

## # Agenda:

### ① How to perform text preprocessing

- HTML Tags remove
- punctuation remove
- emoji handle
- lematization
- stemming
- lowe case

### ② Text Representation / word Embeddings / Text vectorization

- One Hot Encoding
- BOW ( Bag of word)
- TFIDF
- Word2Vec

### ③ project → sentiment Analysis using ML approach

- ML Algo
- 

## \* Lower case

① My name is bappy.

② bappy is a Data scientist

Bappy Barry  
bappy bappy →

You chat on WhatsApp, Instagram, Telegram

- \* I want the money as soon as possible
- + I want the money ASAP →

[Cloud bot] Amazon customer → chatbot  
→ training question

This pizza is Amzing → Amazing  
spell mistake

This pizza is very Tasty. I am very happy  
To be Hence in this store, Awesom experience

Increase in  $\downarrow$  time → computation same  
playing, played, plays  
strong → Play  
 $\downarrow$   
positive

## Text Representation

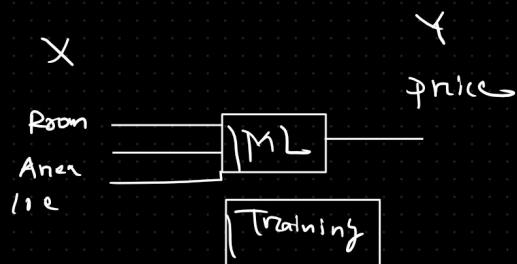
Corpus → R<sub>1</sub>  
Vocabulary → R<sub>2</sub>  
Document → R<sub>3</sub>  
Words → R<sub>4</sub>

Natural language processing is an interdisciplinary subfield of computer science and linguistics. It is primarily concerned with giving computers the ability to support and manipulate speech.

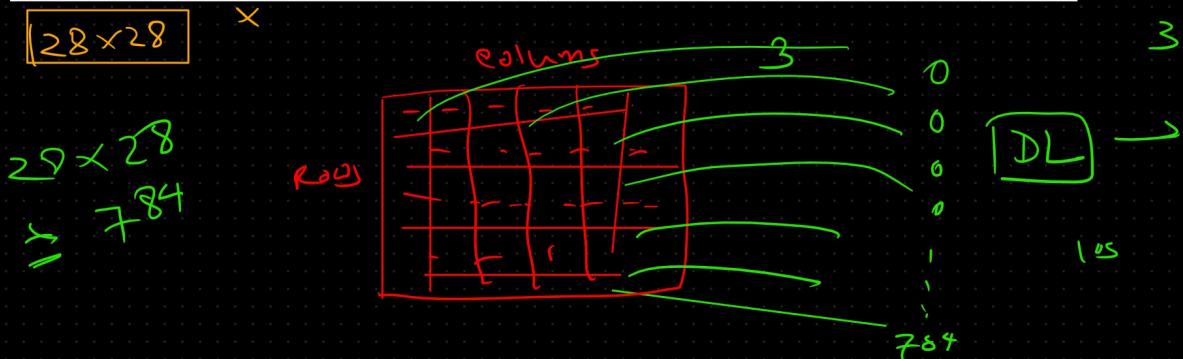
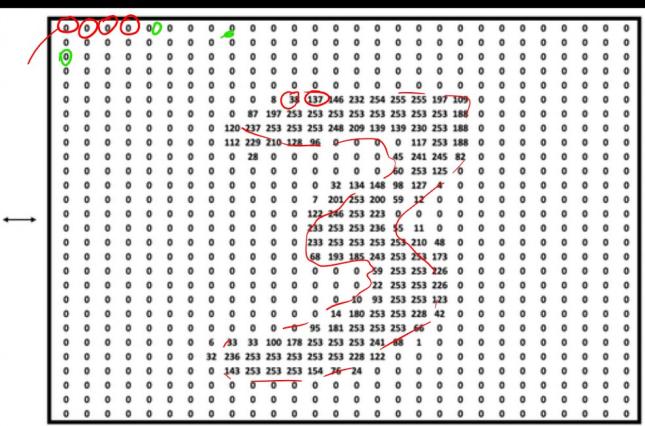
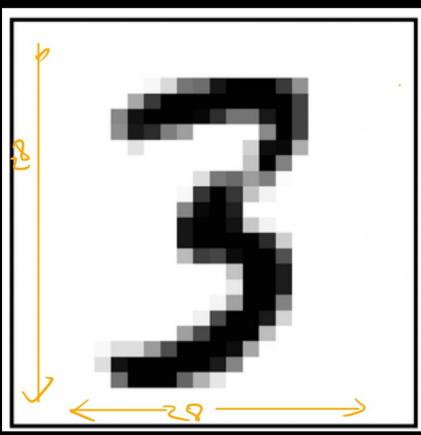
- Q1 What is Feature Extraction from text
- ② Why we need it?
  - ③ Why it is so difficult
  - ④ What is the core idea?
  - ⑤ What are the techniques (ML),

