

# COMP 3605

## Naïve Bayes Classifier

- The Naive Bayes classifier is a family of probabilistic classifiers based on applying Bayes' theorem with the assumption of independence between every pair of features.
- Despite its simplicity, Naive Bayes can be surprisingly accurate and is particularly useful for very large datasets.
- It's often used for text classification, spam filtering, sentiment analysis, and recommendation systems.

### Bayes' Theorem:

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

Where:

- $P(A|B)$  is the posterior probability of class  $A$  given predictor  $B$ .
- $P(A)$  is the prior probability of class  $A$ .
- $P(B|A)$  is the likelihood which is the probability of predictor  $B$  given class  $A$ .
- $P(B)$  is the prior probability of predictor  $B$ .

## EXAMPLE 1

<i>TID</i>	<i>age</i>	<i>income</i>	<i>student</i>	<i>credit_rating</i>	<i>buys_computer?</i>
1	youth	high	no	fair	no
2	youth	high	no	excellent	no
3	middle_aged	high	no	fair	yes
4	senior	medium	no	fair	yes
5	senior	low	yes	fair	yes
6	senior	low	yes	excellent	no
7	middle_aged	low	yes	excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	senior	medium	yes	fair	yes
11	youth	medium	yes	excellent	yes
12	middle_aged	medium	no	excellent	yes
13	middle_aged	high	yes	fair	yes
14	senior	medium	no	excellent	no

Predict the class label for this new tuple:

$X = (\text{age} = \text{youth}, \text{income} = \text{medium}, \text{student} = \text{yes}, \text{credit\_rating} = \text{fair})$

i.e.,  $X = (x_1 = \text{youth}, x_2 = \text{medium}, x_3 = \text{yes}, x_4 = \text{fair})$

<i>age</i>	<i>buys_computer</i>	
	yes	no
youth	2	3
middle_aged	4	0
senior	3	2

<i>income</i>	<i>buys_computer</i>	
	yes	no
low	3	1
medium	4	2
high	2	2

<i>student</i>	<i>buys_computer</i>	
	yes	no
yes	6	1
no	3	4

<i>credit_rating</i>	<i>buys_computer</i>	
	yes	no
fair	6	2
excellent	3	3

## 1. Calculate Prior Probabilities:

Calculate the probability of each class in the class\_buy\_computer column.

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

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

## 2. Calculate Likelihoods:

For each feature (age, income, student, credit\_rating), calculate the likelihood of that feature given each class.

For Buy\_Computer = Yes:

$$P(\text{Age} = \text{Youth} \mid \text{Yes}) = 2/9$$

$$P(\text{Income} = \text{Medium} \mid \text{Yes}) = 4/9$$

$$P(\text{Student} = \text{Yes} \mid \text{Yes}) = 6/9$$

$$P(\text{Credit\_Rating} = \text{Fair} \mid \text{Yes}) = 6/9$$

For Buy\_Computer = No:

$$P(\text{Age} = \text{Youth} \mid \text{No}) = 3/5$$

$$P(\text{Income} = \text{Medium} \mid \text{No}) = 2/5$$

$$P(\text{Student} = \text{Yes} \mid \text{No}) = 1/5$$

$$P(\text{Credit\_Rating} = \text{Fair} \mid \text{No}) = 2/5$$

### 3. Calculate Posterior Probabilities:

For the given example X, calculate the probability of each class given X by multiplying the prior probabilities by the likelihoods.

$$P(X | C_i) = P(x_1 | C_i) \times P(x_2 | C_i) \times \dots \times P(x_n | C_i)$$

Compute  $P(X | \text{yes}) = P(\text{age} = \text{youth} | \text{yes}) * P(\text{income} = \text{medium} | \text{yes}) * P(\text{student} = \text{yes} | \text{yes}) * P(\text{credit\_rating} = \text{fair} | \text{yes})$

$$= 2/9 \times 4/9 \times 6/9 \times 6/9$$

Compute  $P(X | \text{no}) = P(\text{age} = \text{youth} | \text{no}) * P(\text{income} = \text{medium} | \text{no}) * P(\text{student} = \text{yes} | \text{no}) * P(\text{credit\_rating} = \text{fair} | \text{no})$

$$= (3/5) * (2/5) * (1/5) * (2/5)$$

To find the class  $C_i$  that maximizes  $P(X | C_i)P(C_i)$ , we compute:

$P(\text{Yes} | X) = P(\text{Yes}) * P(\text{Age} = \text{Youth} | \text{Yes}) * P(\text{Income} = \text{Medium} | \text{Yes}) * P(\text{Student} = \text{Yes} | \text{Yes}) * P(\text{Credit\_Rating} = \text{Fair} | \text{Yes})$

$P(\text{No} | X) = P(\text{No}) * P(\text{Age} = \text{Youth} | \text{No}) * P(\text{Income} = \text{Medium} | \text{No}) * P(\text{Student} = \text{Yes} | \text{No}) * P(\text{Credit\_Rating} = \text{Fair} | \text{No})$

Substitute the probabilities calculated above:

$$P(\text{Yes} | X) = (9/14) * (2/9) * (4/9) * (6/9) * (6/9) = 0.0282$$

$$P(\text{No} | X) = (5/14) * (3/5) * (2/5) * (1/5) * (2/5) = 0.0069$$

**4.Prediction:** The class with the highest posterior probability will be the prediction.

Since  $P(\text{Yes} | X) > P(\text{No} | X)$ , the Naive Bayes Classifier predicts that the computer will be bought.(yes)

## EXAMPLE 2:

The "Play Tennis" dataset is a well-known toy dataset often used to illustrate the concepts of classification algorithms. The dataset consists of 14 instances, each representing a day with corresponding weather conditions and whether or not tennis was played on that day.

