

Learning classifiers based on Bayes' rule

$X \rightarrow \langle x_1, x_2, x_3, \dots, x_n \rangle \rightarrow$ Assume each of x_1, x_2, \dots are boolean variables
 $Y \rightarrow$ Boolean variable

$f: X \rightarrow Y \Rightarrow$ Given x predict $Y \Rightarrow P(Y/x)$

$$P(Y = y_i | X = x_k) = \frac{P(X = x_k | Y = y_i) P(Y = y_i)}{\sum_j P(X = x_k | Y = y_j) P(Y = y_j)} \rightarrow \text{(As per Bayes theorem)}$$

Estimating number of possibilities:-

If $n \rightarrow$ vector
 \downarrow

$$P(X = x_k | Y = y_i) \Rightarrow$$

$$P(x_1/x_2, x_3, \dots, x_n, y_i) * \underbrace{P(x_2, x_3, \dots, x_n * y_i)}_{P(x_2/x_3, \dots, x_n, y_i) * P(x_3, \dots, x_n, y_i)}$$

2^n possible vector combinations for x

2 possible values for y_i

Total $\rightarrow 2^{n+1}$ parameters must be estimated

$$n=3 \rightarrow P(x = x_k | y = y_i) \quad (x_k = \langle x_1, x_2, x_3 \rangle)$$

$$P(x_1/x_2, x_3, y_i) * P(x_2, x_3, y_i)$$

$$\Rightarrow P(x_1/x_2, x_3, y_i) * P(x_2/x_3, y_i)$$

$$\Rightarrow P(x_1/x_2, x_3, y_i) * P(x_2/x_3, y_i) * P(x_3/y_i) * P(y_i)$$

Assume this is the only training data

x_1	x_2	x_3	y_i
0	1	0	0
1	0	1	1

Given $\rightarrow (0, 0, 0)$

$$P(0/x_2=0, x_3=0, y_i=0) \rightarrow 0$$

$$P(x_2=0/x_3=0, y_i=0) \rightarrow 0$$

* you need to have at least 2^n datapoints.

* If $n=30 \rightarrow 2^{30} \rightarrow 500$ billion datapoints
Assuming $X \rightarrow$ Binary

It is really hard to calculate probability from training data as we

don't get that many datapoints in the training data. And if one probability

returns zero entire multiplication will return zero.

Naive Bayes

Naive bayes make a naive assumption that the features are conditionally independent

Coming to our previous e.g.:-

$$x_k = \langle x_1, x_2, x_3 \rangle$$

$$P(x = x_k | y = y_i)$$

$$P(x_1 | x_2, x_3, y_i) * P(x_2 | x_3, y_i) * P(x_3 | y_i) * P(y_i)$$

↓ Assuming conditional independent

$$P(x_1 | y_i) * P(x_2 | y_i) * P(x_3 | y_i) * P(y_i) \Rightarrow \boxed{P(y_i) \prod_{j=1}^n P(x_j | y_i)} \quad \text{--- ②}$$

With this assumption we only need to calculate 2n probabilities

$$P(Y = y_k | X_1 \dots X_n) = \frac{P(Y = y_k) P(X_1 \dots X_n | Y = y_k)}{\sum_j P(Y = y_j) P(X_1 \dots X_n | Y = y_j)}$$

Sub ② in the above equation

$$P(Y = y_k | X_1 \dots X_n) = \frac{P(Y = y_k) \prod_i P(X_i | Y = y_k)}{\sum_j P(Y = y_j) \prod_i P(X_i | Y = y_j)}$$

To convert this into a classification problem we need to find the maximum probability.

$$Y \leftarrow \arg \max_{y_k} \frac{P(Y = y_k) \prod_i P(X_i | Y = y_k)}{\sum_j P(Y = y_j) \prod_i P(X_i | Y = y_j)}$$

As the denominator does not change based on the value of y_k , it is redundant to calculate everytime.

$$Y \leftarrow \arg \max_{y_k} P(Y = y_k) \prod_i P(X_i | Y = y_k)$$

Understanding Train and test phase in NB

Sample data

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

$$Y \leftarrow \arg \max_{y_k} P(Y = y_k) \prod_i P(X_i | Y = y_k)$$

$y_k \rightarrow \{ \text{yes, no} \}$

$x_i \rightarrow \{ \text{outlook, temp, humidity, windy} \}$
 { Rainy, overcast, sunny }
 { high, normal }
 { hot, mild, cool }

$$* P(x_i = \text{rainy} / y = \text{yes}) = \frac{P(x_i = \text{rainy and } y = \text{yes})}{P(y = \text{yes})}$$

$$= \frac{2/14}{9/14} = \frac{2}{9}$$

In training, we calculate all the possible combinations of probability

$$P(y = \text{yes}) = \frac{9}{14} \quad P(y = \text{no}) = \frac{5}{14}$$

We can repeat the process for all the data.

