

NAIVE BAYESIAN CLASSIFICATION

Pradnya Deorukhkar
MSc(Computer Science)
Sem II
pradnya.15d@gmail.com

Mrs. Dipali Meher MCS,
MPhil, NET
mailto:meher@gmail.com

Agenda

- Introduction
- Solved Example
- Advantage
- Disadvantages
- References

Introduction

Naive Bayesian classification algorithm is very simple which assumes that the classification attributes are independent and they do not have any correlation between them. Let D be a training set of tuples and their associated class labels, and each tuple is represented by an n -D attribute vector $X = (x_1, x_2, \dots, x_n)$. Suppose there are m classes C_1, C_2, \dots, C_m . Classification is to derive the maximum $P(C_i|X)$.

This can be derived from Bayes' theorem

$$P(C_i/X) = \frac{P(X/C_i) * P(C_i)}{P(X)}$$

$$P(X)$$

EXAMPLE 1

AGE	INCOME	STUDENT	CREDIT_RATING	BUYS_COMPUTER
<=30	High	No	Fair	No
<=30	High	No	Excellent	No
31...40	High	No	Fair	Yes
>40	Medium	No	Fair	Yes
>40	Low	Yes	Fair	Yes
>40	Low	Yes	Excellent	No
31...40	Low	Yes	Excellent	Yes
<=30	Medium	No	Fair	No
<=30	Low	Yes	Fair	Yes
>40	Medium	Yes	Fair	Yes
<=30	Medium	Yes	Excellent	Yes
31...40	Medium	No	Excellent	Yes
31...40	High	Yes	Fair	Yes
>40	Medium	No	Excellent	No

C1 : buys_computer = "yes"

C2 : buys_computer = "no"

Data sample X = (age <=30, Income = medium, Student = yes Credit_rating = Fair)

A data sample is given to us here and we have to find whether the person buys a computer or no using Naive Bayesian classification. So first we calculate the probability of buys_computer. So the numbers of yes values divided upon total records that is 14 gives us the probability of buys_computer with yes and same with no.

$P(C_i) : P(\text{buys_computer} = \text{"yes"}) = 9/14 = 0.643$

$P(\text{buys_computer} = \text{"no"}) = 5/14 = 0.357$

The total number of records in the table are 14.

Compute $P(X|C_i)$ for each class

$$P(\text{age} = "<=30" \mid \text{buys_computer} = \text{"yes"}) = 2/9 = 0.222$$

$$P(\text{age} = "<= 30" \mid \text{buys_computer} = \text{"no"}) = 3/5 = 0.6$$

$$P(\text{income} = \text{"medium"} \mid \text{buys_computer} = \text{"yes"}) = 4/9 = 0.444$$
$$P(\text{income} = \text{"medium"} \mid \text{buys_computer} = \text{"no"}) = 2/5 = 0.4$$

$$P(\text{student} = \text{"yes"} \mid \text{buys_computer} = \text{"yes"}) = 6/9 = 0.667$$

$$P(\text{student} = \text{"no"} \mid \text{buys_computer} = \text{"no"}) = 1/5 = 0.2$$

$$P(\text{credit_rating} = \text{"fair"} \mid \text{buys_computer} = \text{"yes"}) = 6/9 = 0.667$$
$$P(\text{credit_rating} = \text{"fair"} \mid \text{buys_computer} = \text{"no"}) = 2/5 = 0.4$$

Multiplying all the probabilities with yes values and no values from above separately.

$P(X|C_i)$:

$$\begin{aligned} P(X|\text{buys_computer} = \text{"yes"}) \\ = 0.222 \times 0.444 \times 0.667 \times 0.667 \\ = 0.044 \end{aligned}$$

$$\begin{aligned} P(X|\text{buys_computer} = \text{"no"}) \\ = 0.6 \times 0.4 \times 0.2 \times 0.4 \\ = 0.019 \end{aligned}$$

Now, multiply these probabilities with the above calculated $P(C_i)$

$P(X|C_i) * P(C_i)$:

$$\begin{aligned} P(X|\text{buys_computer} = \text{"yes"}) * P(\text{buys_computer} = \text{"yes"}) = \\ 0.044 \times 0.643 = 0.028 \end{aligned}$$

$$\begin{aligned} P(X|\text{buys_computer} = \text{"no"}) * P(\text{buys_computer} = \text{"no"}) = \\ 0.019 \times 0.357 = 0.007 \end{aligned}$$

The maximum value is 0.028.

