

NLP - by mentor Dr. Ayan Debnath, IIT Delhi + Harvard University Alumni

### Recap

- ① Stemming → sentiment analysis  
→ base word  
may/may not have meaning
- ② Lemmatization → root base  
→ always have meaning  
→ chatbot
- ③ Code

Topic → text → No / vectors (Neumencal representation)

i/p love NLP → o/p (1)

21 236 35

model

o/p (1 or 0)

### Bag of Words (BOWs)

binary BOW

sent1: He is intelligent boy  
sent2: she is intelligent girl  
sent3: Both of them are intelligent boy & girl

& good boy

lower  
→ stopwords

sent1: intelligent boy  
sent2: intelligent girl  
sent3: intelligent boy & girl

text Processing → ①

→ He is ! intelligent boy !#

step-1: ↓ data cleaning

He is intelligent boy

step-2: ↓ lowering

he is intelligent boy

step-3: ↓ remove stopwords

intelligent boy

step-4: ↓ stemming / Lemmatization

step-5: Tokenization

token ID

Vocabulary → 20K words

Dictionary  
He love NLP  
256 512 1000

He love NLP  
256 512 1000

Token ID

He love NLP & he also love ML  
256 512 1000 25 256 19 512

Ayan loves NLP

?? 512 1

### BOW

step-① Data cleaning

② lowering

③ Remove stopwords

④ Stemming

⑤ Calculate Histogram

Feature Set / Feature Engineering

intelligent	boy	girl	y
f <sub>1</sub>	f <sub>2</sub>	f <sub>3</sub>	y
1	1	0	1
1	0	1	0
1	1	1	1

Keyword  
Import info retrieve

Training Data

### Disadvantage

- ① semantic info X
- ② context X
- ③ weightage X
- ④ Bigger dataset  
BOW doesn't work well
- ⑤ Sparsity

