

Naive Bayes classification

Ques

Consider the following dataset & predict the class of a new instance X using Naive Bayes classification algorithm

ID	Refund	Marital status	Amount	Evide
1	Yes	Singe	125K	No
2	No	Married	100K	No
3	No	Singe	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Singe	85K	Yes
9	No	Married	75K	No
10	No	Singe	90K	Yes.

'Where

$X = (\text{Refund} = \text{No} ; \text{Marital status} = \text{Married} ; \text{Income} = 120\text{K})$.

PREDICT ???

Understanding

Evide is class.

Class is either Yes or No.

ID	Refund	Marital Status	Amount	Evasion
1	Yes	Single	12K	No
2	NO	Married	100K	No
3	NO	Single	70K	No
4	Yes	Married	120K	No
5	NO	Divorced	95K	Yes
6	NO	Married	60K	No
7	Yes	Divorced	220K	No
8	NO	Single	85K	Yes
9	NO	Married	75K	No
10	NO	Single	90K	Yes
→ 11	NO	MARRIED	120K	??
PREDICT				

* Bayesian classification

Problem statement

- Given features x_1, x_2, \dots, x_n .
- Predict a label y .

* Bayes' classifier

- A probabilistic framework for solving classification problems.

- Conditional probability

$$P(C|A) = \frac{P(A, C)}{P(A)} - 1$$

$$P(A|C) = \frac{P(A, C)}{P(C)} \quad - 2$$

Read above as
Probability of C given A. $\rightarrow 1$

Probability of A given C $\rightarrow 2$

* Bayes theorem:

* Remember Bayes classifier uses conditional Probability

Bayes Theorem:-

$$P(C|A) = \frac{P(A|C) P(C)}{P(A)}$$

Probability of C given A is calculated as $\frac{(Probability\ of\ A\ given\ C) * (Probability\ of\ C)}{Probability\ of\ A}$.

In general, compute the posterior probability $P(C | A_1, A_2, \dots, A_n)$ for all values of C using the Bayes Theorem.

Choose the value of C that maximizes $P(C | A_1, A_2, \dots, A_n)$

* class

$$P(C) = \frac{N_C}{N}$$

Eg.

Count the no. of 'No's & 'Yes' from the data set

$$n(\text{No}) = 7, \quad n(\text{Yes}) = 3.$$

No. of 'No's No. of 'Yes'.

$$\text{Probability of 'No'} = \frac{7}{10} = P(\text{No})$$

$$\text{Probability of 'Yes'} = \frac{3}{10} = P(\text{Yes})$$

Attributes here are Marital status
Refund
~~Residence~~

Attributes here are discrete by nature.
Means they are either True / False or Yes / No.

Page No.		
Date		

so for discrete attributes

$$P(A_i | C_k) = \frac{|A_{ik}|}{N_{ke}}$$

where $|A_{ik}|$ is the number of instances having attributes A's and belongs to class C_k .

e.g.

$$P(\text{status} = \text{Married} | \text{No})$$

Means probability of instance who are married and belongs to class No

so the probability will be .

$$P(\text{Married} | \text{No}) = \frac{4}{7}$$

from the dataset

$$n(\text{Married}) = 4.$$

$$n(\text{No}) = 7.$$

similarly

$$P(\text{Refund} = \text{Yes} | \text{Yes})$$

Means probability of instance refund with Yes belonging to class Yes.

$$P(\text{Refund} = \text{Yes} | \text{Yes}) = \frac{0}{3} = 0$$

** simply count the instance with the class stated **

so the question was to predict the CLASS for

Refund = No

Marital status = Married

Amount = 120K

for this Naive Bayes Probability

$$P(X | \text{Class} = \text{No})$$

$$= P(\text{Refund} = \text{No} | \text{Class} = \text{No}) *$$

$$P(\text{Married} | \text{Class} = \text{No}) *$$

$$P(\text{Amount} = 120\text{K} | \text{Class} = \text{No})$$

↳ Means probability of given instance X to belong class No.

Calculated as

Similarly

Probability of instance X belonging to class Yes is given as.

$$P(X \mid \text{class} = \text{Yes})$$

$$= P(\text{Refund} = \text{No} \mid \text{class} = \text{Yes}) * \\ P(\text{Married} \mid \text{class} = \text{Yes}) * \\ P(\text{Income} = 120\text{K} \mid \text{class} = \text{Yes})$$

** Calculate for both class i.e.
YES & NO.

- Probability of Yes & No.

- If whichever probability is with class YES or NO comes greater the prediction falls for that class.

CALCULATION.

$$P(\text{Refund} = \text{No} \mid \text{class} = \text{No})$$

Count Refund as No within class NO.
from the given dataset

$$n(\text{Refund} = \text{No}) = 4.$$

$$n(\text{Class} = \text{No}) = 7.$$

$$\Rightarrow P(\text{Refund} = \text{No} | \text{Class} = \text{No}) = \frac{4}{7}$$

$$P(\text{Married} | \text{Class} = \text{No})$$

$$n(\text{Married}) = 4.$$

$$n(\text{Class} = \text{No}) = 7.$$

$$\Rightarrow P(\text{Married} | \text{Class} = \text{No}) = \frac{4}{7}.$$

$$P(\text{Account} = 120K | \text{Class} = \text{No}).$$

From dataset, Account is not discrete.
i.e. the values are distributive
so we have to go for probability distribution.

for this we need to find normal distribution.