Therefore, X belongs to class ("buys_computer = yes")

Example 2

CHILLS	RUNNY NOSE	HEADACH E	FEVER	FLU?
Y	N	Mild	Y	N
Y	Y	No	N	Y
Y	N	Strong	Y	Y
N	Y	Mild	Y	Y
N	N	No	N	N
N	Y	Strong	Y	Y
N	Y	Strong	N	N
Y	Y	Mild	Y	Y

C1 :Flu = “yes”

C2 : Flu = “no”

Data sample X = (chills=Y, runnynose=N,
headache=mild, fever=Y)

A data sample is given to us here and we have to find whether the person has flu or no using Naive Bayesian classification. So first we calculate the probability of flu. So the numbers of yes values divided upon total records that is 8 gives us the probability of flu with yes and same with no.

$P(C_i) : P(\text{Flu} = \text{“yes”}) = 5/8 = 0.62$

$P(\text{Flu} = \text{“no”}) = 3/10 = 0.37$

The total number of records in the table are 8.

Compute $P(X|C_i)$ for each class

$$P(\text{chills}=\text{"Y"} \mid \text{Flu} = \text{"yes"}) = 3/5 = 0.6$$

$$P(\text{chills} = \text{"Y"} \mid \text{Flu} = \text{"no"}) = 1/3 = 0.33$$

$$P(\text{runny nose}=\text{"N"} \mid \text{Flu} = \text{"yes"}) = 1/5 = 0.2$$

$$P(\text{runny nose}=\text{"N"} \mid \text{Flu} = \text{"no"}) = 2/3 = 0.66$$

$$P(\text{headache}=\text{"mild"} \mid \text{Flu} = \text{"yes"}) = 2/5 = 0.4$$

$$P(\text{headache}=\text{"mild"} \mid \text{Flu} = \text{"no"}) = 1/3 = 0.33$$

$$P(\text{fever}=\text{"Y"} \mid \text{Flu} = \text{"yes"}) = 4/5 = 0.8$$

$$P(\text{fever}=\text{"Y"} \mid \text{Flu} = \text{"no"}) = 1/3 = 0.33$$

Multiplying all the probabilities with yes values and no values from above separately.

$P(X|C_i)$:

$P(X|\text{Flu} = \text{"yes"})$

$$= 0.6 \times 0.4 \times 0.2 \times 0.8$$

$$= 0.0384$$

$P(X|\text{Stolen} = \text{"no"})$

$$= 0.33 \times 0.66 \times 0.33 \times 0.33$$

$$= 0.024$$

Now, multiply these probabilities with the above calculated $P(C_i)$

$P(X|C_i) * P(C_i)$:

$$P(X|\text{Flu} = \text{"yes"}) * P(\text{Flu} = \text{"yes"}) = 0.0384 \times 0.62 = 0.023$$

$$P(X|\text{Flu} = \text{"no"}) * P(\text{Flu} = \text{"no"}) = 0.024 \times 0.37 = 0.0088$$

The maximum value is 0.023.

Therefore, X belongs to class ("Flu=yes")

Example 3

NO	COLOR	TYPE	ORIGIN	STOLEN?
1	Red	Sports	Domestic	Yes
2	Red	Sports	Domestic	No
3	Red	Sports	Domestic	Yes
4	Yellow	Sports	Domestic	No
5	Yellow	Sports	Imported	Yes
6	Yellow	SUV	Imported	No
7	Yellow	SUV	Imported	Yes
8	Yellow	SUV	Domestic	No
9	Red	SUV	Imported	No
10	Red	Sports	Imported	Yes

C1 : Stolen = “yes”

C2 : Stolen = “no”

Data sample X = (color=red, type=SUV,
Origin=Domestic)

A data sample is given to us here and we have to find whether the person has stolen or no using Naive Bayesian classification. So first we calculate the probability of stolen. So the numbers of yes values divided upon total records that is 10 gives us the probability of stolen with yes and same with no.

$P(C_i) : P(\text{Stolen} = \text{“yes”}) = 5/10 = 0.5$

$P(\text{Stolen} = \text{“no”}) = 5/10 = 0.5$

The total number of records in the table are 10.

