

BAYES CLASSIFIER

Prior Probability:- It is the probability of event before collection of new data.

Posterior Probability:- It is the updated (revised) probability that considers additional new information.

e.g:- Let we have 3 pieces of land.



Out of these three pieces of land, one has oil inside the surface.

∴ Probability of getting oil from land c before digging starts is $P(c) = 1/3$. This is called as Prior Probability.

Now, let land A is digged, & we didn't get oil. Hence now the probability of getting oil from land c becomes $P(c) = 1/2$ as now only B & c are remaining & we got information that A doesn't have oil. This is called as Posterior Probability.

BAYES RULE:

It gives posterior probability if prior probabilities are known & class conditional probability is known.

It is given by

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

e.g.: In a class of students, 60% are boys & 40% are girls. All boys carry black color mobiles whereas girls carry black color & white color mobiles in equal numbers. Let a student is selected randomly & found that the mobile carried was black. What is the probability that the student is girl?



$$\rightarrow P(\text{Boys}) = 0.6 \quad P(\text{Girls}) = 0.4$$

$$P(\text{Black} | \text{Girl}) = 0.5 \quad P(\text{Black}) = 0.8$$

$$\therefore P(G | \text{Black}) = ?$$

$$\therefore P(G | \text{Black}) = \frac{P(\text{Black} | \text{Girl}) \cdot P(\text{Girl})}{P(\text{Black})}$$

$$= \frac{0.5 \times 0.4}{0.8}$$

$$\boxed{P(G | \text{Black}) = 0.25.}$$

On similar basis, let us say we have two classes C_1 & C_2 .

From Bayes rule.

$$P(w_i | x) = \frac{P(x | w_i) \cdot P(w_i)}{P(x)}$$

where

$P(x) \xrightarrow{\text{prob.}} \text{event has occurred.}$

$$\frac{(A)_q \cdot (A|B)_q}{(B)_q} = (B|A)_q$$

let us say we want to check, the given vector x (vector feature vector) belongs to which class. For this we check what is the probability of given vector x (vector term is used because it may have more than one variable) lying in either class C_1 & in class C_2

i.e $P(C_1|x) \rightarrow$ probability of given x lying in class C_1 ,

$$P(C_2|x) \rightarrow -11 - \text{class } C_2.$$

whichever probability is more, that x vector belongs to that class.

$$\text{i.e if } P(C_1|x) > P(C_2|x)$$

x belongs to C_1 ,

$$\text{or if } P(C_1|x) < P(C_2|x)$$

x belongs to C_2 .

Now by Bayes theorem. $P(C_1|x)$ can be determined by

$$P(C_1|x) = \frac{P(x|C_1) \cdot P(C_1)}{P(x)}$$

$$P(C_2|x) = \frac{P(x|C_2) \cdot P(C_2)}{P(x)}$$

Comparing

$$P(c_1|x) > P(c_2|x)$$

i.e.

$$\frac{P(x|c_1) \cdot P(c_1)}{P(x)} > \frac{P(x|c_2) \cdot P(c_2)}{P(x)}$$

The denominator $P(x)$ can be cancelled out as the result of comparison does not change without it.

It is called as normalizing constant.

Thus for comparison we can use only

$$P(x|c_1) \cdot P(c_1) > P(x|c_2) \cdot P(c_2)$$

x belongs to c_1

↳ $P(x|c_1) \cdot P(c_1) < P(x|c_2) \cdot P(c_2)$

x belongs to c_2 .

Ex:- consider the following data set. There are total 4 variables (outlook, temperature, humidity & windy). There are two classes play - Yes & play - No. Each data instance will have 4 attributes.

$$x = \{x_1, x_2, x_3, x_4\}$$

\downarrow \downarrow \downarrow \downarrow
 outlook temp num windy

Note:- Hence x is called as feature vector.

NOW e.g. consider 1st instance

$$x = \{ \text{sunny, hot, high, false} \}. This belongs$$

to class - play - NO.

Similarly for $X = \{\text{Rainy, cool, normal, False}\}$, the class is Yes. {see data serial no. ⑤}.

But we want to find the class of

$X = \{\text{Rainy, cool, High, True}\}$ i.e. this

if these are the weather conditions whether to play or not to play. From the data set we cannot determine the class. Hence we have to find its class i.e under such conditions whether to play or not to play.

Data set - DS1

S.NO	outlook	Temperature	Humidity	windy	Play
1	sunny	Hot	high	False	NO
2	sunny	Hot	high	True	NO
3	overcast	Hot	high	False	Yes
4	Rainy	mild	high	False	Yes
5	Rainy	cool	normal	False	Yes
6	Rainy	cool	normal	True	NO
7	overcast	cool	normal	True	Yes
8	sunny	mild	high	False	NO
9	Sunny	cool	normal	False	Yes
10	Rainy	mild	normal	False	Yes
11	sunny	mild	normal	True	Yes
12	overcast	mild	high	True	Yes
13	overcast	hot	normal	False	Yes
14	Rainy	mild	high	True	NO

