

① frequency based (encoding)

stemming

② Neural Network based (embedding)

lemmatization

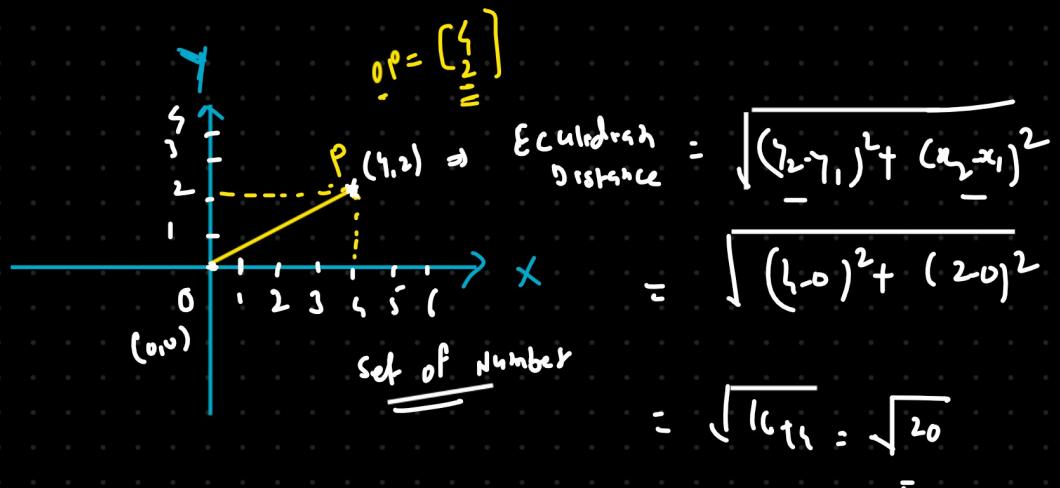
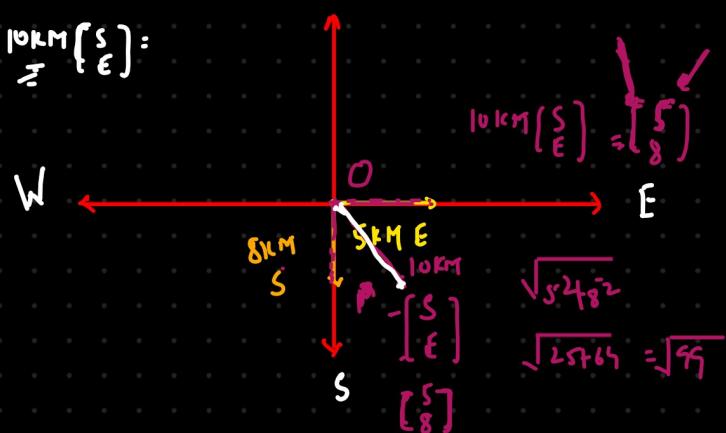
Vector :-

mathematically \Rightarrow magnitude + direction

$$\overrightarrow{OP} = \sum_{i=1}^N s_i e_i$$

vector :-

$$\overrightarrow{OP} = \begin{bmatrix} s_1 \\ s_2 \\ \vdots \\ s_N \end{bmatrix}$$



$$OP = \sqrt{20} \begin{bmatrix} x \\ y \end{bmatrix}$$

$$\overrightarrow{OP} = \sqrt{20} \begin{bmatrix} 4 \\ 2 \end{bmatrix}$$

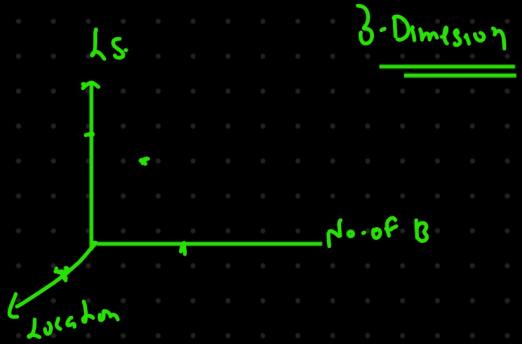
$$\therefore \overrightarrow{OP} = \hat{x} + \hat{y} + \hat{z}$$



House Price

	Number of bedrooms	housesize	Location
Set of No:	5	1000	B
1-D array	6	1200	D
	8	1600	M