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Let's say we have a new instance with the following features:

Outlook: Sunny

Temperature: Mild

Humidity: High

Wind: Strong

We want to predict whether tennis will be played under these conditions (Play Tennis: Yes/No).

Temperature	Yes	No
Hot	2	2
Mild	4	2
Cool	3	1

Humidity	Yes	No
High	3	4
Normal	6	1

Outlook	Yes	No
Sunny	2	3
Overcast	4	0
Rain	3	2

Wind	Yes	No
Strong	3	3
Weak	6	2

Predict a class label using the Naive Bayes Classifier for a new instance with the following features:

- Outlook: Sunny
- Temperature: Mild
- Humidity: High
- Wind: Strong

We want to predict whether tennis will be played under these conditions (Play Tennis: Yes/No).

## 1. Calculate Prior Probabilities

$P(\text{Yes}) = 9/14$  (9 days of playing tennis out of 14)

$P(\text{No}) = 5/14$  (5 days of not playing tennis out of 14)

## 2. Calculate Likelihoods

For the class label Yes:

$P(\text{Outlook}=\text{Sunny}|\text{Yes}) = 2/9$  (2 sunny days out of 9 days of playing tennis)

$P(\text{Temperature}=\text{Mild}|\text{Yes}) = 4/9$

$P(\text{Humidity}=\text{High}|\text{Yes}) = 3/9$

$P(\text{Wind}=\text{Strong}|\text{Yes}) = 3/9$

For the class label No:

$P(\text{Outlook}=\text{Sunny}|\text{No}) = 3/5$  (3 sunny days out of 5 days of not playing tennis)

$P(\text{Temperature}=\text{Mild}|\text{No}) = 2/5$

$P(\text{Humidity}=\text{High}|\text{No}) = 4/5$

$P(\text{Wind}=\text{Strong}|\text{No}) = 3/5$

## 2. Calculate Posterior Probabilities

$$P(X | C_i) = P(x_1 | C_i) \times P(x_2 | C_i) \times \dots \times P(x_n | C_i)$$

$P(X|\text{Yes}) = P(\text{Outlook}=\text{Sunny}|\text{Yes}) * P(\text{Temperature}=\text{Mild}|\text{Yes}) *$

$P(\text{Humidity}=\text{High}|\text{Yes}) * P(\text{Wind}=\text{Strong}|\text{Yes})$

$P(X|\text{Yes}) = P(2/9) * (4/9) * (3/9) * (3/9)$

$P(X|\text{No}) = P(\text{Outlook}=\text{Sunny}|\text{No}) * P(\text{Temperature}=\text{Mild}|\text{No}) *$

$P(\text{Humidity}=\text{High}|\text{No}) * P(\text{Wind}=\text{Strong}|\text{No})$

$P(X|\text{No}) = (3/5) * (2/5) * (4/5) * (3/5)$

For the given example  $X = (\text{Outlook}=\text{Sunny}, \text{Temperature}=\text{Mild}, \text{Humidity}=\text{High}, \text{Wind}=\text{Strong})$ , calculate the posterior probabilities:

$$P(\text{Yes} | X) = P(\text{Yes}) * P(\text{Outlook}=\text{Sunny} | \text{Yes}) * P(\text{Temperature}=\text{Mild} | \text{Yes}) * P(\text{Humidity}=\text{High} | \text{Yes}) * P(\text{Wind}=\text{Strong} | \text{Yes})$$

$$P(\text{No} | X) = P(\text{No}) * P(\text{Outlook}=\text{Sunny} | \text{No}) * P(\text{Temperature}=\text{Mild} | \text{No}) * P(\text{Humidity}=\text{High} | \text{No}) * P(\text{Wind}=\text{Strong} | \text{No})$$

Substitute the calculated probabilities:

$$P(\text{Yes} | X) = (9/14) * (2/9) * (4/9) * (3/9) * (3/9) = 0.0053$$

$$P(\text{No} | X) = (5/14) * (3/5) * (2/5) * (4/5) * (3/5) = 0.0206$$

**3. Prediction:** The class with the highest posterior probability will be the prediction.

Since  $P(\text{No} | X) > P(\text{Yes} | X)$ , the Naive Bayes Classifier predicts that tennis will not be played under the given conditions (Play Tennis: No).



## EXERCISE 1:

### Dataset:

Age	Income	Gender	Owns House	Buy Product
Young	High	M	N	No
Young	High	M	N	No
Middle	High	M	Y	Yes
Old	Medium	F	Y	Yes
Old	Low	F	Y	No
Old	Low	F	N	Yes
Middle	Low	F	Y	Yes
Young	Medium	M	N	No
Young	Low	F	Y	Yes
Old	Medium	F	Y	Yes
Young	Medium	F	Y	Yes
Middle	Medium	M	Y	Yes
Middle	High	F	Y	Yes
Old	Medium	M	N	No
Middle	Low	F	N	Yes

Predict if an individual with the following attributes will buy the product:

Age: Middle

Income: Low

Gender: F

Owns House: N