

NLP pipeline

① Data Acquisition → Pre-processing → Feature Engineering / Extraction

* Feature Extraction

- ① What is feature extraction & why do we need it?
- ② Why is it difficult & What is the core idea behind it?
- ③ What are the techniques required for this?

Text → Vector

↓
ML/DL - Algorithm → Matrix - numbers

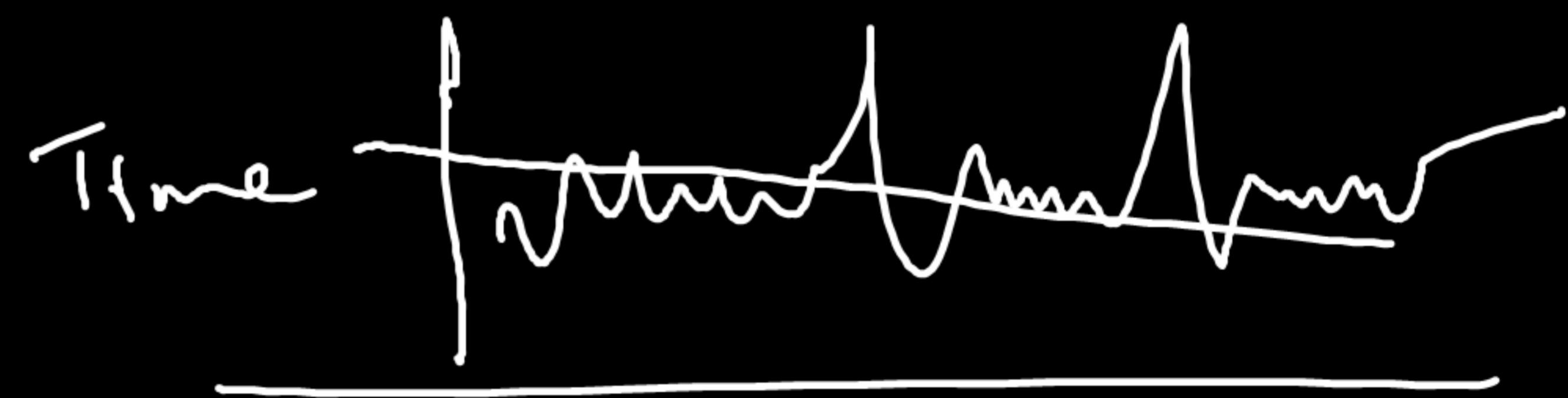
$y = f(x)$
↓
No

ML - ✓ Tabular data -

→ Time Series forecasting

Date & Time - Numerical Target ✓

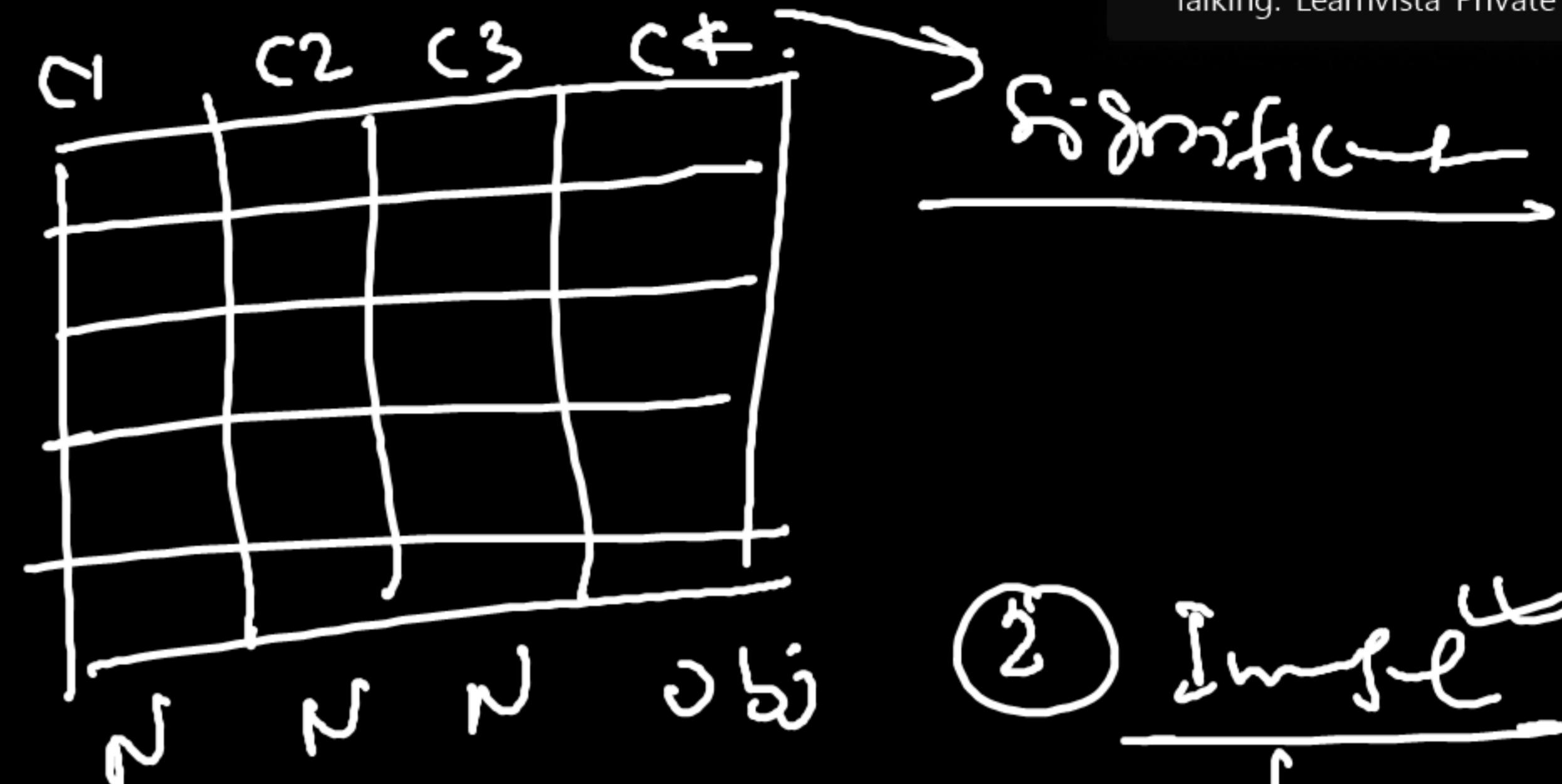
→ Heart rate



→ Text data

→ Hi, my name is Kumar Sandesh .

