

# Natural Language Processing (NLP) - Day 2

Recap

- 1) Application - Alexa, Voice assistant, Text Summarization
- 2) Roadmap -
- 3) Data Pre-processing - cleaning
- 4) Stemming, Lemmatization, Stop words
- 5) Hands on

BOW Bag of words

text  $x \rightarrow$  No.

boy girl

Sent1: He is intelligent boy  $\rightarrow$  intelligent boy

Sent2: she is intelligent girl  $\rightarrow$  intelligent girl

Sent3: Both of them are intelligent boy & girl

Food is good  $\rightarrow +ve / +1$   
" " bad  $\rightarrow -ve / 0$

word	Freq
boy	2
girl	2
intelligent	3

vector  
 $\Rightarrow$  No.s

	words	features	
	$f_1$ boy	$f_2$ girl	$f_3$ intelligent
sent1	2	0	1
2	0	1	1
3	1	1	1

=  $t/6$

model  $\rightarrow$  predict

## Disadvantage

- 1) Semantic information  $\times$
- 2) Context  $\times$
- 3) Huge dataset  $\times$

	features	
	HR	BMR
sl1		
sl2		
sl3		
sl4		

$$y = f(x) = f(x_1, x_2, x_3, \dots)$$

TF-IDF

independent variables

	intelligent	boy
words	words1	words2
sent1	0.9	0.01
2	1	
3	0	

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

TF (term Frequency)

IDF

$$TF = \frac{\text{No. of ref. words in sent}}{\text{Total no. of word.}}$$

$$IDF = \log\left(\frac{\text{No. of sentences}}{\text{No. of " containing that word}}\right)$$

words	IDF
boy	$\log(3/2)$
girl	$\log(3/2)$
intelligent	$\log(3/3)$

TF

	sent1	sent2	sent3
boy	1/2	0	1/2
girl	0	1/2	1/2
intelligent	1/3	1/3	1/3

words	Freq
boy	2
girl	2
intelligent	3

	word1	w2	w3
sent1	1 $\rightarrow$ 0.9	0	1
sent2	0	1	1 $\rightarrow$ 0.01
sent3	1	1	1

max. features = 3000

w1	w2	...	w2000
----	----	-----	-------

100  
1000  
50,000

✓  
✓  
✓