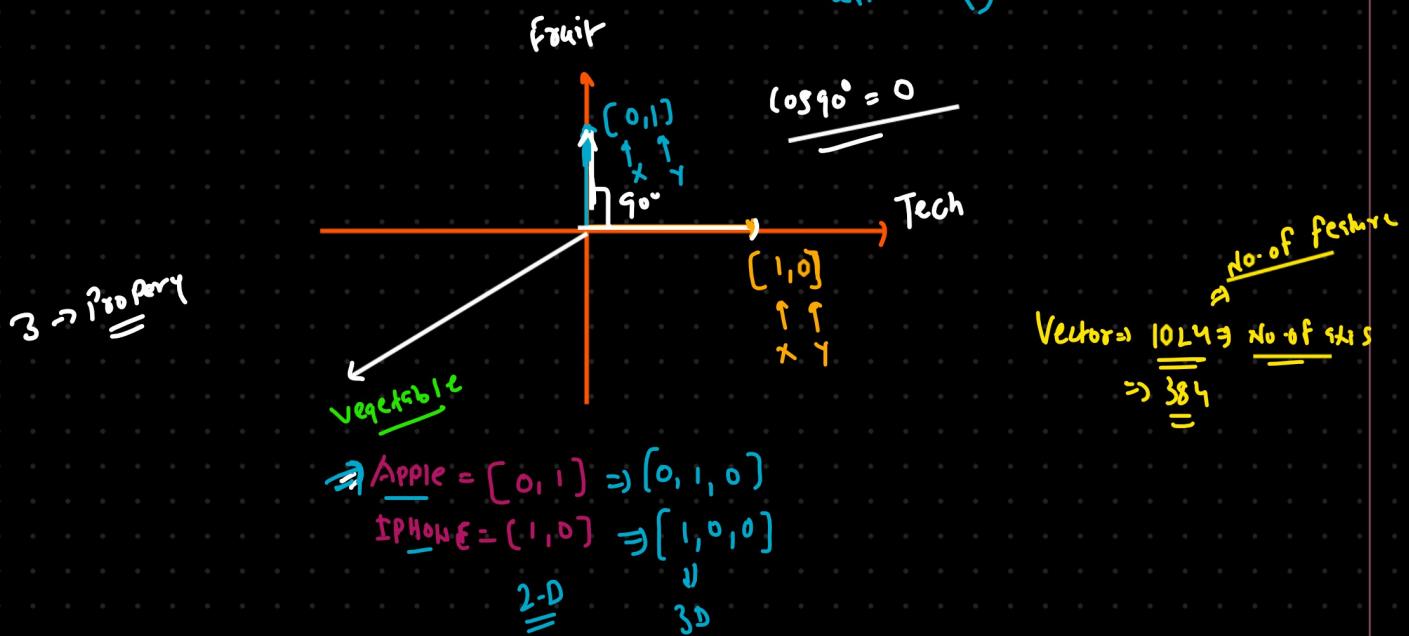
Price
sol



Matrix

Scalability search

Apple \Rightarrow
apple \Rightarrow



Apple \Rightarrow iPhone

$[0, 1] \cdot [1, 0]$

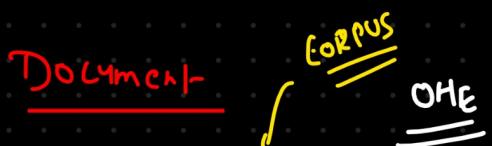
Dot Product \Rightarrow

$$\underline{0 \times 1 + 1 \times 0 = 0}$$

Cosine Similarity

$$= \frac{\bar{A} \cdot \bar{B}}{|A| \cdot |B|} = \frac{0}{\square} = 0$$

- 1 OHE
 - 2 BOW
 - 3 TF-IDF
 - 4 N-grams
- $\left. \begin{array}{l} \\ \\ \\ \end{array} \right\}$ frequency

