

$$\begin{aligned}
 & P(\text{doesn't like 'T' } \cap \text{ likes 'RJ'}) \\
 &= P(\text{doesn't like 'T' } \cap \text{ likes 'RJ'} | \text{ doesn't like 'T')} \\
 &\quad \times P(\text{likes 'RJ'}) \\
 &\hline
 & P(\text{doesn't like 'T'})
 \end{aligned}$$

So, we derived Bayes theorem.

Naive Bayes :

It is a classification technique based on Bayes theorem with an assumption of independence among predictors. In simple terms a Naive Bayes classifier assumes that the presence of a particular feature of a class is unrelated to the presence of any other feature.

For example a fruit may be considered to be an apple if it is red, round and about 3 inches in diameter. Even if these features depend upon each other, all of these properties independently contribute to the probability that this fruit is an apple and that is why it is known as 'naive'.

$$P(y|x) = \frac{P(x|y) P(y)}{P(x)}$$

likelihood → class Prior Probability
↓
 Posterior Probability

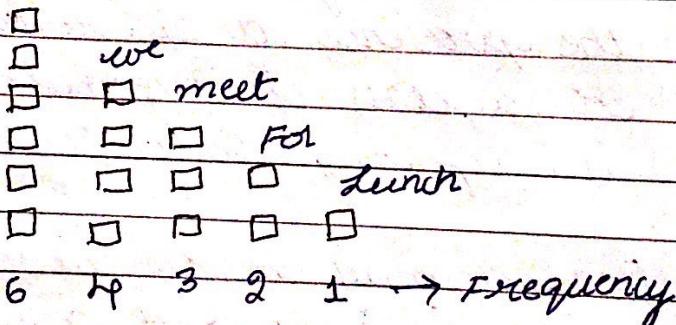
Predictor Probability

Let's take an example, imagine we received normal email from friends and family and also received spam (unwanted messages that are usually solicitous [sic] scams). Now we want to filter out the spam messages.

To solve this let's make a histogram of all the words that occur in normal emails from friends and family.

Message from family and friends : 10 mails

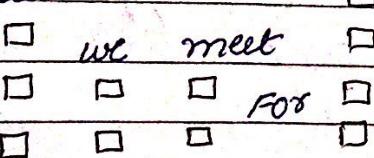
can



spam Mail : 5 mails

can

lunch



contingency / likelihood table: These are constructed by listing all the levels of one variable as rows in a table and the levels of other variables as columns, then finding the joint or cell frequency of each cell.

	Normal Emails (N)	Spam (S)
can	6 [0.315]	3 (0.27)
we	4 [0.25]	2 (0.18)
meet	3 [0.19]	2 (0.18)
For	2 [0.13]	0 (0)
lunch	1 [0.06]	4 (0.36)

$$P(\text{can}|N) = \frac{6}{16} = 0.375 \quad P(\text{can}|S) = \frac{3}{11} = 0.27$$

$$P(\text{we}|N) = \frac{4}{16} = 0.25 \quad P(\text{we}|S) = \frac{2}{11} = 0.18$$

$$P(\text{meet}|N) = \frac{3}{16} = 0.1875 \quad P(\text{meet}|S) = \frac{2}{11} = 0.18$$

$$P(\text{For}|N) = \frac{2}{16} = 0.13 \quad P(\text{For}|S) = 0$$

$$P(\text{lunch}|N) = \frac{1}{16} = 0.06 \quad P(\text{lunch}|S) = \frac{4}{11}$$

Now let's say you got an email 'meet for lunch'. we want to decide if it is a normal mail or spam.

The guess is estimated from the training data.

$$\text{Total messages} = 10 \text{ Normal} + 5 \text{ Spam} = 15$$

$$P(N) = \frac{10}{15} = 0.66$$

Now we multiply the probability with initial guess.

$$P(N) \times P(\text{can} | N) \times P(\text{we} | N) \times P(\text{meet} | N) \\ \times P(\text{for} | N) \times P(\text{lunch} | N)$$

$$= 0.000092$$

Likewise do for spam and we get the value as 0

since $0.000092 > 0$ it goes into the normal mail.

Let's take another example, 'FOR lunch Lunch Lunch'. Now calculate like above we get

spam mail = 0. but the mail should be classified as spam mail. The work around is add 1 for each word. So if we recompute [$P(FOR|S) = \frac{1}{11+4} = \frac{1}{15}$]

$$\begin{aligned}\text{Normal mail} &= 0.66 \times 0.14 \times (0.09)^3 \\ &= 0.024\end{aligned}$$

$$\begin{aligned}\text{Spam mail} &= 0.50 \times 0.06 \times (0.33)^3 \\ &= 0.029\end{aligned}$$

So, with this work around we are able to classify the mails correctly.

Naive Bayes treats language like it just a bag full of words. It ignores all the rules because keeping track of every single reasonable phrase in a language would be impossible. A Naive Bayes ignores relationship it has high bias but it works well in practice so, it has low variance.

Different variants in Naive Bayes:

(i) Bernoulli Naive Bayes:

It is useful if your feature vectors are binary. One application would be text classification where 1 is 'word occurs' and 'word doesn't occur' in document respectively.

Bernoulli distribution

$$P(\text{True}) = p$$

$$P(\text{False}) = q$$

Let's say x is a random variable

$$x=1 \text{ [True]}$$

$$x=0 \text{ [False]}$$

$$P(x=x) = p^x (1-p)^{1-x}$$

$$P(x) \begin{cases} p & \text{if } x=1 \\ q & \text{if } x=0 \end{cases}$$

(ii) Multinomial dist Naive Bayes: It can be used to find the occurrence of word i.e., the frequency in a document i.e., discrete count.

Multinomial distribution:

$$P(x_1=x_1, \dots, x_n=x_n)$$

$$= \frac{n!}{x_1! \dots x_n!} p_1^{x_1} \dots p_n^{x_n}$$

(iii) Gaussian NB: If the data is continuous in the features then it is useful.

PDF:

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

$$-\infty < x < \infty$$