By Bayes theorem

$$P(c_i|x) = \frac{P(x|c_i) \cdot P(c_i)}{P(x)}$$

we can ignore $P(x)$

i.e

$$P(c_1|x) = P(x|c_1)P(c_1)$$

$$\therefore P(c_2|x) = P(x|c_2)P(c_2)$$

We have to compare whether

$$P(x|c_1) \cdot P(c_1) > P(x|c_2) \cdot P(c_2)$$

where

$$P(c_1) \rightarrow P(\text{Yes})$$

let $c_1 \rightarrow \text{class Yes}$

$$P(c_2) \rightarrow P(\text{No}).$$

$c_2 \rightarrow \text{class No}$

From the data there are total 9 Yes & 5 No
out of 14 values

$$\therefore P(c_1) = P(\text{Yes}) = 9/14 \quad \} - \text{equation } ①$$

$$P(c_2) = P(\text{No}) = 5/14 \quad \}$$

Now when we have to find $P(x|c_i)$, it is very expensive activity. Because we have to consider all the combinations.

If we consider outlook as sunny then temperature can be hot, mild or cool also for these combn again humidity can be high, normal & again for all these combn windy can be false or true. Similarly for outlook = overcast again so many combinations. Thus the

⑦

complexity again increases if n (variables) increases.
 In order to simplify the calculations, we assume
 For each class the attributes are independent.
 The classifier resulting from this assumption is
 called as Naive Bayes classifier.

with Naive Bayes classifier we can write

$$P(x|c_i) = \prod_{k=1}^n P(x_k|c_i)$$

i.e The product of probabilities of each of values
 of attributes of x for given class c .

i.e The probability of outlook = sunny for
 play = No will be

$$P(x_1 = \text{sunny}|c_2) = 3/5$$

i.e out of 5 Nos there are 3 sunny. outlook.

let us verify with naive classifier the class out for
 vector $x = \{\text{sunny, hot, high, false}\}$.

→ For this we have to compare

$$P(x|c_1) \cdot P(c_1) \quad \& \quad P(x|c_2) \cdot P(c_2)$$

we know $P(c_1) = 9/14$ $P(c_2) = 5/14$

NOW for class c_1 , i.e play = Yes

$$P(x|c_1) = ?$$

$$P(x = \{\text{sunny, hot, high, false}\}|c_1) = ?$$

: finding separate probabilities & then taking their product. finding for class C_1 , i.e Yes

$$\text{i.e } P(x = \text{sunny} | c_1) = 2/9$$

: There are 9 Yes & out of 9 there are 2 sunny.

Similarly

$$P(x = \text{hot} | c_1) = 2/9$$

$$P(x = \text{high} | c_1) = 3/9$$

$$P(x = \text{False} | c_1) = 6/9$$

$$\therefore P(x | c_1) = P(x = \text{sunny} | c_1) * P(x = \text{hot} | c_1) \\ * P(x = \text{high} | c_1) * P(\text{False} | c_1)$$

$$= \frac{2}{9} * \frac{2}{9} * \frac{3}{9} * \frac{6}{9}$$

$$\Rightarrow P(x | c_1) = 0.0109$$

$$\therefore P(c_1 | x) = P(x | c_1) * P(c_1)$$

$$= 0.0109 * \frac{9}{14}$$

$$\therefore \boxed{P(c_1 | x) = 0.0070} \quad - \textcircled{A}$$

i.e Probability of x belonging to class C_1 , i.e Play = Yes is 0.0070.

similarly now finding for class (C_2) i.e NO

$$P(x = \text{sunny} | c_2) = 3/15$$

$$P(x = \text{hot} | c_2) = 2/15$$

$$P(x = \text{high} | c_2) = 4/15$$

$$P(x = \text{False} | c_2) = 2/15$$

$$\therefore P(x|c_2) = P(x = \text{sunny} | c_2) * P(x = \text{hot} | c_2) \\ * P(x = \text{high} | c_2) * P(x = \text{False} | c_2)$$

$$= \frac{3}{5} * \frac{2}{5} * \frac{4}{5} * \frac{2}{5}$$

$$P(x|c_2) = 0.0768$$

$$P(c_2|x) = P(x|c_2) * P(c_2)$$

$$= 0.0768 * \frac{5}{14}$$

$$\boxed{P(c_2|x) = 0.02742} - \textcircled{B}$$

i.e Probability of x belonging to class c_2 i.e Play = No is 0.02742.

From \textcircled{A} & \textcircled{B}

$$P(c_2|x) > P(c_1|x)$$

\therefore The result is class play = No. i.e for given x vector there should not be any play.

We can see that our calculation is correct as we got the same result what is given in the data set.

But you can find the class for vect

$$x = \{\text{Rainy, cool, high, True}\} \text{ by same}$$

method. Note the class is not given in the dataset. You have to find it.

Now if the variables are with numeric values.