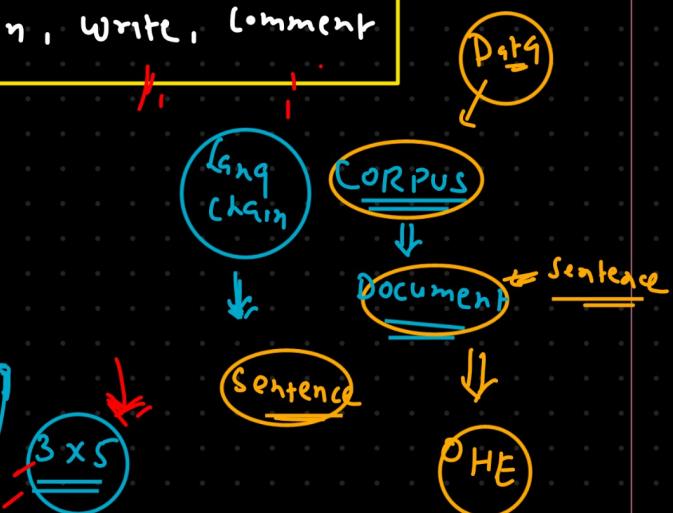


$\left\{ \begin{array}{l} \text{People watch neuron neuron watch theory} \\ \text{People write comment theory write comment} \end{array} \right\}$

$$\begin{aligned} D_1 &= \boxed{\text{people}} \quad \boxed{\text{watch}} \quad \boxed{\text{theory}} \\ D_2 &= \text{neuron} \quad \text{watch} \quad \text{neuron} \\ D_3 &= \text{people} \quad \boxed{\text{write}} \quad \text{comment} \\ D_4 &= \text{theory} \quad \text{write} \quad \boxed{\text{comment}} \end{aligned}$$

using words from corpus Vocabulary \Rightarrow $\text{People, watch, neuron, write, comment}$ as many AS = Dimension

$$D_1 = \left[\begin{bmatrix} 1, 0, 0, 0, 0 \\ 0, 1, 0, 0, 0 \\ 0, 0, 1, 0, 0 \end{bmatrix} \right]$$



$$\begin{aligned} D_2 &= \left[\begin{bmatrix} 0, 0, 1, 0, 0 \\ 0, 1, 0, 0, 0 \\ 0, 0, 1, 0, 0 \end{bmatrix} \right] \end{aligned}$$

$$D_3 = \left[\begin{matrix} [1, 0, 1, 0, 0, 0] \\ [0, 0, 0, 1, 0] \\ [0, 0, 0, 0, 1] \end{matrix} \right]$$

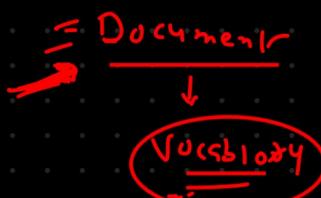
Word.

old \rightarrow 0HE

$$D_4 = \left[\begin{matrix} [0, 0, 1, 0, 0] \\ [0, 0, 0, 1, 0] \\ [0, 0, 0, 0, 1] \end{matrix} \right]$$

$[0, 0, 0, 1, 0] \in$ another word
 $[0, 0, 0, 0, 1] \in$ another

$$D_4 = \begin{bmatrix} 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$



Disadvantage $(\cancel{2010}, 2014)$

Sparse (Zero) $\underline{\underline{DL}}$

$\cancel{\text{It}}$
 $\cancel{\text{Sustain context}}$

② BOW (Bag of words)

- $D_1 \rightarrow$ People watch neuron
- $D_2 \rightarrow$ neuron watch neuron
- $D_3 \rightarrow$ people write comment
- $D_4 \rightarrow$ neuron write comment

BOW \Rightarrow ① Vocabulary
 ② based on the Document
 we calculate frequency
 $\text{No} \rightarrow (\cancel{\text{sum}}) \Rightarrow \cancel{\text{count}} \times \text{no. } X$

	People	watch	neuron	write	comment
D_1	1	1	1	0	0
D_2	0	1	2	0	0
D_3	1	0	0	1	1
D_4	0	0	1	1	1

Bow

I Simple

(2) Some sort of context

Few

(Item no of zero)

Dense

OHE \Rightarrow Sparse (more zero)

- Disadvantage \Rightarrow
- (1) if the size of vocab is too long Sparse (zero)
 - (2) OOV (it can't handle)

DOC \rightarrow Vocabulary \Rightarrow encode (frequency)



(3) N-GRAMS

$N=1$

$N=2$

$N=3$

$N=4$

$N=1 \Rightarrow$ Unigram

$N=2 \Rightarrow$ bigram

$N=3 \Rightarrow$ trigram

$N=4 \Rightarrow$ pair of two

" " three

$D_1 \Rightarrow$ People watch neuron

$D_2 \Rightarrow$ neuron water neuron

$D_3 \Rightarrow$ People write comment

$D_4 \Rightarrow$ neuron write comment

$\Rightarrow N=2 \Rightarrow$ bigram

(Pair of 2 word) unique

\Rightarrow customized.