$$P(x | c) = P(\text{Sunny} | \text{Yes}) = 3/9 = 0.33$$

Frequency Table		Play Golf	
		Yes	No
Outlook	Sunny	3	2
	Overcast	4	0
	Rainy	2	3

Likelihood Table		Play Golf		
		Yes	No	
Outlook	Sunny	3/9	2/5	5/14
	Overcast	4/9	0/5	4/14
	Rainy	2/9	3/5	5/14
		9/14	5/14	

$$P(x) = P(\text{Sunny}) = 5/14 = 0.36$$

$$P(c) = P(\text{Yes}) = 9/14 = 0.64$$

Train time complexity:-

$$O(n * d * c)$$

n = no of data points
 d = features
 c = classes

why n ?

In worst case a feature can have n -different possibilities

Space complexity

$$O(n * d * c)$$

Store $n * d * c$ cells (from lookup tables)

Frequency Table		Play Golf	
		Yes	No
Outlook	Sunny	3	2
	Overcast	4	0
	Rainy	2	3

Likelihood Table		Play Golf	
		Yes	No
Outlook	Sunny	3/9	2/5
	Overcast	4/9	0/5
	Rainy	2/9	3/5

Frequency Table		Play Golf	
		Yes	No
Humidity	High	3	4
	Normal	6	1

Likelihood Table		Play Golf	
		Yes	No
Humidity	High	3/9	4/5
	Normal	6/9	1/5

Frequency Table		Play Golf	
		Yes	No
Temp.	Hot	2	2
	Mild	4	2
	Cool	3	1

Likelihood Table		Play Golf	
		Yes	No
Temp.	Hot	2/9	2/5
	Mild	4/9	2/5
	Cool	3/9	1/5

Frequency Table		Play Golf	
		Yes	No
Windy	False	6	2
	True	3	3

Likelihood Table		Play Golf	
		Yes	No
Windy	False	6/9	2/5
	True	3/9	3/5

Testing phase

Given a point $x_q \rightarrow \langle \text{overcast}, \text{hot}, \text{high}, \text{True} \rangle \rightarrow \text{play Golf}$

$$p(y=\text{yes}/x_q) = p(\text{yes}) * p(\text{overcast}/\text{yes}) * p(\text{hot}/\text{yes}) * p(\text{high}/\text{yes}) * p(\text{True}/\text{yes})$$

$$= \frac{9}{14} \times \frac{4}{9} \times \frac{2}{9} \times \frac{3}{9} \times \frac{3}{9} = 0.3456$$

$$p(y=\text{NO}/x_q) = p(\text{NO}) * p(\text{overcast}/\text{NO}) * p(\text{hot}/\text{NO}) * p(\text{high}/\text{NO}) * p(\text{True}/\text{NO})$$

$$= \frac{5}{14} \times 0 \times \dots = 0$$

$$p(y=\text{yes}/x_q) > p(y=\text{NO}/x_q) \rightarrow x_q \rightarrow \text{Target Yes.}$$

Test time complexity

$$O(d * c)$$

↓
dimensions

no. of classes

Wait!! → Just because one feature resulted in zero probability -
entire posterior probability came to zero.
↓
This is not fair, what if the other features are important?

To solve this problem, we have Laplace smoothing.

$$p(\text{overcast}/\text{NO}) = \frac{\text{no. of points where Outlook is overcast And play golf = NO}}{\text{no. of points where play golf is NO}} = \frac{0}{5} = 0$$

↓ modified

$$= \frac{\text{no. of points where Outlook is overcast And play golf = NO} + \alpha}{\text{no. of points where play golf is NO} + c \alpha}$$

$\alpha \rightarrow$ hyper-param
 $c \rightarrow$ no. of classes
here $(c=2)$

$$= \frac{0+1}{5+2 \times 1} = \frac{1}{7} = 0.14 (\neq 0)$$

Here $\alpha \rightarrow$ hyperparam
 α is small \rightarrow overfitting
 α is large \rightarrow underfitting

As you can see the example on the right
higher values of alpha gives us almost a constant value ($\frac{1}{n}$)

(without Laplace smoothing)

$p(x/y)$	$\alpha=1$	$\alpha=100$
$2/3$	$\frac{2+1}{3+2} = \frac{3}{5}$	$\frac{2+100}{3+200} = \frac{102}{203} \approx 0.5$
1	$\frac{1+1}{1+2} = \frac{2}{3}$	$\frac{1+100}{1+200} = \frac{101}{201} \approx 0.5$