

NAIVE BAYES:-

→ Supervised ML Algorithm (classification)

→ Text classification.

ex. Mail $\begin{cases} \text{spam} \\ \text{not spam} \end{cases}$

Probability → chances of Outcome.

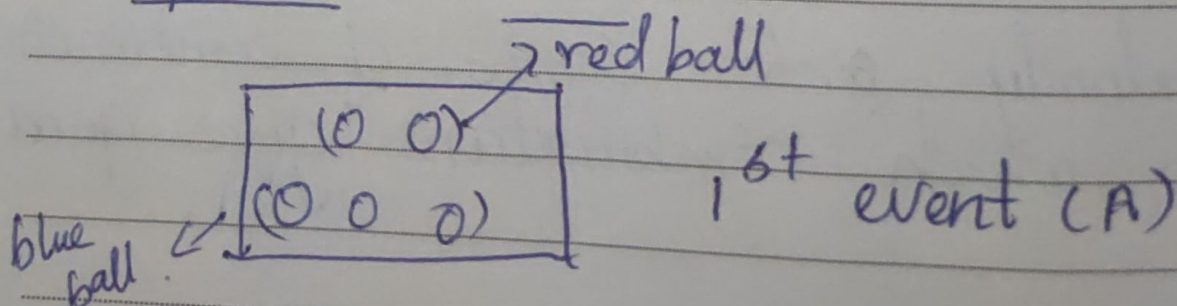
Independent Event:-

Toss a coin: H/T.

$$P(H) = 0.5$$

$$P(T) = 0.5$$

Dependent Event:-



Prob of getting a blue

$$P(A) = 3/5$$

Event (B):- Prob of getting red when event A is already taken place

$$P(B/A) = 1/4$$

NB → Conditional Probability.

↳ Bayes Theorem:-

$$P(A/B) = \frac{P(A) \times P(B/A)}{P(B)}$$

$$P(A) = 2/5, \quad P(B/A) = 1/4.$$

Prob. of both the events $[P(A) \& P(B/A)]$ to happen.

$$P(B \cap A) \rightarrow P(A) \times P(B/A)$$

$$= 2/5 \times 1/4 = 1/10.$$

$$P(B/A) = \frac{P(B \cap A)}{P(A)} = \frac{(1/10)}{(2/5)} = 1/4$$

Conditional probability.

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

↓
Prob. of event A when event B is always happened.

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

$$P(A) \times P(B/A) = P(B) \times P(A/B)$$

$$P(B/A) = \frac{P(B) \times P(A/B)}{P(A)}$$

$$P(A/B) = \frac{P(A) \times P(B/A)}{P(B)}$$

Bayes Theorem
↓
Used by NB classifier

why called Naive Bayes?

Bayes → Bayes Theorem.

Naive → It assume that each i/p

variables are independent.

let (x_1, x_2, x_3, x_4, y)

$$P(y | x_1, x_2, x_3, x_4) = P(x_1/y) P(x_2/y) P(x_3/y) P(x_4/y)$$

$$= \frac{P(x_1) P(x_2) P(x_3) P(x_4) \times P(y)}{P(x_1) P(x_2) P(x_3) P(x_4)}$$

↪ constant (denominator)

$$= \prod_{i=1}^n P(x_i/y) \times P(y)$$

probability of all the features

O/p will be $y = \arg \max \left(\prod_{i=1}^n P(x_i/y) \times P(y) \right)$

Ex:

Outlook	i/p		x_1	
	Yes	No	$P(Y)$	$P(NO)$
Sunny	2	3	2/9	3/5
Overcast	4	0	4/9	0/5
Rainy	3	2	3/9	2/5
	9	5		

Temperature x_2 i/p

	Yes	No	$P(y)$	$P(N)$
Hot	2	2	$2/9$	$2/5$
mild	4	2	$4/9$	$2/5$
Cold	3	1	$3/9$	$1/5$
	9	5		

Play y o/p

		$P(y)$
Yes	9	$9/14$
No	5	$5/14$

Prblm Stmt:-

Today is sunny & Hot, whether will play or not?

Today (sunny, hot)

$P(y) \rightarrow$

$P(N)$

$P(y)$ for sunny = $2/9$

$P(N)$ for ~~sunny~~ hot = $2/9$

$$P(y \text{ / today}) = P(\text{sunny / yes}) \times P(\text{hot / yes}) \times P(yes)$$

$P(\text{Today}) \rightarrow \text{Constant}$

$$= 2/9 \times 2/9 \times 9/14 = 0.031$$

P(N0):

$$P(N0) - \text{Sunny} = 3/5, \quad P(N0) - \text{Hot} = 2/5$$

$$P(N0/\text{today}) = 3/5 \times 2/5 \times 5/14$$

$$P(N0) = 5/14$$

$$= 0.08571$$

$P(y)$ when sunny & hot $\rightarrow 0.031$

$$P(N) \rightarrow 0.08571$$

Normalize it $P(y) = \frac{0.031}{0.031 + 0.08571} \approx 0.27$

$$P(N) = 1 - P(y) = 1 - 0.27 = 0.73$$

O/P \Rightarrow No. (max value is $P(N)$)

Player will not play when its hot

& Sunny.

Text Classification:-

Preprocessing for text data:-

- Remove punctuation
- Remove stop words (is, was, were, are...)
- Change everything into lower case
- Convert Text data to numerical vectors

Text to numerical vectors:-

i) Bag of words : (BOW)

→ Way of extracting features from text.

Ex: Sentence 1 : He is a good boy.
S2 : She is a good girl
S3 : Boy & girl are good

Step 1: Remove stop words:-

stop words are words which does not add much meaning to the sentence.

S₁: good boy

S₂: good girl

S₃: good boy girl

step 2: Convert text to numerical vectors
(Feature vector creation)

good boy girl.

S ₁	1	1	0
S ₂	1	0	1
S ₃	1	1	1

Disadvantage of Bow.

→ It cannot identify the semantics of a word in a sentence.

→ Feature Dimension.

→ does not capture the relationship between features.

(2) TF - IDF (Term Frequency \times Inverse Document frequency)

→ Quantifies the importance or relevance of string representation.

TF: measures how frequently a word appears in a document.

= $\frac{\text{no. of repetition of words in sentence}}{\text{Total: no of words in sentence}}$

Step 1: Remove stop words.

S_1 : good boy

S_2 : good girl.

S_3 : good boy girl.

Step 2: Calculate TF

	S_1	S_2	S_3
good	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{3}$
boy	$\frac{1}{2}$	0	$\frac{1}{3}$
girl	0	$\frac{1}{2}$	$\frac{1}{3}$

IDF : measure how important a word is.

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

word	IDF
good	$\log(3/3) = 0$
boy	$\log(3/2) =$
girl	$\log(3/2) =$

step 3: Final Vector Creation.

$\Rightarrow \text{TF} \times \text{IDF}$

	good	boy	girl
s_1	0	$\frac{1}{2} \times \log(3/2)$	0
s_2	0	0	$\frac{1}{2} \times \log(3/2)$
s_3	0	$\frac{1}{3} \times \log(3/2)$	$\frac{1}{3} \times \log(3/2)$

Types:

① Gaussian NB:-

- Used for Continuous features.
- It is used in classification and it assumes that ^{all} features follow a normal distribution.

② Multinomial NB:-

- It is used for discrete counts.
- Used when the features represents the frequency.
- Used for text classification.

③ Bernoulli NB:-

- the input features are binary (0's and 1's)
- Used for text classifications.