$N=1 \Rightarrow \{ \text{People, watch, neuron, write, comment} \}$

$N=2 \Rightarrow \{ \text{People watch, watch neuron, neuron watch, People write, write comment, neuron write} \}$

$N=3$ \Rightarrow Pair of three words.

	People watch	watch iheuroh	inuron.watch	People write	write comment	i belmon white
D ₁	1	1	0	0	0	0
D ₂	0	1	1	0	0	0
D ₃	0	0	0	0	1	1
D ₄	0	0	0	0	0	1

- Advantage
- ① Capture more context
 - ② Reducing the ambiguity

this movie is very good

this movie is **hot** good

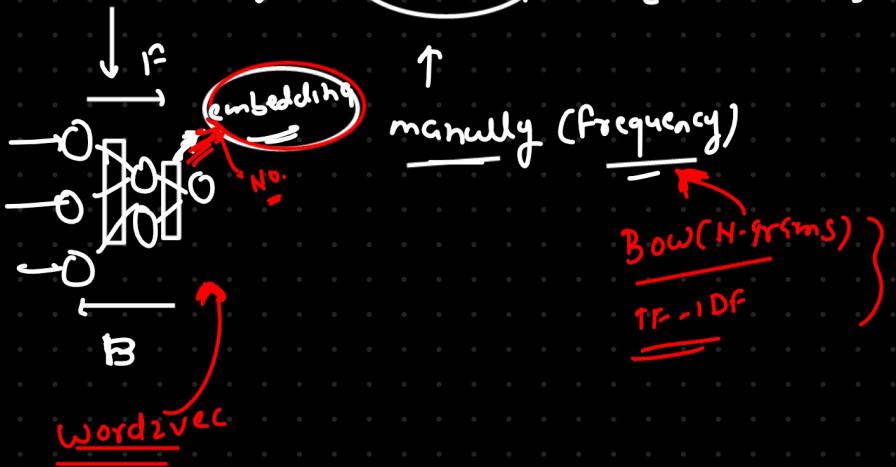
context

Unigram $N=1$ { This, movie, is, very, good, Not } \Rightarrow -

Bow $N=1$ { This movie, movie is, is very, very good, is hot, hot good }

$N=2$

Embedding & Encoding = []



4

TF-IDF

↳ ~~q u o g q e~~

TF \Rightarrow term frequency

IDF \Rightarrow Inverse Document frequency

$$\begin{aligned}
 &\text{No. of sentences} \\
 &D_1 \Rightarrow \boxed{\text{People}} \quad \boxed{\text{watch}} \quad \boxed{\text{Ineuron}} \quad . \quad \frac{1}{3} \cdot \frac{1}{3} \cdot \frac{1}{3} \log\left(\frac{4}{2}\right) \cdot \log\left(\frac{4}{3}\right) \\
 &D_2 \Rightarrow \text{Ineuron} \quad \text{watch} \quad \text{Ineuron} \quad . \quad \frac{2}{3} \cdot \frac{1}{3} \cdot \frac{2}{3} \quad " \quad " \\
 &D_3 \Rightarrow \text{People} \quad \text{watch} \quad \text{comment} \quad . \quad \frac{1}{3} \cdot \frac{1}{3} \cdot \frac{1}{3} \quad " \quad \log\left(\frac{4}{2}\right) \\
 &D_4 \Rightarrow \text{Ineuron} \quad \text{watch} \quad \text{comment} \quad . \quad \frac{1}{3} \cdot \frac{1}{3} \cdot \frac{1}{3} \quad " \quad "
 \end{aligned}$$

TF \Rightarrow No. of occurrence of the words in given D

(word, D)
 \downarrow
 Ineuron D_1

total word in the document

$$\text{IDF} = \log \left(\frac{\text{Total no. of Documents}}{\text{No. of Document which contains given term}} \right)$$

$$\text{Ineuron} = \log\left(\frac{4}{3}\right)$$

TF x IDF

TF x ID^F

D₁

D₂

D₃

D₄

People

watch

sheeon

wite

comment

$\frac{1}{3} \times \log(4/2)$