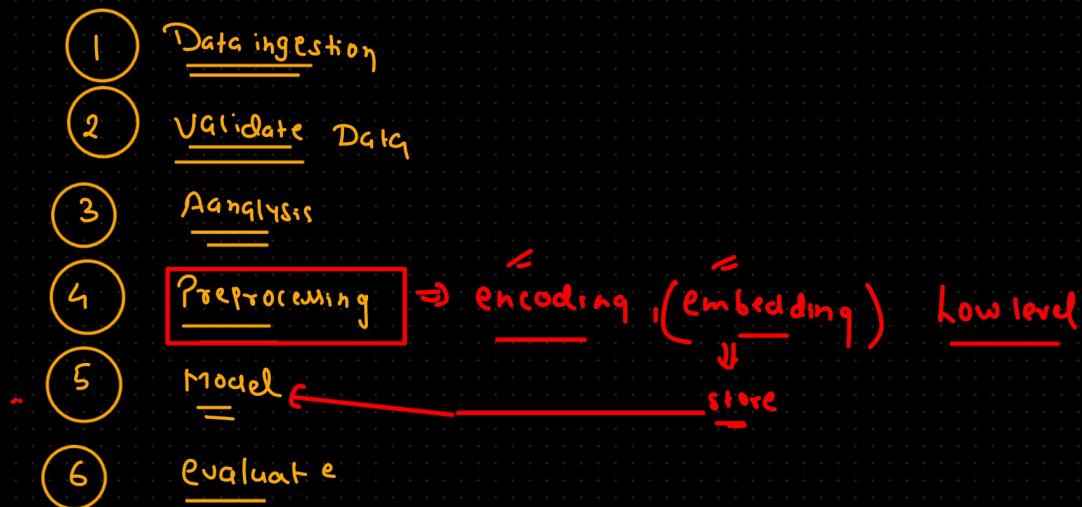


Next → (embedding) → word2vec
(Neural Network)

Transformer

- Similarity
- 1 Dot Product
- 2 Cosine Similarity }



- 1 Corpus → entire text - Paragraph
- 2 Document → Sentence
- 3 token → word

Vocabulary

Unique Collection
of word (token)

Corpus
Doc

Paragraph Corpus

Improvements in transformer-based deep neural networks enabled an AI boom of generative AI systems in the early 2020s. These include large language model (LLM) chatbots such as ChatGPT, Copilot, Bard, and LLaMA, and text-to-image artificial intelligence art systems such as Stable Diffusion, Midjourney, and DALL-E.[7][8][9] Companies such as OpenAI, Anthropic, Microsoft, Google, and Baidu as well as numerous smaller firms have developed generative AI models.[3][10][11]

OHE

CORPUS

= { People watch inuron neuron watch inuron
People write comment- inuron write comment }

- D₁

People Watch inuron

Sunny

$\Rightarrow 4 \times 6$

\Rightarrow padding

- D₂

inuron Watch neuron $\Rightarrow 3 \times 6$

- D₃

People Write Comment $\Rightarrow 3 \times 6$

- D₄

neuron write comment- $\Rightarrow 3 \times 6$

Document

(000000)
(---)
[---]

testing

training

(000000)

50000
Sunny

Vocabulary

{ People, watch, inuron, write, comment- }

V=5

OHE

People	Watch	inuron	write	comment-
--------	-------	--------	-------	----------

5

$$D_1 = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \end{bmatrix} (3 \times 5)$$

$$D_4 = \begin{bmatrix} [00100] \\ [00010] \\ [00001] \end{bmatrix} (3 \times 5)$$

most of the values
zero

Pros

- Easy to implement.

