

## Mathematical formulation

data  $\rightarrow x_1, x_2, \dots, x_n$  Label ( $k$  no of classes)

Query pt  $\rightarrow \langle x_1, x_2, x_3, \dots, x_n \rangle$

• Since, there are  $k$  no of classes, Naive Bayes will find  $k$  no of probabilities, one for each class for new query point

$\rightarrow$  • Let's call this query pt  $x^T$  collectively

• Since there are  $k$  no of classes, let's say  $(y_1, y_2, y_3, \dots, y_k)$ , Naive Bayes will find probabilities for each class

• If there were 2 classes, ex Yes or No this will be called Binary

• Since there are  $k$  no of classes, the probab. will be

$$\rightarrow P(y_1 | x^T) = P(x^T | y_1) P(y_1) / P(x^T)$$

$$\rightarrow P(y_2 | x^T) = P(x^T | y_2) P(y_2) / P(x^T)$$

$$\rightarrow P(y_3 | x^T) = P(x^T | y_3) P(y_3) / P(x^T)$$

$$\rightarrow P(y_k | x^T) = P(x^T | y_k) P(y_k) / P(x^T)$$

- We can remove denominator, for simplification
- In among <sup>class</sup> whose probabilities is most, that class will be assigned to our query  $x^T$

• Let's focus on general case

$$P(Y_K | X_T) = P(X_T | Y_K) P(Y_K)$$

• putting back  $X_T$ , and our equation will be

$$P(Y_K | X_T) = P(\underline{X_1 \cap X_2 \cap X_3 \cap \dots \cap X_n} | Y_K) P(Y_K)$$

• Intersection is there, as in query it represents events are happening at same time that's how query was made

•  $\cap$  can be represented as, as both are same

$\Rightarrow$  can be rewrite as

$$P(Y_K | X_T) = P(\underline{X_1, X_2, X_3, \dots, X_n} | Y_K) P(Y_K)$$

We know 
$$P(A|B) = \frac{P(A, B)}{P(B)}$$

Therefore  $\Rightarrow P(Y_K | X_T) = \frac{P(X_1, X_2, X_3, \dots, X_n, Y_K)}{P(Y_K)} \cdot P(Y_K)$

$$P(Y_K | X_T) = P(X_1, X_2, X_3, \dots, X_n, Y_K)$$

We know 
$$P(A|B) = \frac{P(A, B)}{P(B)}$$

$$\Rightarrow \underline{P(A, B) = P(A|B) \times P(B)}$$

We have 
$$P(Y_K | X_T) = P(\underbrace{X_1}_A, \underbrace{X_2, X_3, \dots, X_n, Y_K}_B)$$

Considering  $A \rightarrow X_1$  &  $B \rightarrow X_2, X_3, \dots, X_n, Y_K$  apply



underline rule  $\Rightarrow P(A, B) = P(A|B) \times P(B)$

here,

$$\Rightarrow P(Y_k, X_1) = P(X_1 | X_2, X_3, \dots, X_n, Y_k) \times \underbrace{P(X_2, X_3, \dots, X_n, Y_k)}_{2^{nd}}$$

Applying same rule again on 2<sup>nd</sup> term

$$= P(X_1 | X_2, X_3, \dots, X_n, Y_k) \times P(X_2 | X_3, X_4, \dots, X_n, Y_k) \times \underbrace{P(X_3, X_4, \dots, X_n, Y_k)}_3$$

$\Rightarrow$  Applying same again on 3<sup>rd</sup> term

$$= P(X_1 | X_2, X_3, \dots, X_n, Y_k) \times P(X_2 | X_3, X_4, \dots, X_n, Y_k) \times P(X_3 | X_4, X_5, \dots, X_n, Y_k)$$

$\Rightarrow$  Applying till last

$$\Rightarrow \underbrace{P(X_1 | X_2, X_3, \dots, X_n, Y_k) \times P(X_2 | X_3, X_4, \dots, X_n, Y_k) \times \dots \times P(X_{n-1} | X_n, Y_k)}_{P(X_n | Y_k) \times P(Y_k)} \rightarrow \text{main eqn}$$

Main cheez: Naive Bayes takes Naive assumption, that in your data there are  $X_1, X_2, \dots, X_n$  columns, and you take assumption that all the features/columns are independent of each other

- That means  $X_1$  doesn't depend on  $X_2, X_3, X_4, \dots, X_n$
- same goes for  $X_2$  " " on  $X_3, X_4, \dots, X_n$