X my is Hi Sunday how kumar



② Image

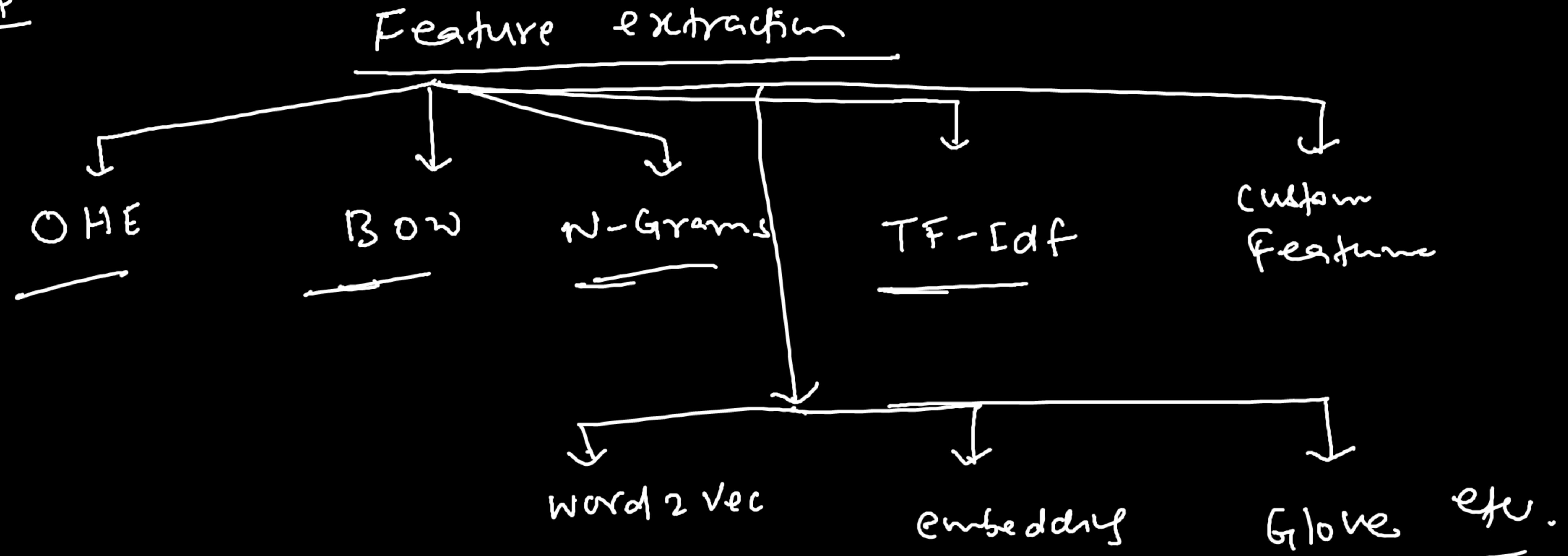
Pixel value

1
No

Frequency - Numer



NLP



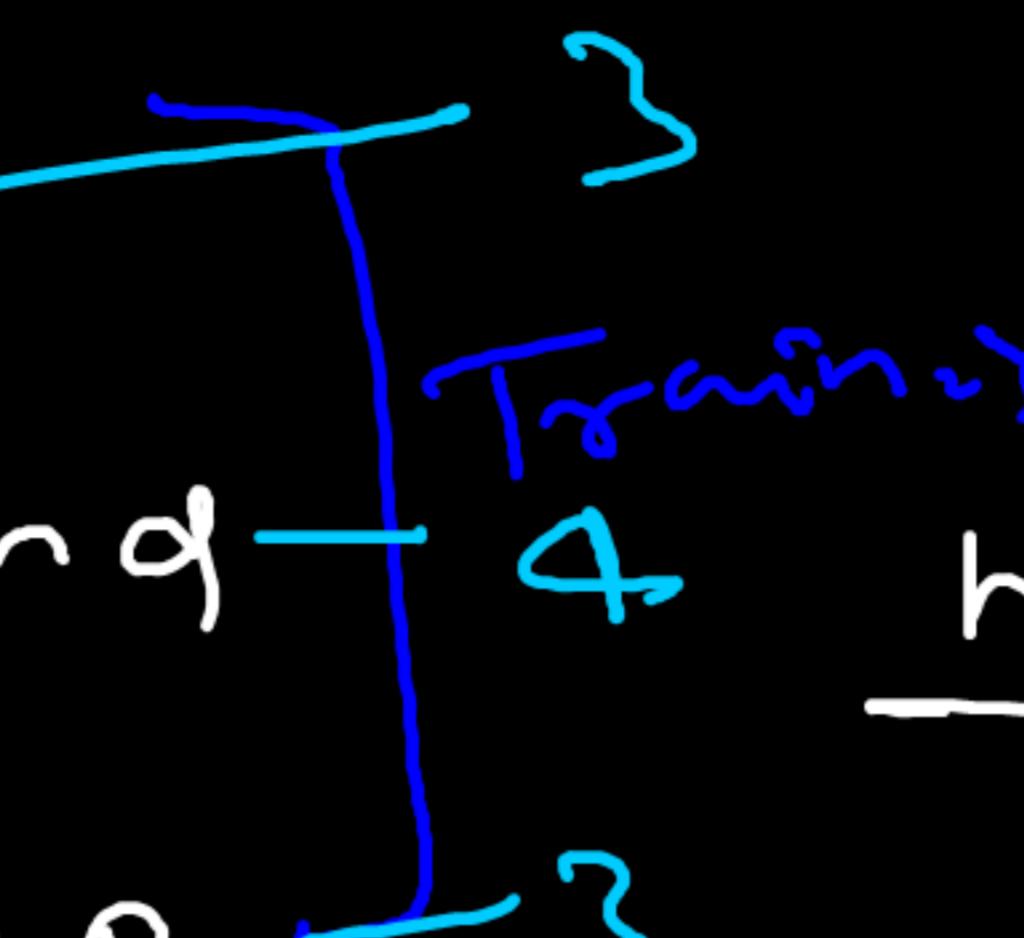
* Common Terms

- Corpus - consolidated ^{DOCs} words - each Sentence/DOC has many words from
- Vocabulary = unique words in all Documents
- Documents - 11 / Sentence / Review

S1

I like NLPS2 NLP is in Demand

S3 Learnby teaches NLP



$$\text{word} - S1 = 3$$

$$S2 = 4$$

$$S3 = 3$$

S4 → NLP ~~Learning - RNN concept~~

$$\underline{\text{Documents}} = 3$$

$$\underline{\text{vocab}} = \underline{\text{unique}} = 8$$

I like NLP is in Demand Learnby teaches

Corpus = 1 { I like NLP NLP is in
Demand Learnby teaches }
NLP

Objective : Text → Number / vector

How will we do it?



8 ↵
Vocab

I like NLP is in demand learn b.

HÉ: [1000000] [0100000] [0010000] [0001000] [0000100] [0000010]

Real time Problem - 50K vocab \rightarrow complex

Sparse matrix - more zeros

Approach 1 = OHE

OHE
↓

only works when
we have some
text

Locality	D	M	P	K	C	A	
De	1	0	0	0	0	0	(10000)
rus	0	1	0	0	0	0	(01000)
Per	0	0	1	0	0	0	(00100)
Kot	0	0	0	1	0	0	(000100)
che	0	0	0	0	1	0	(000010)
Amr	0	0	0	0	0	1	(000001)

zm Webinar Chat

ashish raiya to Hosts and panelists

AN Yes

Urmil Shah to Hosts and panelists

US less complex

Virat to Hosts and panelists

V yes

Rashmika Saravanan to Hosts and panelists

RS yes

Urmil Shah to Hosts and panelists

US sparsity and OOV

Kavitha to Hosts and panelists

K don't - will be removed

yes sir

Who can see your messages? Recording on

To: Hosts and panelists

Type message here...

