

## Naive Bayes Classifier

Bayes Theorem:

A, B events with prob.  $P(A)$  &  $P(B)$

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

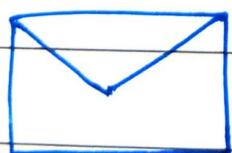
Posterior Probability

Likelihood : conditional prob of B given A

(Prob of B)  
(Normalisation factor)

Prior Prob of A

Baye's Theorem for Classification :



Spam (0)

Not Spam (1)

$x = \text{Mail}$

$y = 1 \ (\text{Not spam})$

$$\begin{aligned} P(Y=1|x) &= \frac{P(X \cap Y=1)}{P(x)} \\ &= \frac{P(Y=1) * P(x|Y=1)}{P(x)} \end{aligned}$$

// Posterior Prob that your message belongs to class 1.

$$P(Y=0|x) = \frac{P(x|Y=0) P(Y=0)}{P(x)}$$

Final prediction =  $\operatorname{argmax}(P(y_i|x))$   
 i.e. class  $y_i$  for which posterior prob is maximised.

Say you have 100 emails (60 spam, 40 not spam).

$$P(Y=0) = 0.6 \quad \& \quad P(Y=1) = 0.4$$

Note:  $P(x)$  is same for both & we are taking max, hence we can discard it.

Meaning :  $P(Y=1/X) \propto P(X/Y=1) P(Y=1)$

OR

$P(Y=c/X)$   $\propto$  likelihood  $\times$  prior

Say you have :

$x_1$  = [ "get unlimited discount, 95% off" ]

$x_2$  = [ "I am in meeting" ]

So, each word acts as feature.

$P(X/Y=0)$  = In how many spam mails these words occur.

= 0.8 (say) (means 80%).

of the time we see these words ~~not~~ in spam)

Hence,  $P(Y=0/X) = P(Y=0) P(X/Y=0)$

$$= 0.6 \times 0.8$$

$$= 0.48$$

$$P(Y=1/X) = 0.2 \times 0.4 = 0.08$$

Hence, message is spam.

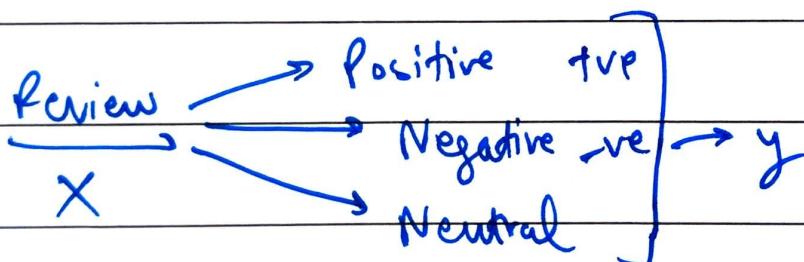
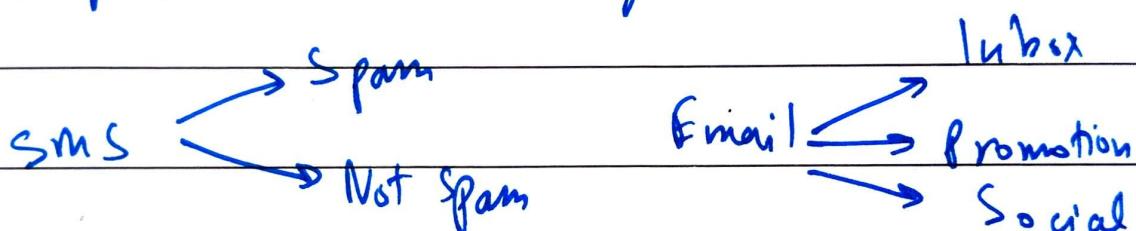
$$P(\text{spam} / \text{"lottery"}) = \frac{P(\text{"lottery"} / \text{spam}) \cdot P(\text{spam})}{P(\text{"lottery"})}$$

↑ likelihood      ↑ prior  
↓

$$P(\text{"lottery"}) = P(\text{"lottery"} / \text{spam}) + P(\text{"lottery"} / \text{Not spam})$$

Naive Bayes Classifier :

Example : Text Classification.



You will have Discrete features. You can convert text into features using bag of words model.

Vector: like ↑

0	0	1	0	1	0	0
---	---	---	---	---	---	---

$x = x_1, x_2, \dots, x_{|v|}$

$$P(\text{+ve} / x) \quad P(y/x) = \frac{P(y) \cdot P(x/y)}{P(x)}$$

$P(x|y) \rightarrow$  Compute from training data.

$P(y) \rightarrow$  Prior  $\rightarrow$  what prob of total texts are positive, etc.

Posterior

$$P(Y=1|x) = \frac{P(Y=1) P(x|Y=1)}{P(x)}$$

$$P(Y=0|x) = \frac{P(Y=0) P(x|Y=0)}{P(x)}$$

$\text{argmax } P(Y_i|x) =$  hypothesis.

i.e. that value of  $i$  for which posterior probability is maximized.