Compute $P(X|C_i)$ for each class

$$P(\text{color}=\text{"red"} \mid \text{Stolen} = \text{"yes"}) = 3/5 = 0.6$$

$$P(\text{color}=\text{"red"} \mid \text{Stolen} = \text{"no"}) = 2/5 = 0.4$$

$$P(\text{type}=\text{"SUV"} \mid \text{Stolen} = \text{"yes"}) = 1/5 = 0.2$$

$$P(\text{type}=\text{"SUV"} \mid \text{Stolen} = \text{"no"}) = 3/5 = 0.6$$

$$P(\text{origin}=\text{"Domestic"} \mid \text{Stolen} = \text{"yes"}) = 2/5 = 0.4$$

$$P(\text{origin}=\text{"Domestic"} \mid \text{Stolen} = \text{"no"}) = 3/5 = 0.6$$

Multiplying all the probabilities with yes values and no values from above separately.

$P(X|C_i)$:

$P(X|\text{Stolen} = \text{"yes"})$

$= 0.6 \times 0.4 \times 0.2$

$= 0.048$

$P(X|\text{Stolen} = \text{"no"})$

$= 0.6 \times 0.4 \times 0.4$

$= 0.096$

Now, multiply these probabilities with the above calculated $P(C_i)$

$P(X|C_i) * P(C_i)$:

$P(X|\text{Stolen} = \text{"yes"}) * P(\text{Stolen} = \text{"yes"}) = 0.048 \times 0.5 = 0.024$

$P(X|\text{Stolen} = \text{"no"}) * P(\text{Stolen} = \text{"no"}) = 0.096 \times 0.5 = 0.048$

The maximum value is 0.048.

Therefore, X belongs to class ("Stolen= no")

Example 4

OUTLOOK	TEMP	HUMIDITY	WINDY	PLAY GOLF?
Rainy	Hot	High	False	No
Rainy	Hot	High	True	No
Overcast	Hot	High	False	Yes
Sunny	Mild	High	False	Yes
Sunny	Cool	Normal	False	Yes
Sunny	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Rainy	Mild	High	False	No
Rainy	Cool	Normal	False	Yes
Sunny	Mild	Normal	False	Yes
Rainy	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Sunny	Mild	High	True	No

C1 : Play Golf = “yes”

C2 : Play Golf = “no”

Data sample X = (outlook=rainy, temp=hot, humidity=high, windy=false)

A data sample is given to us here and we have to find whether the person has play golf or no using Naive Bayesian classification. So first we calculate the probability of play golf. So the numbers of yes values divided upon total records that is 14 gives us the probability of playing golf with yes and same with no.

$P(C_i) : P(\text{Play golf} = \text{“yes”}) = 9/14 = 0.643$

$P(\text{Play golf} = \text{“no”}) = 5/14 = 0.357$

The total number of records in the table are 14.

Compute $P(X|C_i)$ for each class

$$P(\text{outlook}=\text{"rainy"} \mid \text{Play Golf} = \text{"yes"}) = 2/9 = 0.222$$

$$P(\text{outlook}=\text{"rainy"} \mid \text{Play Golf} = \text{"no"}) = 3/5 = 0.6$$

$$P(\text{temp}=\text{"hot"} \mid \text{Play Golf} = \text{"yes"}) = 2/9 = 0.222$$

$$P(\text{temp}=\text{"hot"} \mid \text{Play Golf} = \text{"no"}) = 2/5 = 0.4$$

$$P(\text{humidity}=\text{"high"} \mid \text{Play Golf} = \text{"yes"}) = 3/9 = 0.33$$

$$P(\text{humidity}=\text{"high"} \mid \text{Play Golf} = \text{"no"}) = 4/5 = 0.8$$

$$P(\text{windy}=\text{"false"} \mid \text{Play Golf} = \text{"yes"}) = 6/9 = 0.666$$

$$P(\text{windy}=\text{"false"} \mid \text{Play Golf} = \text{"no"}) = 2/5 = 0.4$$

Multiplying all the probabilities with yes values and no values from above separately.

$P(X|C_i)$:

$P(X|\text{play golf} = \text{"yes"})$

$= 0.222 \times 0.222 \times 0.333 \times 0.666$

$= 0.0109$

$P(X|\text{play golf} = \text{"no"})$

$= 0.6 \times 0.4 \times 0.8 \times 0.4$

$= 0.0768$

Now, multiply these probabilities with the above calculated $P(C_i)$

$P(X|C_i) * P(C_i)$:

$P(X|\text{play golf} = \text{"yes"}) * P(\text{play golf} = \text{"yes"}) =$

$0.0109 \times 0.643 = 0.007$

$P(X|\text{play golf} = \text{"no"}) * P(\text{play golf} = \text{"no"}) =$

$0.0768 \times 0.357 = 0.0274$

The maximum value is 0.0274.