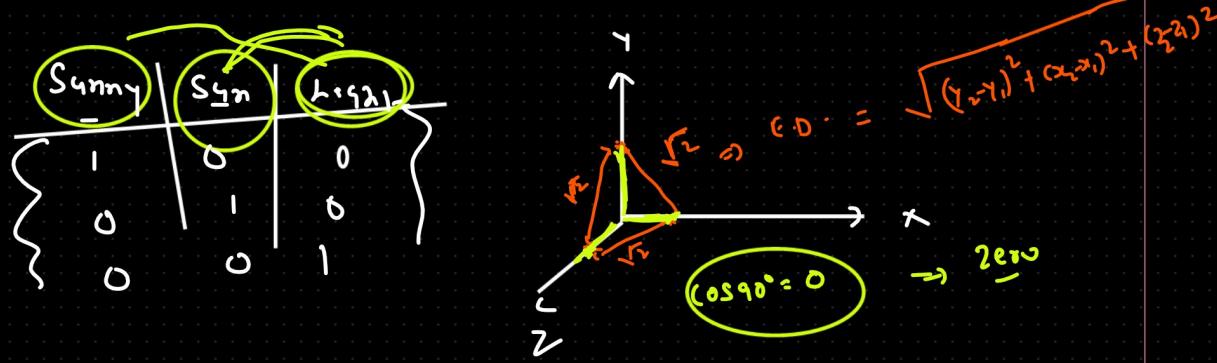
Cons

1 Sparsity (sparse matrix)

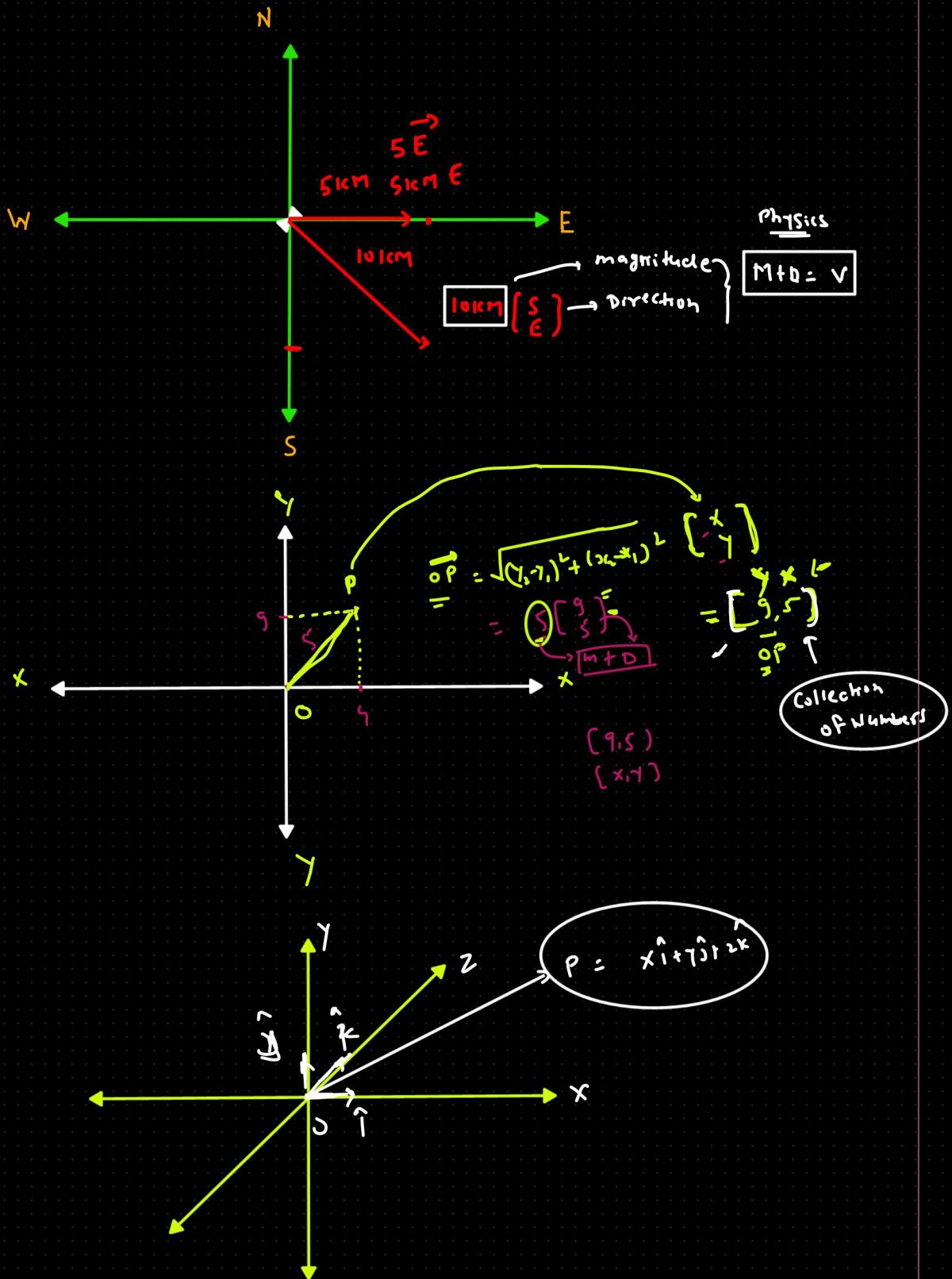
2 No fixed size

3 OOV

4 Content will not be there



Vector



Beg of word

Watch Sunny Sun

D₁ - People Watch inuron write comment

D₂ - inuron Watch inuron $\Rightarrow 3 \times 6$

D₃ - People write Comment $\Rightarrow 3 \times 6$

D₄ - inuron write comment $\Rightarrow 3 \times 6$

trainin

	People	watch	inuron	write	comment
D ₁	1	2	1	0	0
D ₂	0	1	2	0	0
D ₃	1	0	0	1	1
D ₄	0	0	1	1	1

Adv

Simple inflection

Conform

What you are doing?