$P(Y=1) = \frac{\text{Count all tve Reviews}}{\text{Total Review}}$

$$= \sum_{i=1}^m \frac{\mathbb{1}[y^{(i)}=1]}{m} // \text{Add 1 every time } y=1$$

$$X = \langle n_1 \ n_2 \ n_3 \ \dots \ n_{uv} \rangle$$

To compute:

$$P(n_1, n_2, \dots, n_{|w|} / y=1) =$$

$$P(n_1 / y=1) \cdot P(n_2 / y=1, n_1) \cdot P(n_3 / y=1, n_1, n_2)$$

$$\dots P(n_{|w|} / y=1, n_1, n_2, \dots, n_{|w|-1})$$

-  $n_1$  ... good ...

$n_2$  - awesome ...

$n_3$  - like ...

$P(\text{awesome/tre, good})$



already seen  
good

This is becoming complex.

$$P(n_i / y, n_1, \dots, n_{i-1})$$

~~Prob.~~: Prob. that word  $n_i$  is coming given you know the class,  $n_i$  is only conditioned on  $y$  & does not depend on if you have seen other words before.

$$= P(n_i / y)$$

\* Used because it works pretty well in practice.

\* This is Naive Baye's Assumption.

$$\begin{aligned}
 P(x|y=1) &= P(n_1|y=1) P(n_2|y=1) \dots \\
 &\quad \dots P(n_m|y=1) \\
 &= \prod P(n_i|y=1)
 \end{aligned}$$

$$P(y=1|x) = \frac{\prod P(n_i|y=1) P(y=1)}{P(x)}$$

$$P(x) = P(y=0) P(x|y=0) + P(y=1) P(x|y=1)$$

### Mushroom Classification:

Given mushrooms with say 3 features :

$$X = \begin{bmatrix} n_1 & n_2 & n_3 \\ - & - & - \\ - & - & - \\ - & - & - \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

Naive Baye's : Compute Posterior prob. of mushroom of each type

$$P(y=1|x), P(y=2|x) \text{ & } P(y=3|x)$$

Given some features  $x$  of Test point  $(n, y)$

single example.

Now, find class having max posterior prob

$$P(y=1/x) = \frac{P(x/y=1) P(y=1)}{P(x)}$$

$$x = \langle n_1, n_2, n_3 \rangle$$

Naive Bayes Assumption:

$$P(x/y=1) = P(n_1/y=1) + P(n_2/y=1) + P(n_3/y=1)$$

$$P(y=1/x) \propto \prod_{i=1}^n P(n_i/y=1)$$

General :

$$P(y=c/x) \propto \prod_{i=1}^n P(n_i/y=c) P(y=c)$$

(n = no. of features)

↓      ↓      ↓

This is a single example training / testing      likelihood      Prior

$$= \frac{\prod_{i=1}^n P(n_i/y=c) P(y=c)}{\sum_{c=1}^k \left( \prod_{i=1}^n P(n_i/y=c) P(y=c) \right)}$$

$(k = \text{no. of classes})$

$\sum_{c=1}^k \left( \prod_{i=1}^n P(n_i/y=c) P(y=c) \right)$

Remains same so ignore

Naive Bayes for text:

$$P(x_j | y_i = c) = \frac{\text{count}(x_j, y_i = c)}{\sum_{\text{words } w} \text{count}(w, y_i = c)}$$

(j<sup>th</sup> feature) (i<sup>th</sup> class)      sum of words count(w, y<sub>i</sub> = c)  
 = count of words  $x_j$  that occurred in class c

\* ~~def~~: for every word in vocab, no. of times it appears in  $y = c$ , whole sum.

$$P(y = c | n) \propto \prod P(x_i | y = c) P(y = c)$$

↓ [multiply of all features  
of given n]

But there is a problem with this:

$$P(y = \text{+ve} | n) = \prod P(x_i | y = \text{+ve}) P(y = \text{+ve})$$

(Review)

Now suppose training data had words

like:

— good —  
 — happy —  
 — awesome —  
 — . . . —  
 — liked ...

Training

I was  
 overjoyed  
 after movie.

Testing data

We know we will treat each word as a feature. Now,

$$P(\text{overjoyed} \mid Y = \text{tre}) = 0$$

because "overjoyed" is not in training set.

\* If you have not seen something in your finite training set, it does not mean that its probability is 0.

Solution:

Laplace smoothing: Say that the particular word appeared at least 1 time by adding 1.

$$P(\eta_j \mid y=c) = \frac{\text{count}(\eta_j, y=c) + 1}{\sum_{w \in V} (\text{count}(w, y=c) + 1)}$$

$$= \frac{\text{count}(\eta_j, y=c) + 1}{\sum_{w \in V} (\text{count}(w, y=c)) + 1 \uparrow}$$

$\uparrow$   
(vocab size)

## Multivariate Bernoulli Naive Bayes

- \* We assume each feature will take value 0 or 1.
- \* We will not define freq like in vocab, we will define if feature was present or not.