Therefore, X belongs to class ("play golf = no")

Example 5

PATIENT	DISEASE	SUGAR LEVEL	SURVIVAL CHANCES
Small	Serious	High	Yes
Medium	Normal	Low	Yes
Senior	Lifetime	Normal	Yes
Small	Lifetime	High	No
Small	Normal	High	Yes
Senior	Serious	Normal	No
Medium	Serious	Low	Yes
Senior	Normal	Low	No
Medium	Lifetime	Normal	Yes
Medium	Serious	High	No
Senior	Normal	Low	No

C1 : survival chances = “yes”

C2 : survival chances = “no”

Data sample X = (age=senior, disease=normal, sugar level=normal)

A data sample is given to us here and we have to find whether the person has survival chances or no using Naive Bayesian classification. So first we calculate the probability of survival chances. So the numbers of yes values divided upon total records that is 11 gives us the probability of survival chances with yes and same with no.

$P(C_i) : P(\text{survival chances} = \text{“yes”}) = 6/11 = 0.545$

$P(\text{survival chances} = \text{“no”}) = 5/11 = 0.455$

The total number of records in the table are 11.

Compute $P(X|C_i)$ for each class

$$P(\text{age}=\text{"senior"} \mid \text{survival chances} = \text{"yes"}) = 1/6 = 0.166$$

$$P(\text{age}=\text{"senior"} \mid \text{survival chances} = \text{"no"}) = 3/5 = 0.6$$

$$P(\text{disease}=\text{"normal"} \mid \text{survival chances} = \text{"yes"}) = 2/6 = 0.33$$

$$P(\text{disease}=\text{"normal"} \mid \text{survival chances} = \text{"no"}) = 2/5 = 0.4$$

$$P(\text{sugar level}=\text{"normal"} \mid \text{survival chances} = \text{"yes"}) = 2/6 = 0.33$$

$$P(\text{sugar level}=\text{"normal"} \mid \text{survival chances} = \text{"no"}) = 1/5 = 0.2$$

Multiplying all the probabilities with yes values and no values from above separately.

$P(X|C_i)$:

$P(X|\text{survival chances} = \text{"yes"})$

$$= 0.166 \times 0.33 \times 0.33$$

$$= 0.018$$

$P(X|\text{survival chances} = \text{"no"})$

$$= 0.6 \times 0.4 \times 0.2$$

$$= 0.048$$

Now, multiply these probabilities with the above calculated $P(C_i)$

$P(X|C_i) * P(C_i)$:

$P(X|\text{survival chances} = \text{"yes"}) * P(\text{survival chances} = \text{"yes"}) =$

$$0.018 \times 0.545 = 0.0098$$

$P(X|\text{survival chances} = \text{"no"}) * P(\text{survival chances} = \text{"no"}) =$

$$0.048 \times 0.455 = 0.02184$$

The maximum value is 0.02184.

Therefore, X belongs to class ("survival chances = no")

Advantages

- Very simple, easy to implement and fast. Need less training data.
- It is highly scalable.
- It handles continuous and discrete data.
- Not sensitive to irrelevant features.
- It is useful in natural processing language.
- It can be easily updateable if we want to add training data.
- It is not a complicated algorithm.

Disadvantages

- The first disadvantage is that the Naive Bayes classifier makes an assumption on the shape of your data distribution, i.e. any two features are independent given the output class. Due to this, the result can be very bad.
- Another problem happens due to data scarcity. For any possible value of a feature, you need to estimate a likelihood value by an approach.
- This can result in probabilities going towards 0 or 1, which in turn leads to numerical instabilities and worse results will be seen.
- With Naive-Bayes, if you do have a class label and a certain attribute value together then the probability estimate will be zero.

References

- 1) www.ijserd.com/articles/IJSRDV4I70116.pdf
- 2) <https://www.researchgate.net>
- 3) *GEEKS FOR GEEKS*
- 4) Indian Journal of Science and Technology, Vol 8(16), DOI: 10.17485/ijst/2015/v8i16/62055, July 2015.
- 5) www.ijraset.com

THANK YOU