-then the normal distribution PDF should be evaluated
 e.g: consider the same previous data set with variables
 outlook, Temperature, Humidity & windy, only difference
 is now consider the variable outlook & windy same
 but temperature & humidity are with numeric values
 as follows:

Dataset - DS2

S.N	OUTLOOK	temperature	Humidity	windy	play
1	sunny	85	85	false	NO
2	sunny	80	90	true	NO
3	overcast	83	86	false	YES ←
4	Rainy	70	96	false	YES ←
5	Rainy	68	80	false	YES ←
6	Rainy	65	70	true	NO
7	overcast	64	65	true	YES ←
8	sunny	72	95	false	NO
9	sunny	69	70	false	YES ←
10	Rainy	75	80	false	YES ←
11	sunny	75	70	true	YES ←
12	overcast	72	90	true	YES ←
13	overcast	81	75	false	YES ←
14	Rainy	71	91	false	NO

find the class for vector $x = \{ \text{Rainy}, \text{temp} = 60, \text{Humidity} = 62, \text{windy} = \text{false} \}$

we have to compare

$$P(c_1|x) \& P(c_2|x)$$

i.e.

$$P(x|c_1) \cdot P(c_1) \& P(x|c_2) \cdot P(c_2)$$

we know:

$$P(c_1) = 9/14 \quad P(c_2) = 5/14$$

For

$P(x|c_1)$ i.e for class c_1 , i.e Play = Yes.

$$X = \{ \text{rainy}, \text{temp} = 60, \text{Humidity} = 62, \text{windy} = \text{false} \}$$

$$P(x|c_1) = P(x = \text{rainy}|c_1) * P(x = \text{temp} = 60 | c_1) \\ * P(x = \text{humid} = 62 | c_1) * P(x = \text{windy} = \text{false})$$

$$P(x = \text{rainy}|c_1) = 3/9 \quad \text{---(1)}$$

$P(x = \text{temp} = 60 | c_1)$ Now here it is numeric value

Hence probability cannot be taken directly. Therefore

finding PDF. For different classes.

for c_1

\therefore PDF is given by $-\frac{(x - \mu)^2}{2\sigma^2}$

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x - \mu)^2}{2\sigma^2}}$$

\therefore for variable temp $-\frac{(x - \mu_{\text{temp}})^2}{2\sigma_{\text{temp}}^2}$

$$f(x = \text{temp}) = \frac{1}{\sigma_{\text{temp}} \sqrt{2\pi}} e^{-\frac{(x - \mu_{\text{temp}})^2}{2\sigma_{\text{temp}}^2}}$$

$$\mu_{\text{temp}} = \frac{83 + 70 + 68 + 64 + 69 + 75 + 75 + 72 + 81}{9} \rightarrow \text{for Yes only.}$$

$$\boxed{\mu_{\text{temp}}|_{c_1} = 73}$$

$$\text{similarly } \sigma_{\text{temp}|c_1}^2 = \frac{(83 - 73)^2 + (70 - 73)^2 + \dots + (81 - 73)^2}{8}$$

$$\therefore \boxed{\sigma_{\text{temp}} \approx 6.16}$$

similarly

$$\therefore f(x = \text{temp} = 60) = \frac{1}{6.16 \sqrt{2\pi}} e^{-\frac{(60 - 73)^2}{2(6.16)^2}}$$

$$f(x = \text{temp} = 60 | c_1) = 0.00698. \quad - \textcircled{2}$$

similarly finding humidity = 62 for c_1 , we get.

$$f(x = \text{humidity} = 62 | c_1) = 0.0096. \quad - \textcircled{3}$$

$$P(x | c_1) =$$

$$P(x = \text{windy} = \text{false} | c_1) = 6/9 = \textcircled{4}$$

$$P(x | c_1) = \textcircled{1} \times \textcircled{2} \times \textcircled{3} \times \textcircled{4}$$

$$= \frac{3}{9} \times 0.00698 \times 0.0096 \times \frac{6}{9}$$

$$P(x | c_1) = 0.01489$$

$$P(c_1 | x) = P(x | c_1) \times P(c_1)$$

$$= 0.01489 \times \frac{9}{14}$$

$$P(c_1 | x) = 9.572 \times 10^{-6} \quad - \textcircled{A}$$

Now proceeding similarly way for class $c_2 = \text{NO}$

we get

$$P(x = \text{rainy} | c_2) = 2/5 + 8/2 + 0/5 + 8/8 = 9/15$$

$$P(x = \text{temp} | c_2) = 0.0094$$

$$P(x = \text{hum} = 62 | c_2) = 0.0018$$

$$P(x = \text{windy} = \text{false} | c_2) = 2/5 = 0.4$$

$$\therefore P(x|c_2) = \frac{2}{5} \times 0.0094 \times 0.0018 \times \frac{2}{5}$$

$$= 2.7072 \times 10^{-6}$$

$$\therefore P(c_2|x) = P(x|c_2) \cdot P(c_2)$$

$$= 2.7072 \times 10^{-6} \times \frac{5}{14}$$

$$\boxed{P(c_2|x) = 0.9668 \times 10^{-6}} - \textcircled{B}$$

Comparing \textcircled{A} & \textcircled{B} we get

$$P(c_1|x) > P(c_2|x)$$

\therefore Given x belongs to class c_1 , i.e play = Yes.