Dis adv.

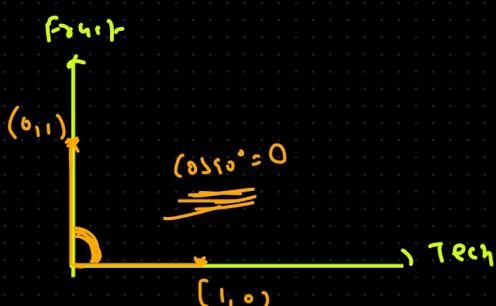
- 1 Sparsity still is there
- 2 Order
- 3 GOV (Dim. issue)

~~watch inuron~~

Apple iPhone Apple

\Rightarrow

Tech	Fruit
1	0
0	1



Dot product

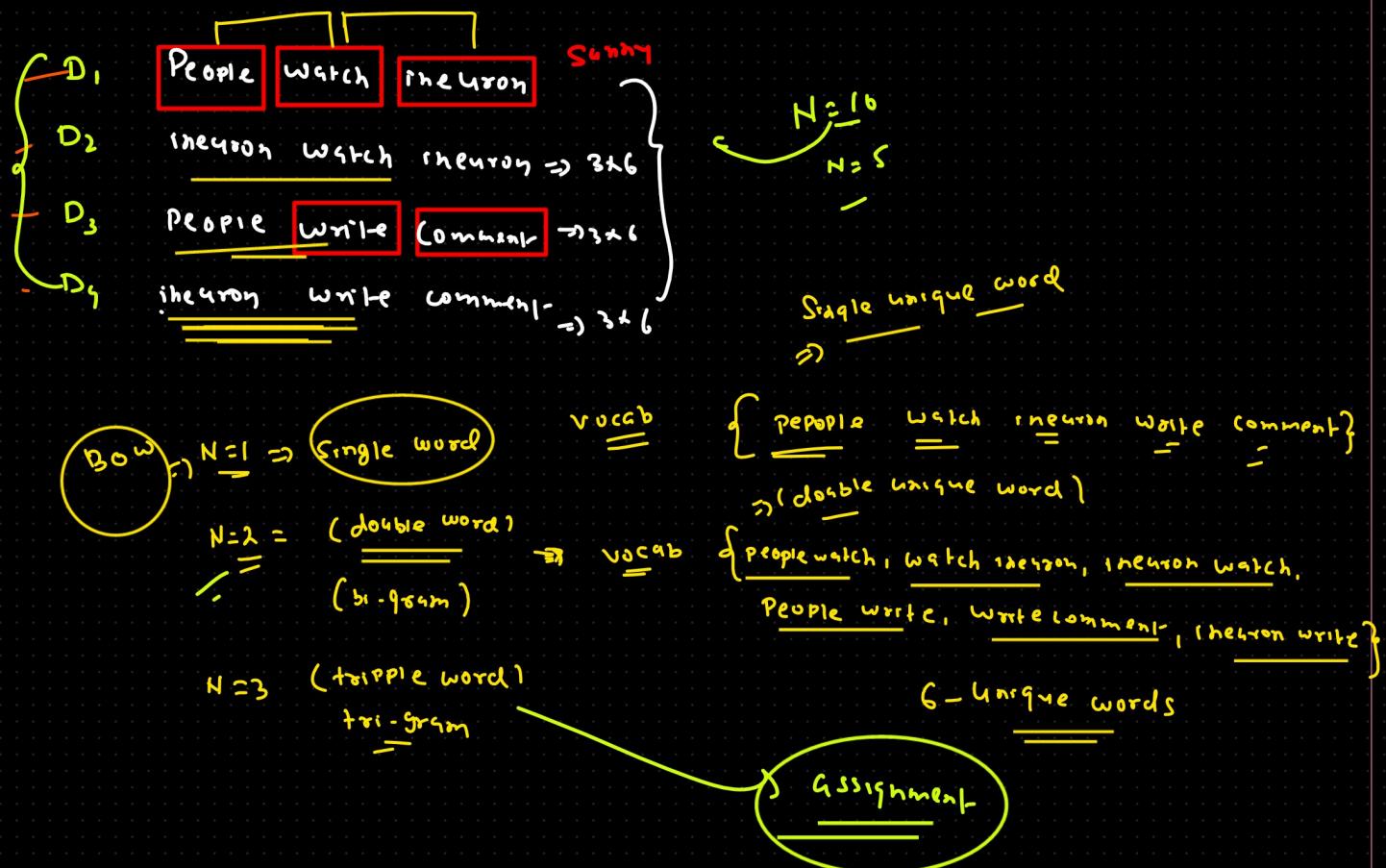
$$(x_1, y_1) \cdot (x_2, y_2)$$

$$(x_1 + x_2) + (y_1 + y_2)$$

