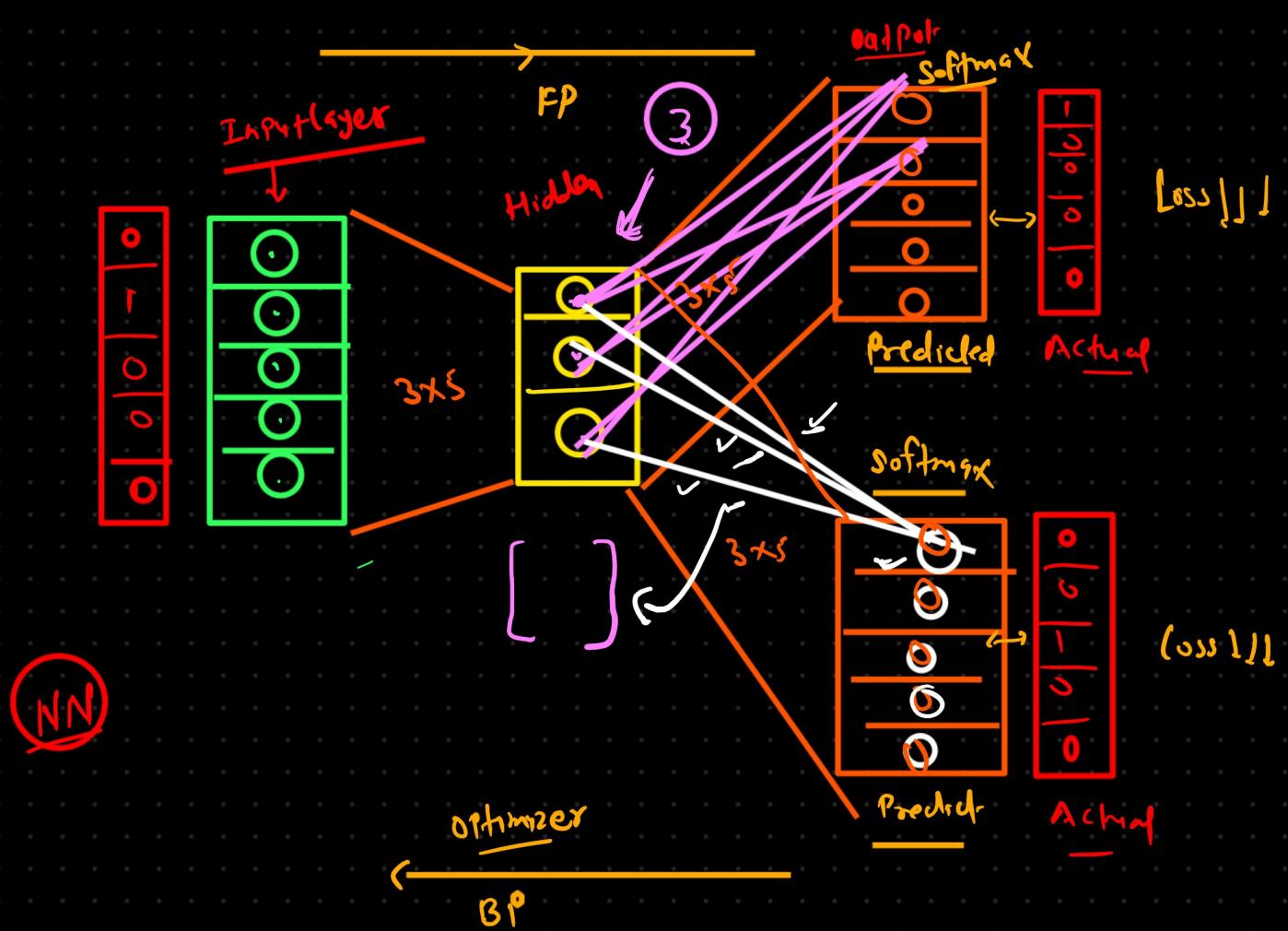
CBOW

Window Size = 3
Context word = 7

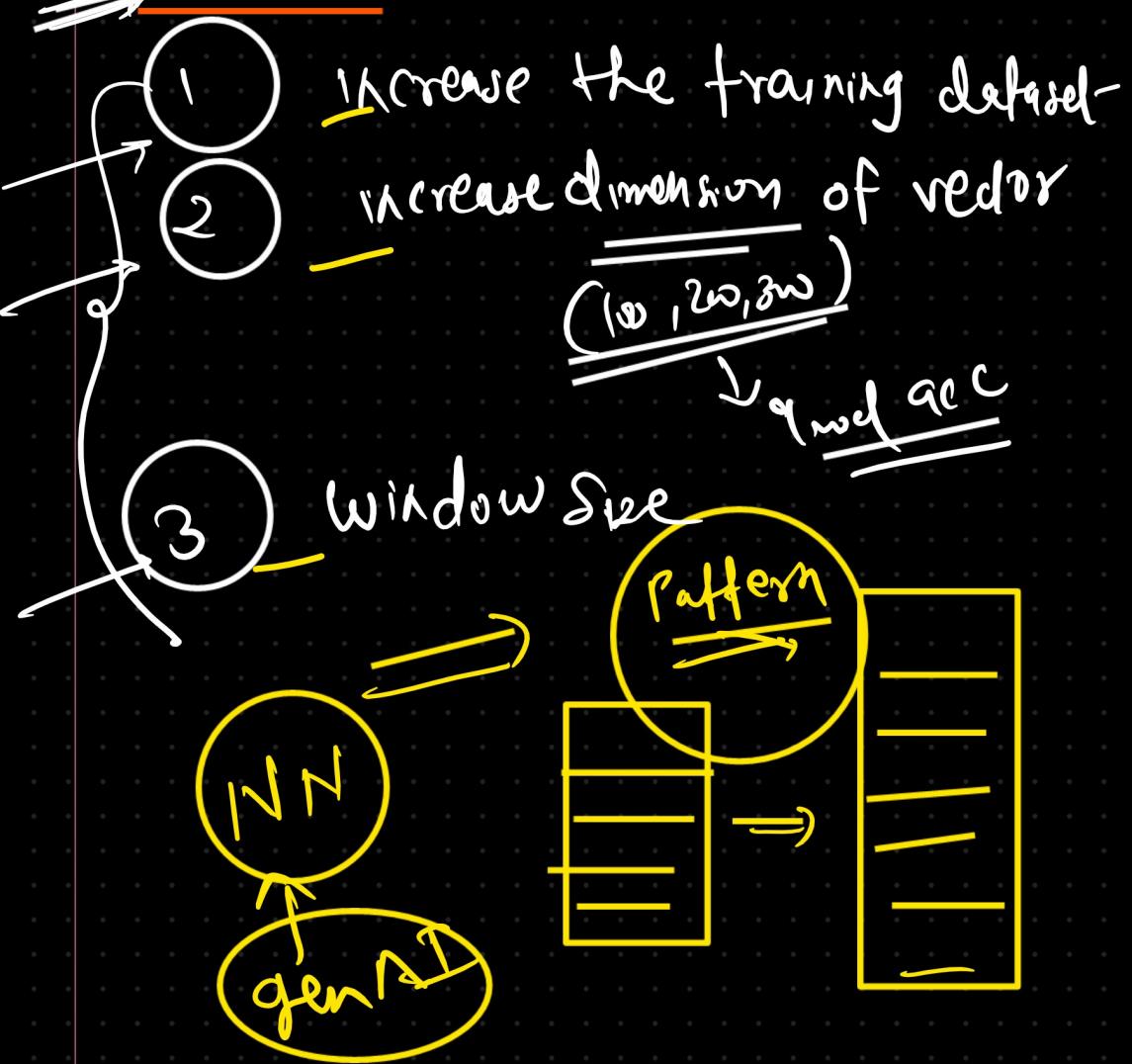