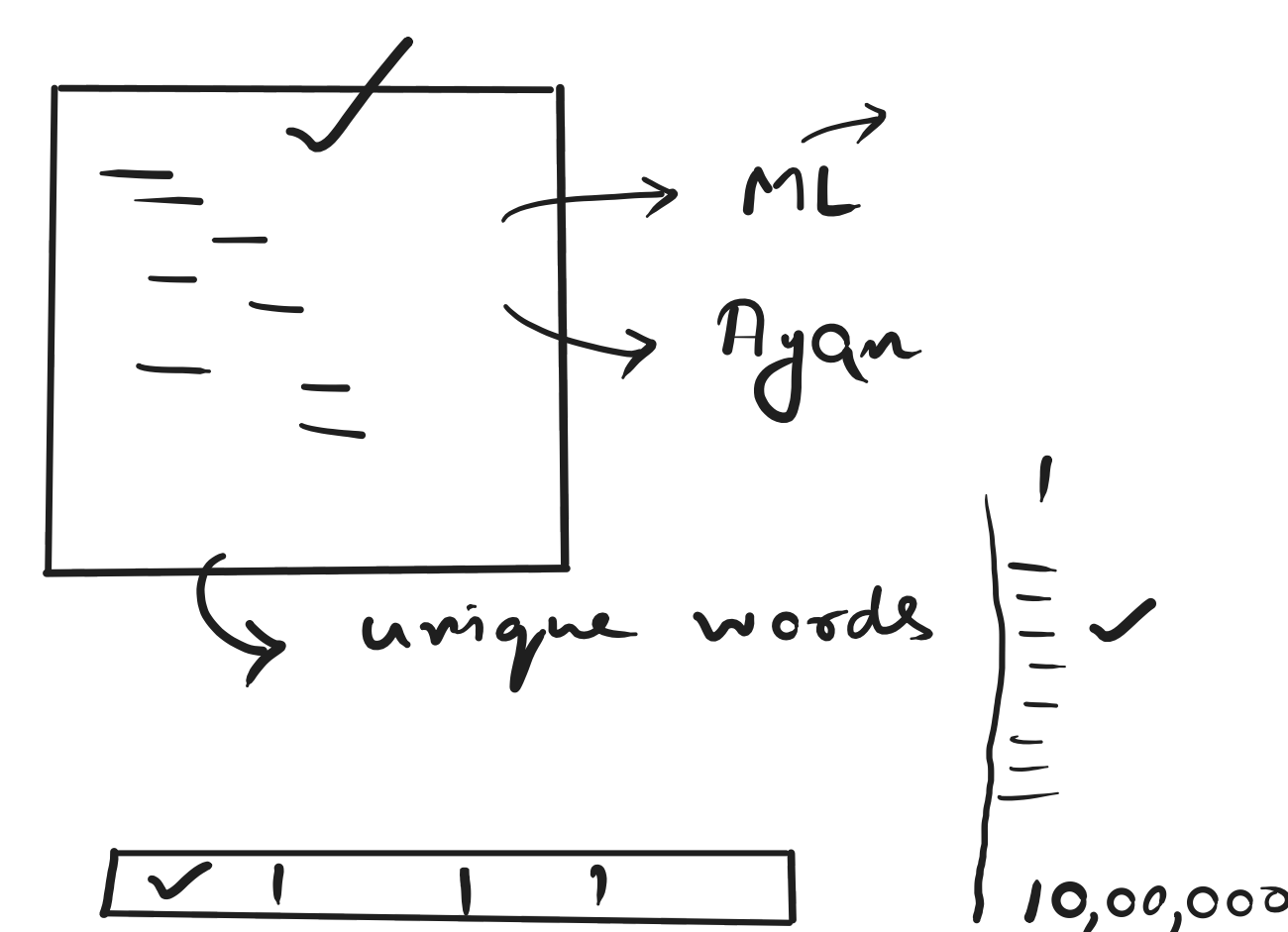
### advantage

- ① Easy to implement
- ② Processing time ↓
- ③ works well with small dataset

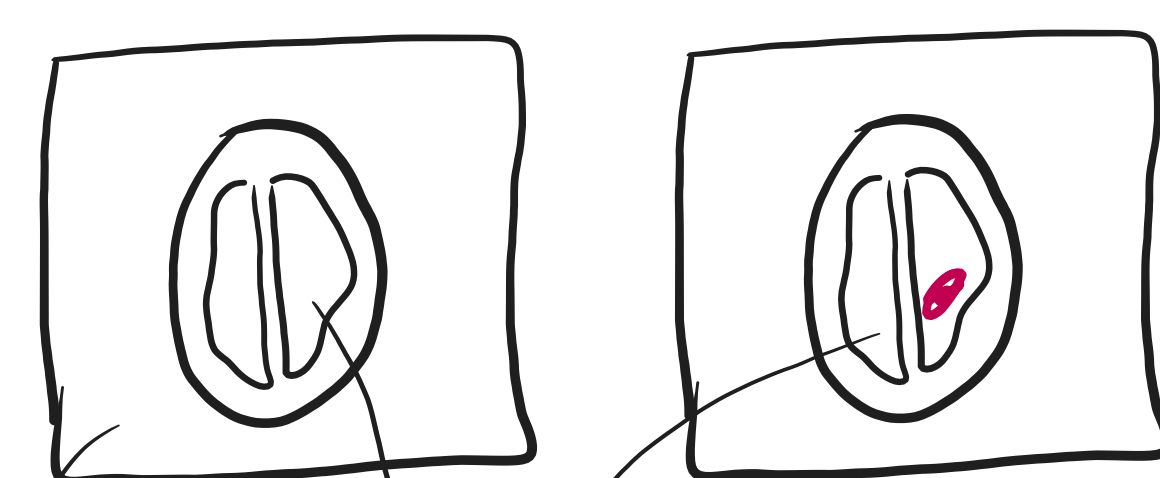
sent1  
sent2  
...  
sent 10,000

unique words

10,00,000



text → model →



Pt no.	BMR	Age	mean	median	f <sub>4</sub>	y/output
1						NT(0)
2						T(1)
...						
10,000						

Training Dataset