NORMAL DISTRIBUTION

$$P(A_i | c_j) = \frac{1}{\sqrt{2\pi \sigma_{ij}^2}} e^{\left(\frac{(A_i - \mu_{ij})^2}{2\sigma_{ij}^2} \right)}$$

Where

μ_{ij} = mean.

σ_{ij}^2 = variance.

e.g. How to calculate mean & variance

Suppose given data set

$$X = (5, 8, 7, 6, 9)$$

$$\sigma^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2$$

\bar{x} = mean of samples.

n = no. of samples.

$$\therefore \bar{x} = \frac{5+8+7+6+9}{5}$$

$$\boxed{\bar{x} = 7}$$

$$\sigma^2 = \frac{(5-7)^2 + (8-7)^2 + (7-7)^2 + (6-7)^2 + (9-7)^2}{5-1}$$

$$\boxed{\sigma^2 = 2.5}$$

In our case .

$$P(Amount = 120K \mid Class = No)$$

$$= \frac{1}{\sqrt{2\pi\sigma_{ij}^2}} e^{-\frac{(A_i - \mu_{ij})^2}{2\sigma_{ij}^2}}$$

$$A_i = 120K$$

μ_{ij} = Mean of all the amounts belonging to class No

$$= \frac{125 + 100 + 70 + 120 + 60 + 220 + 75}{7}$$

Total no. of Class No = 7.

$$\mu_{ij} = \frac{770}{7} = 110K.$$

(~~Ex~~ Amount corresponding to class No is taken for mean)

$$\begin{aligned} \sigma_{ij}^2 &= (125 - 110)^2 + (100 - 110)^2 + (70 - 110)^2 \\ &\quad + (120 - 110)^2 + (60 - 220)^2 + (75 - 110)^2 \\ &\quad + (60 - 110)^2 + (220 - 110)^2 \end{aligned}$$

$$\begin{aligned} &= (15)^2 + (-10)^2 + (-40)^2 + (10)^2 + \\ &\quad (-50)^2 + (110)^2 - (-35)^2 \end{aligned}$$

Page No.			
Date			

$$= \underline{17850}$$

6

$$\sigma_{ij}^2 = 2975$$

∴ $P(\text{Amount } \neq 120 \text{K} \mid \text{Class} = \text{No})$

$$= \frac{1}{\sqrt{2\pi(2975)}} e^{\left(\frac{(120-110)}{2(2975)}\right)}$$

$$= 0.0072$$

$\Rightarrow P(X \mid \text{Class} = \text{No})$

$$= \frac{4}{7} \times \frac{4}{7} \times 0.0072$$

$$= 0.0023$$

So the conditional probability
for instance X for class No. is
0.0023.

Now, calculate for class Yes 150

Page No.		
Date		

$$P(X \mid \text{Class} = \text{Yes})$$

$$= P(\text{Refund} = \text{No} \mid \text{Class} = \text{Yes}) *$$

$$P(\text{Married} \mid \text{Class} = \text{Yes}) *$$

$$P(\text{Amount} = 125K \mid \text{Class} = \text{Yes})$$

$$\Rightarrow P(\text{Refund} = \text{No} \mid \text{Class} = \text{Yes})$$

$$= \frac{3}{3} = n(\text{Yes}) = 3.$$

$$P(\text{Married} \mid \text{Class} = \text{Yes})$$

$$= \frac{0}{3} = 0.$$

$$P(\text{Amount} = 125K \mid \text{Class} = \text{Yes})$$

$$A_i = 125K$$

$$N_{ij} = \frac{95 + 85 + 90}{3} = 90.$$

$$\sigma_{ij}^2 = \frac{(95 - 90)^2 + (85 - 90)^2 + (90 - 90)^2}{3}.$$

$$= 25$$

$$\Rightarrow P(\text{Amount} = 125K \mid \text{Class} = \text{Yes})$$

$$= \frac{1}{\sqrt{2\pi\sigma_{ij}^2}} e^{-\left(\frac{A_i - \mu_{ij}}{2\sigma_{ij}^2}\right)^2}$$

$$= \frac{1}{\sqrt{2\pi} 25} e^{\left(\frac{125 - 90}{2 \times 25} \right)}$$

$$= \frac{1}{\sqrt{2\pi} 25} e^{\frac{-125 + 90}{50}}$$

$$= 1.2 \times 10^{-9}$$

~~Thus~~ $P(X | \text{Class} = \text{Yes})$

$$= \frac{3}{3} \times \frac{0}{2} \times 1.2 \times 10^{-9}$$

$$= 0$$

$$P(X | N_0) \cdot P(N_0)$$

$$= 0.0023 \times \frac{7}{10}$$

$$=$$

$$P(X | \text{Yes}) \cdot P(\text{Yes})$$

$$= 0 \times \frac{3}{10}$$

$$= 0$$

$$\therefore P(\text{No} | X) > P(\text{Yes} | X)$$

\Rightarrow The prediction for the instance X is No

i.e. Refund = No.

Marital status = Married

Amount = 120K.

Eva/de/Class = No.

Assignment
(Do submit when the college reopens.)

Ques

For the above dataset find the prediction for instance X

$X = (\text{Refund} = \text{Yes}; \text{status} = \text{single}; \text{Amount} = 100\text{K})$

Ques

Give the Naïve Bayes algorithm

Ques

What are the pros & cons of Naïve Bayes algorithm.