skipgram \Rightarrow Inverse of CBOW





Window C



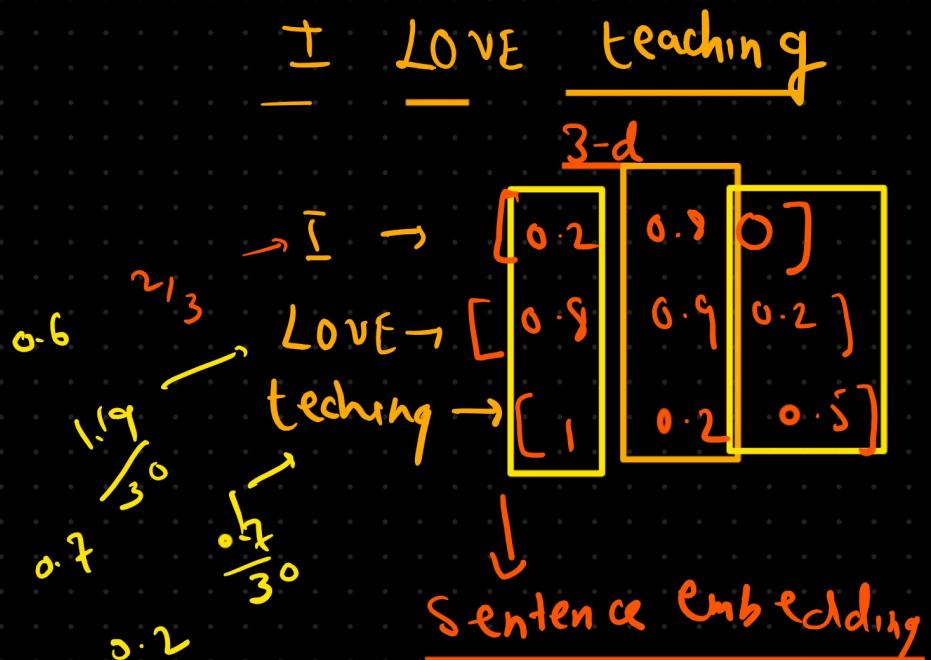
(word \rightarrow vector)