$$(0,1) \cdot (1,0)$$

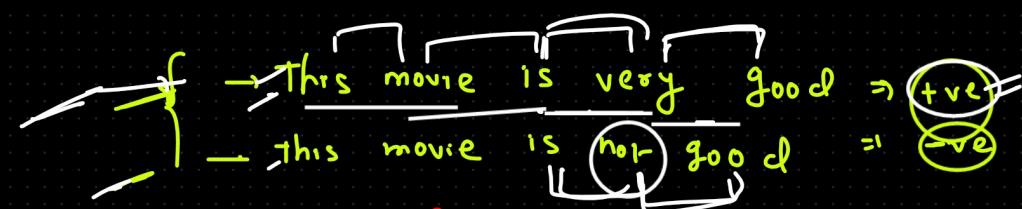
$$0 \times 1 + 1 \times 0 = 0$$

N-gram

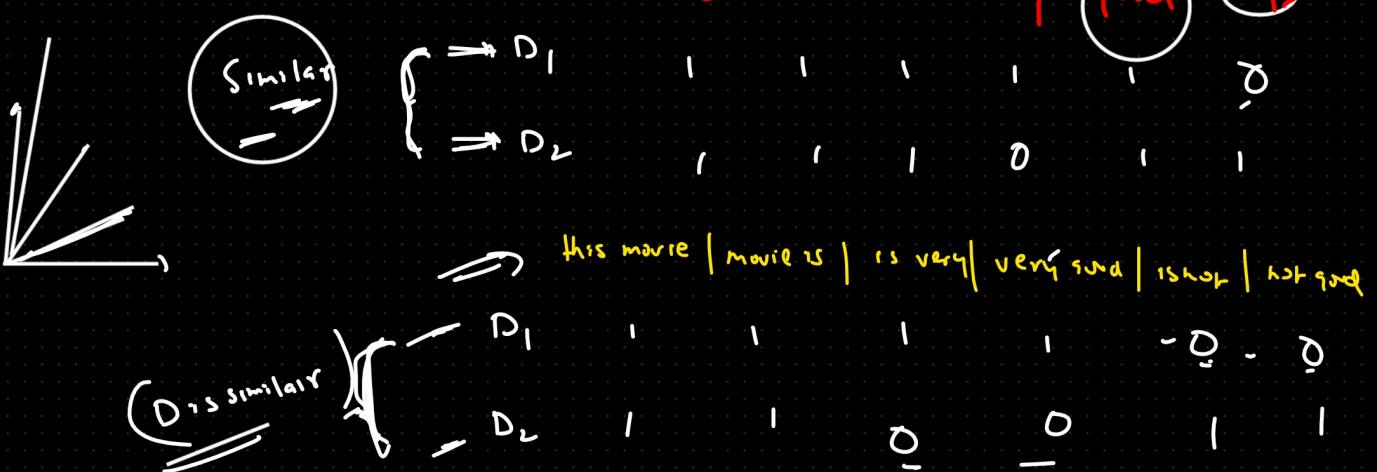


	Peoplewatch	watchInuron	neuronwatch	people write	write comment	Inuron write
D ₁	1	1	0	0	0	0
D ₂	0	1	1	0	0	0
D ₃	0	0	0	0	1	0
D ₄	0	0	0	0	1	1