D1 : I like NLP $\underline{[10000000] [01000000] [00100000]}$

Bag of words (BOW)

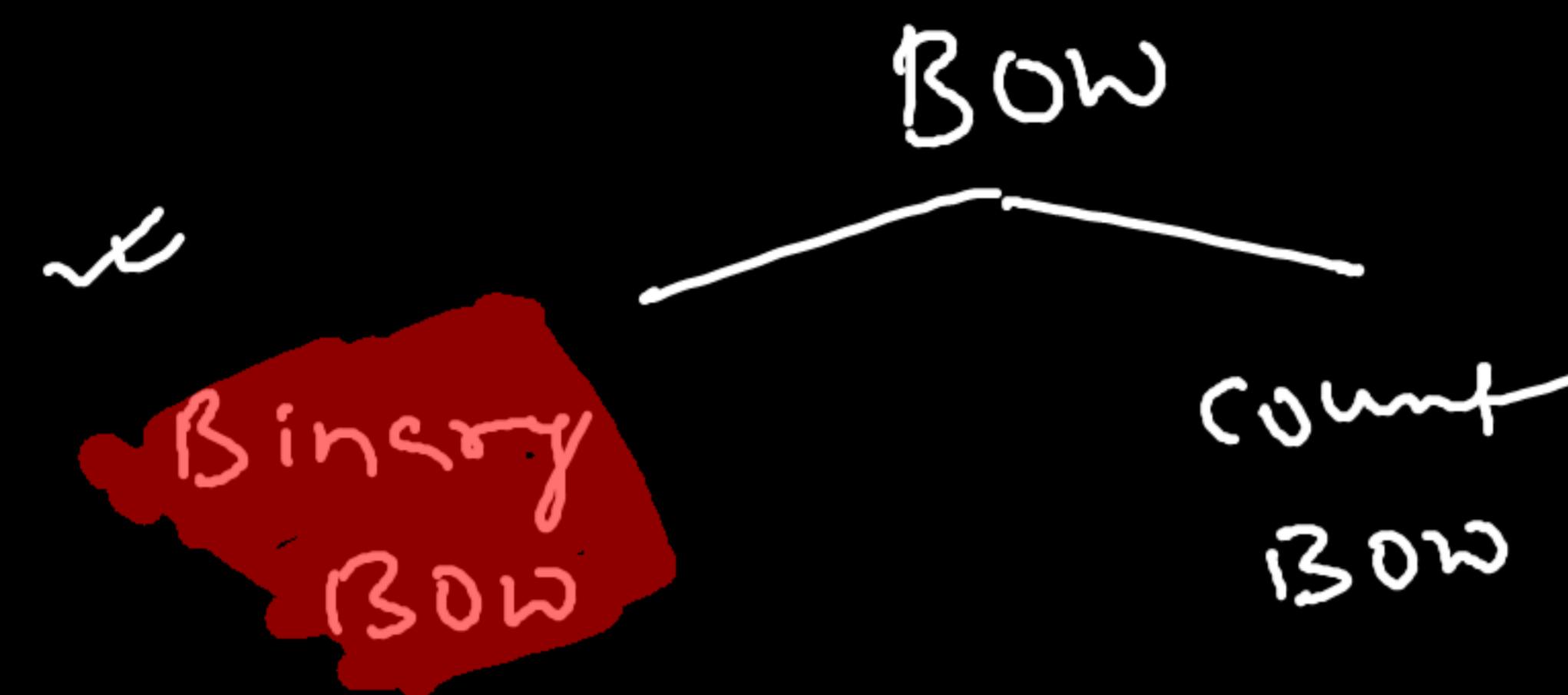
vocab = 8

	I	like	NLP	is	in	Demand	Learnvista	teaches
D1	1	1	1	0	0	0	0	0
D2	0	0	1	1	1	1	0	0
D3	0	0	1	0	0	0	1	1

D1 :- I like NLP : $\underline{[1 \ 1 \ 1 \ 0 \ 0 \ 0 \ 0]}$

→ This concept is called BOW idea ; Text → vector





100
NLP is really in Demand

NLP & NLP concept if Required

NLP is really in demand and concept require

D 1 -

| | } | } | } | } | } | } | }

Conf
Box

D2 =

2 1 0 0 0 1 1 1

12

1 1 0 0 0 1 1 1

J Binswanger
BOW

Pros

- ① Bow is extremely simple & intuitive

Cons :-

① sparsity

② Oov - Out of Vocab

I like pizza → ①
I ~~don't~~ like pizza ②

Remove stopw^{er}

↳ ignore the new vocab