Sentences.

(collection of words)

\Rightarrow (Avg word2vec)



$$\frac{0.2 + 0.9 + 1}{3} = \frac{2}{3}$$

$$\frac{0.9 + 0.9 + 0.2}{3} = \frac{1.9}{3}$$

$$\frac{0.2 + 0.2 + 0.5}{3} = \frac{0.9}{3}$$

\Rightarrow (I LOVE teaching) = [0.6, 0.7, 0.2]

avg word2vec

Follow Sunny For gen AI

intermediate



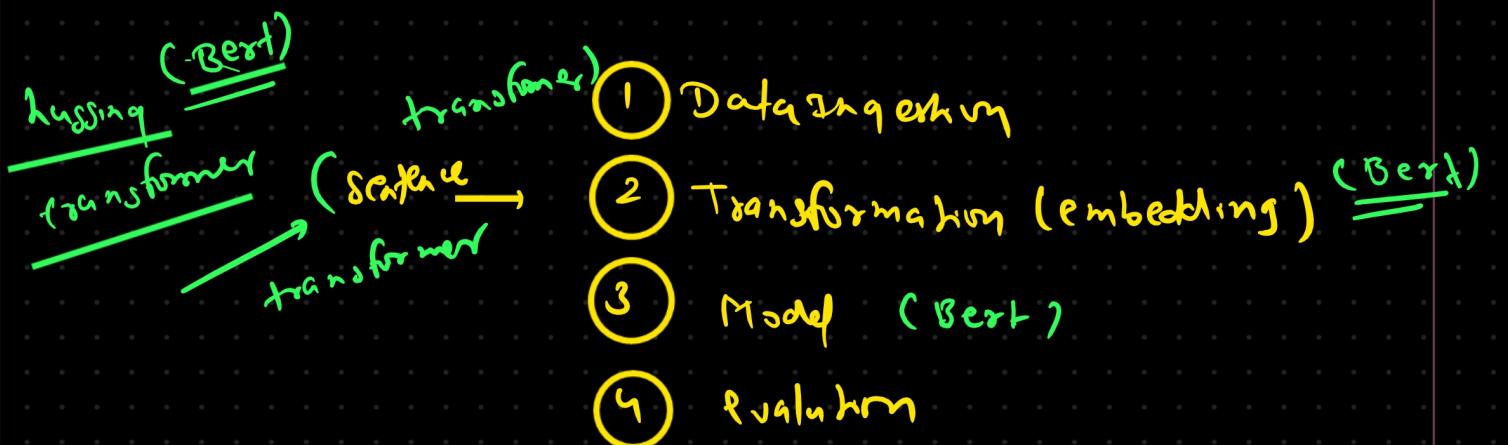
transformer

[]

(mid) A B C

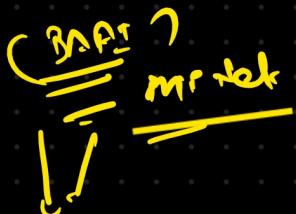
100.1.
50-60-1.
rest
BASE

mentor
student



- OpenAI → GPT
- OpenAI Embedding

1024



Dataset size: Larger datasets generally benefit from more powerful models like MPNet.

Computational resources: If you have limited resources, BGE Small En or MiniLM might be better options.

Task complexity: For complex tasks like question answering or text summarization, MPNet is often preferred.

Embedding dimensionality: Different models produce embeddings of varying dimensions. Choose based on downstream task requirements.

Performance vs. efficiency trade-off: Decide if you prioritize high accuracy or faster processing

Experimentation is key. Try different models and evaluate their performance on your specific task and dataset to find the best fit.

Preprocessing
(embedding)

model

