

improving improve improved

## Spam Mail Detection

The movie is good and scary. →

1.) Bag of Words

2.) TF-IDF

- 1.) Text processing is necessary.  $\rightarrow D_1$
- 2.) Text processing is necessary and <sup>is</sup> important.  $\rightarrow D_2$
- 3.) " " easy.  $\rightarrow D_3$

$\{D_1, D_2, D_3\} \rightarrow \text{corpus}$

Vocabulary  $\rightarrow$  Unique words in a corpus.

'Text'      'text'  
lower()

	Text	Processing	is	necessary	and	important	easy	-
D1	1	1	1	1	0	0	0	
D2	1	1	2	1	1	1	0	
D3	1	1	1	0	0	0	1	

TF - IDF

TF → Term Frequency

IDF → Inverse Document Frequency

$TF = \frac{\text{No. of times term 't' appears in a doc.}}{\text{No. of terms in a document}}$

$$TF(\text{Text}) = \frac{1}{4}$$

$$TF(\text{and}) = 0$$

$$TF(\text{pro}) = 1/4$$

$$TF(\text{imp}) = 0$$

$$TF(\text{is}) = 1/4$$

$$TF(\text{easy}) = 0$$

$$TF(\text{net}) = 1/4$$

Terms	TF(D1)	TF(D2)	TF(D3)
Text	1/4	1/6	1/4
pro	1/4	1/6	1/4
is	1/4	2/6	1/4
new	1/4	1/6	0
and	0	1/6	0
imp	0	1/6	0
easy	0	0	1/4

$$IDF = \log \frac{\text{no. of documents}}{\text{no. of document with term 't'}}$$

$$= \log \frac{3}{3} \Rightarrow \log(1) = 0$$

$$IDF(\text{necessary}) = \log \frac{3}{2}$$

Terms	$D_1$	$D_2$	$D_3$	IDF	$TF \times IDF(D_1)$	$TFI(D_2)$	$TF(D_3)$
Text	1	1	1	0	0	0	0
Proc.	1	1	1	0	0	0	0
is	1	2	1	0	0	0	0
nece.	1	1	0	$\log(3/2)$	$1/4 \times \log(3/2)$	$1/6 \times \log(3/2)$	0
and	0	1	0	$\log(3)$	0	$1/6 \times \log(3)$	0
imp.	0	1	0	$\log(3)$	0	$1/6 \times \log(3)$	0
easy	0	0	1	$\log(3)$	0	0	$1/4 \times \log(3)$

## Naive Bayes Classifier

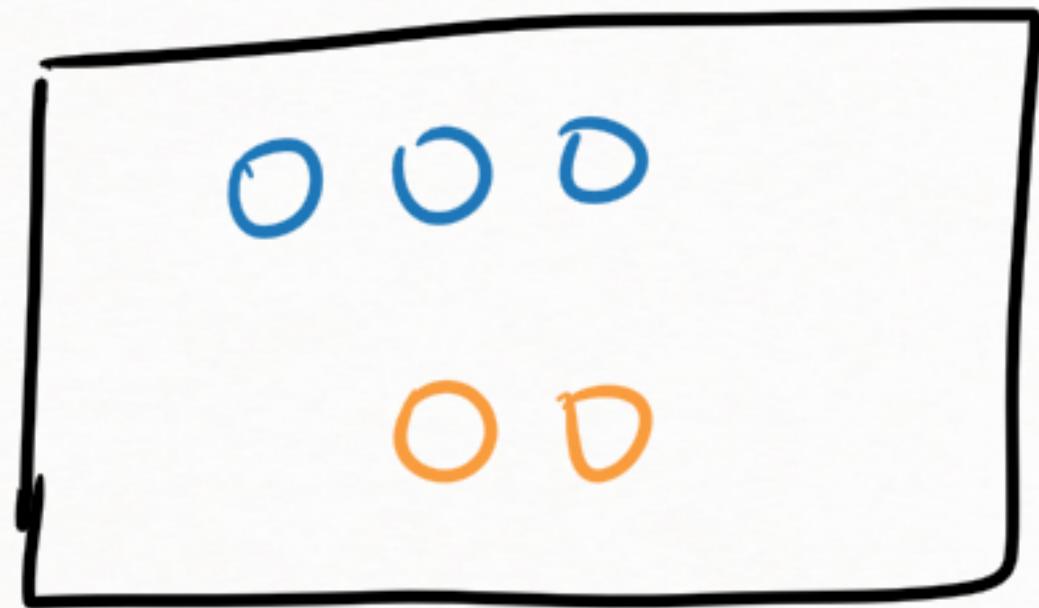
Independent Event → What is probability of getting  
head after tossing coin

Dependent "

$$P(H) = \frac{1}{2}, P(T) = \frac{1}{2}$$

$$P(1) = \frac{1}{6}$$

$$P(2) = \frac{1}{6}$$



$$P(B) = \frac{3}{5}$$

$$P(O) = \frac{2}{4} = \frac{1}{2}$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

Conditional  
-(i)  
Probability

$$A \cap B = B \cap A \quad - \text{ (ii)}$$

$$P(B/A) = \frac{P(B \cap A)}{P(A)} \quad - \text{ (iii)}$$

$$P(B \cap A) = P(B/A) \cdot P(A)$$

$$P(A|B) = \frac{P(A) \cdot P(B|A)}{P(B)}$$

↑ likelihood

↑ Prior

↓ Probability

Posterior

Probability

Outlook	Temperature	Humidity	Windy	PlayTennis
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

$$P(\text{Yes}) = \frac{9}{14}$$

$$P(\text{No}) = \frac{5}{14}$$

Outlook

	Yes	No	P(Yes)	P(No)
Sunny	2	3	2/5	3/5
Overcast	4	0	4/5	0/5
Rainy	3	2	3/5	2/5

## Temperature

	Yes	No	$P(Y)$	$P(N)$
Hot	2	2	2/9	2/5
Mild	4	2	4/9	2/5
Cool	3	1	3/9	1/5

$$P(\text{Yes} \mid \text{Rainy, Cool}) = ?$$

$$\begin{aligned}
 P(\text{Yes} / \overbrace{\text{Rainy}, \text{Cool}}^{A \quad B}) &= P(\text{Yes}) * P(\text{Rainy} / \text{Yes}) * P(\text{Cool} / \text{Yes}) \\
 &= \frac{5}{14} * \frac{2}{5} * \frac{1}{2} \\
 &= 0.047
 \end{aligned}$$

$$\begin{aligned}
 P(\text{No} / \text{Rainy}, \text{Cool}) &= P(\text{No}) * P(\text{Rainy} / \text{No}) * P(\text{Cool} / \text{No}) \\
 &= \frac{5}{14} * \frac{2}{5} * \frac{1}{5} = 0.020
 \end{aligned}$$

$$P(\text{Yes} | R, c) = \frac{0.047}{0.047 + 0.020} \\ = 0.62 \Rightarrow 62\%$$

$$P(\text{No} | R, c) = 1 - 0.62 \quad P(\text{sunny, not}) = ? \\ = 0.38 \approx 38\%.$$