- $\rightarrow$  • we have this equation
- let's take one term out of it just for explanation

- $P(X_1 | X_2, X_3, \dots, X_n, Y_k)$  (remember, represents  $n$ )
- We know that  $x_1$  doesn't depend on  $x_2, x_3, \dots, x_n$ , but depend on  $Y_k$
- We know that when events are independent :-  $P(A|B) = P(A)$
- But in our case,  $x_1$  ~~depend~~ independent of  $x_2, x_3, \dots, x_n$  but depend on  $Y_k$
- Something like this  $\rightarrow P(A|B \cap C) = P(A|B, C)$
- where  $A$  is independent of  $B$ , but depends on  $C$

In this situation, you can write

- $P(A|B, C)$  as  $\Rightarrow P(A|B \cap C)$
- because of  $P(A|B) = P(A)$  when  $A$  &  $B$  are independent

- Apply same logic to main eq<sup>n</sup>
- Our eq<sup>n</sup> will be modified to :-

$$P(Y_k | X_1, X_2, \dots, X_n) = P(X_1 | Y_k) \cdot P(X_2 | Y_k) \cdot P(X_3 | Y_k) \cdot \dots \cdot P(X_{n-1} | Y_k) \cdot P(Y_k)$$

Ex for class  $y=1$ , the prob will be :-

$$P(Y_1 | X_1, X_2, \dots, X_n) = P(X_1 | Y_1) \cdot P(X_2 | Y_1) \cdot P(X_3 | Y_1) \cdot \dots \cdot P(X_n | Y_1) \cdot P(Y_1)$$

- Similarly, you will find all class probabilities
- Compare each class prob, whichever is highest that will be assigned that label



- That's why the eq<sup>n</sup> was broken down into simpler terms (in intuition part) (This was the proof)

## \* Overview of training and testing phase of Naive Bayes

### Training Phase:

- In training phase you calculate all possible probabilities

Ex. In a dataset we used, we have two classes  $\rightarrow$  'Yes' and 'No'

- for every input column, you check how many categories you have

Ex In ~~output~~ <sup>outlook</sup> col we had categories: - Sunny, Overcast, rainy

- we have two classes Yes and No, and in outlook  $\rightarrow$  3 categories
- We will calculate  $3 \times 2 = 6$  probabilities
- i.e.
  - $P(\text{Sunny} | \text{Yes})$
  - $P(\text{Sunny} | \text{No})$
  - $P(\text{Overcast} | \text{Yes})$
  - $P(\text{Overcast} | \text{No})$
  - $P(\text{Rainy} | \text{Yes})$
  - $P(\text{Rainy} | \text{No})$

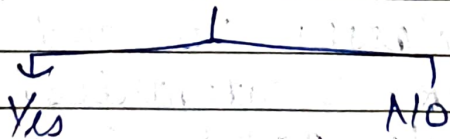
- Same goes for other columns.
- Temperature contains Hot, Mild, Cool.
- 6 conditional prob will be calculated
  - ie  $\rightarrow P(\text{Hot} | \text{Yes})$
  - $\rightarrow P(\text{Hot} | \text{No})$
  - $\rightarrow P(\text{Mild} | \text{Yes})$
  - $\rightarrow P(\text{Cool} | \text{Yes})$
  - $\rightarrow P(\text{Mild} | \text{No})$
  - $\rightarrow P(\text{Cool} | \text{No})$

- During training, it will calculate all these possible probabilities and store it in dictionary.

EX

→ If new query pt occurs, example:

- { Sunny, Hot, mist, False }  $\Rightarrow x^T$
- We have to predict whether it will be Yes or No.
- 2 probabilities will be calculated as there is 2 classes



- $P(\text{Yes} | x^T) \text{ \& } P(\text{No} | x^T)$

- $$P(\text{Yes} | x^T) = P(x_1 | \text{Yes}) \times P(x_2 | \text{Yes}) \times P(x_3 | \text{Yes}) \times P(x_4 | \text{Yes}) \times P(\text{Yes})$$

- $$P(\text{Yes} | x^T) = P(\text{Sunny} | \text{Yes}) \times P(\text{Hot} | \text{Yes}) \times P(\text{mist} | \text{Yes}) \times P(\text{False} | \text{Yes}) \times P(\text{Yes})$$

- And these all probabilities are

already stored in dictionary

- Similarly  $P(N_0 | X_T)$  will be calculated
- Whosever probab is highest, will be assigned that label
- refer to code example