Ex:- Capturing more context, Reducing the ambiguity



{ this movie very good not? }



4 TF-IDF

TF → term frequency

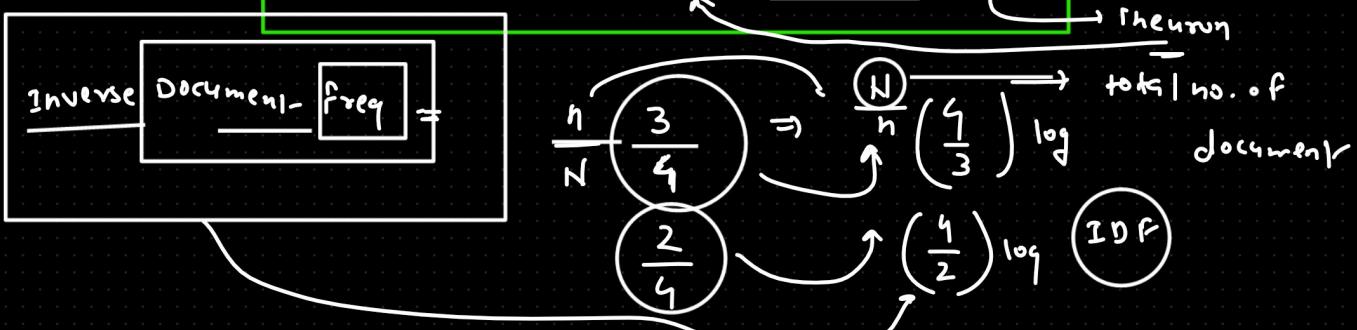
IDF → inverse document frequency

$TF \rightarrow D_1$	$\boxed{\text{People}}$	$\boxed{\text{watch}}$	$\boxed{\text{neuron}}$	Sunny	$\left(\frac{1}{3}\right) \cdot \left(\frac{1}{3}\right) \cdot \left(\frac{1}{3}\right)$
$\rightarrow D_2$	$\boxed{\text{neuron}}$	$\boxed{\text{watch}}$	$\boxed{\text{neuron}}$	$\Rightarrow 3 \times 6$	$\left(\frac{2}{3}\right) \cdot \left(\frac{1}{3}\right) \cdot \left(\frac{2}{3}\right)$
$\rightarrow D_3$	$\boxed{\text{People}}$	$\boxed{\text{write}}$	$\boxed{\text{Comment}}$	$\Rightarrow 3 \times 6$	$\left(\frac{1}{3}\right) \cdot \left(\frac{1}{3}\right) \cdot \left(\frac{1}{3}\right)$
$\rightarrow D_4$	$\boxed{\text{neuron}}$	$\boxed{\text{write}}$	$\boxed{\text{comment}}$	$\Rightarrow 3 \times 6$	$\left(\frac{1}{3}\right) \cdot \left(\frac{1}{3}\right) \cdot \left(\frac{1}{3}\right)$

$$TF(\text{word}, D) = \frac{\text{No of occurrence of term in given } D}{\text{total term in the document } (D)}$$

$$TF(\text{neuron}, D_1) = \frac{1}{3}$$

$$IDF = \log \left(\frac{\text{Total No. of Corpus in Document}}{\text{No. of documents with the term}} \right)$$



- D₁ People watch neuron Sunny
- D₂ neuron watch neuron $\Rightarrow 3 \times 6$
- D₃ people write comment $\Rightarrow 3 \times 6$
- D₄ neuron write comment $\Rightarrow 3 \times 6$

$TF = \frac{\text{freq of given term in doc}}{\text{total term in doc}}$

$IDF = \log \left(\frac{\text{total document}}{\text{Doc No. which contains term}} \right)$

Alphabetical

	(TF × IDF)				
(People)	1/3 × log(4)	1/3 × log(4)	1/3 × log(4)	0	0
D ₁	1/3 × log(4)	1/3 × log(4)	1/3 × log(4)	0	0
D ₂	0	1/3 × log(4)	2/3 × log(4)	0	0
D ₃	1/3 × log(4)	0	0	1/3 × log(4)	1/3 × log(4)
D ₄	0	0	1/3 × log(4)	1/3 × log(4)	1/3 × log(4)

$$TF = ? \Rightarrow \underline{\text{people}} \underline{\text{watch}} \underline{\text{neuron}} \underline{\text{people}}, \underline{\text{people}}$$

\downarrow \downarrow \downarrow \downarrow

$1/3$ $1/3$ $3/3$

D₁ people watch neuron $0.33 \quad (0-1) \Rightarrow 0.6$

D₂ people comment neuron $1DF = TF \times IDF$

$$DF = \frac{x}{x} = 1 \rightarrow \frac{1}{1} = 1 \leftarrow$$

$$DF = \frac{1}{2} = 0.5 = \frac{10}{20} = 0.5 \leftarrow IDF$$

$$DF = \frac{2}{2} = 1 = 1 \leftarrow$$