ML: tabular Data, csv, xsl

# Room	Area	Locality	price
5	20	Delhi	500
-	-	-	-
-	-	-	-
-	-	-	-



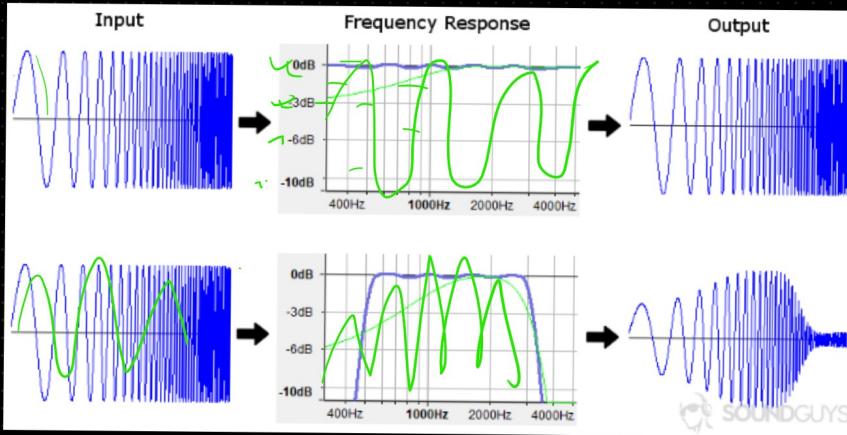
# CV: Images, Videos,



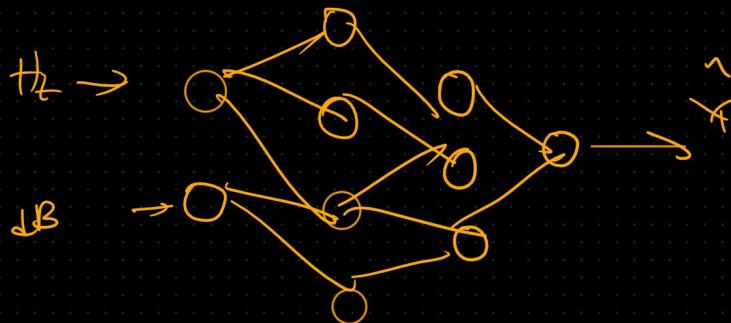
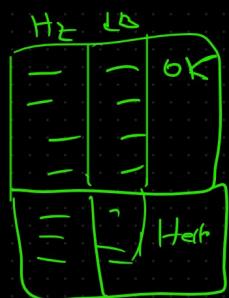
OK

speech

Hello

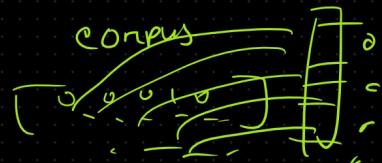


dB



Text data in NLL:

My name is Barry  
 $t_1 \quad t_2 \quad t_3 \quad t_4$



My name is Barry  
 0 1 2 3

# One Hot:

corpus:

D <sub>1</sub>	people watch dswithney
D <sub>2</sub>	dswithney watch dswithney
D <sub>3</sub>	people write comments
D <sub>4</sub>	dswithney write comment

people watch dswithney dswithney  
 watch dswithney people write  
 comments dswithney write comment

## # One Hot:

$R \rightarrow$	$D_1$	people watch dswithbary
-	$D_2$	dswithbary watch dswithbary
-	$D_3$	people write comments
-	$D_4$	dswithbary write comments

### corpus

people watch dswithbary dswithbary  
 watch dswithbary people write  
 comments dswithbary write comments

$$N = 5$$

people	watch	dswithbary	write	comment
--------	-------	------------	-------	---------

1	0	0	0	0
0	1	0	0	0
0	0	1	0	0
0	0	0	1	0
.	.	.	.	.

Sparse-matrix

$$D_1 = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \end{bmatrix}, \quad D_2 = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

$$\begin{matrix} 0 & 0 & : & : & : \\ \vdots & \vdots & & & \end{matrix}$$

$$\text{Sent} \rightarrow \begin{bmatrix} 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}$$

$$D_2 = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{bmatrix},$$

$$3 \times 5 \xrightarrow{\text{3}} \boxed{\text{model}} \rightarrow \text{shape} - \overline{(0 \ 0 \ 1)}$$

### \* Drawbacks:

- ① Sparsity
- ② No fixed size
- ③ OOV
- ④ No capturing semantic meaning

\* pos! easy to Implement

## # Bag of words

	X	Y	
1	D <sub>1</sub>	people watch ds with bry	1
2	D <sub>2</sub>	ds with bry water ds with bry	1
3	D <sub>3</sub>	people write comments	0
4	D <sub>4</sub>	ds with bry write comment	0

## Research project

during sentiment analysis

Bag of word is best technique

This movie is WOW → positive

This movie is Amazing, Amazing performance by

SRK. WOW!

Very bad movie. bad performed by X --  
Negative

This movie is Amazing, Amazing perform by SRK  
But very bad performance by Deepika →

↑  
n-grams → JK+

This movie is very good → Positive

This movie is not very good → Negative