$D_1 = \text{"like I like to swim"} \rightarrow [1 1 0 1]$

$D_2 = \text{"I like looking"} \rightarrow [1 1 1 0 0]$

Vocab = [1, like, looking, swim, to] some size

Say, for  $D_1 \Rightarrow y = \text{"Sports"}$

$D_2 = y = \text{"Code"}$

$$\begin{aligned} P(y=c) &= \prod p(n_i | y=c) (1 - p(n_i | y=c)) \\ &= (0.2)(1) * (0.3)(1) * (1)(0.9) \end{aligned}$$

$\nwarrow$  P of having "like" in  $y=c$

Here, b values are  $D_i$ ; i.e.

[1 1 0 1 1]

Prob of  
not taking  
the word  
coding in  $y=c$ .  
 $= (1 - 0.1)$

So,  $p(n|y=c)$  is probability of this sentence  $n$  being in class  $y=c$ .

$$p(y/n) = \pi p(n; | y=c) p(y=c)$$

$$p(n; | y=c) = d_{y,n; y=c} \quad (\text{frequency})$$

= Count the no. of documents that contain the feature  $n$ ; and they belong to class  $c$ .

= Count of documents having class  $c$  & containing the feature  $n$ .

Total no. of documents in class  $c$ .

$$= \frac{\text{"+1}}{\text{"+2}} \quad // \text{ Laplace smoothing}$$

every word has 2 choices, OR

because it is a bernoulli model.

Now, if  $N=D$ , then we get  $\frac{1}{2}$

We added 2 ex. one has the word & other does not.

allowed as we are taking  $P * (1-P)$ .

In numerator, we add 1 example of having that word.

i.e. if does not exist, then its  $(1-P)$

is also taken.

### Multinomial Event model Naive Bayes

- \* we do not take binary data. (for text data)
- \* we may take frequency of words / term.  
Frequency  $f(t, d) \Rightarrow$  No. of times a term appeared in a particular document.
- \* Normalised term freq =  $\frac{f(t, d)}{n_d}$   $\rightarrow$  (total no. of documents)

$$P(Y/x) = \prod P(n_i | y) P(y)$$

(can also use  $(tf, idf)$  form)

$$P(n_i | y=c) = \frac{\sum f(n_i, d \in c) + 1}{\sum n_{d \in c} + 1}$$

= sum of term freq of word  $n_i$  where

document  $d$  belongs to class  $c$

+ 1

no. of times word  
 $n_i$  is appearing  
in all documents  
belonging to class  $c$

Total no. of words in  
documents that  
belong to class  $c$ .

+ 1

$\alpha = 1 \Rightarrow$  Laplace smoothing -

$\alpha \Rightarrow$  hyper factor which changes amount of smoothing.

Difference between the two:

\* Multivariate Bernoulli Event model:

d = "offer for you"

vector = [  $\underbrace{0 \dots 0}_{\text{vocab size} = |\mathcal{V}|}$  |  $\underbrace{0 \dots 0}_{\text{offer}}$  |  $\underbrace{0 \dots 0}_{\text{for}}$  |  $\underbrace{0 \dots 0}_{\text{you}}$  ]

$P(n_i | y = \text{offer})$

$1 - P(n_i = 0 | y = \text{offer})$

(l means is included in document)

$$P(n_i = 1 | y = \text{spam}) = P(n_i | y = \text{spam})$$

$$P(n_i = 0 | y = \text{spam}) = 1 - P(n_i | y = \text{spam})$$

$$L(\theta) = \prod_{i=1}^{|\mathcal{V}|} P(n_i | y = \text{spam})^{n_i} (1 - P(n_i | y = \text{spam}))^{1 - n_i} P(\text{spam})$$

\* Multinomial Event model naive Bayes:

d = "offer for you" =  $\langle n_1, n_2, n_3 \rangle$

n for d = 3

$$\prod_{i=1}^n P(n_i | y) P(y)$$

(comes from multinomial distribution)

\* Here we do not take entire vector of length  $|\mathcal{V}|$

\* n can differ for every document.

\* Performs better generally when vocab size is large.

\* Also depends on features, you need to experiment.

### Gaussian Naive Bayes:

- \* Multivariate & multinomial are used for discrete features (as boolean, or count)
- \* Gaussian NB: continuous feature
  - In ~~that~~ this, data has gaussian distribution.

Prob. distribution formula:

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

continuous as in:

$x$  can be 1.1, 2.345, 8.99, etc

\* As you go away from mean, the prob density decreases.

\*  $\sigma$ : How quickly density will fall, or, the width of bell curve.

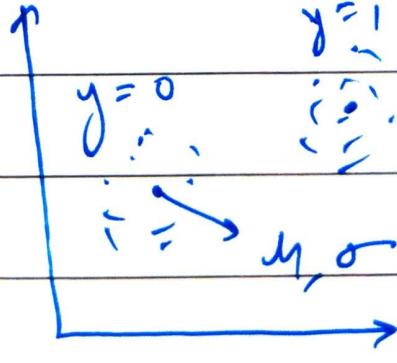
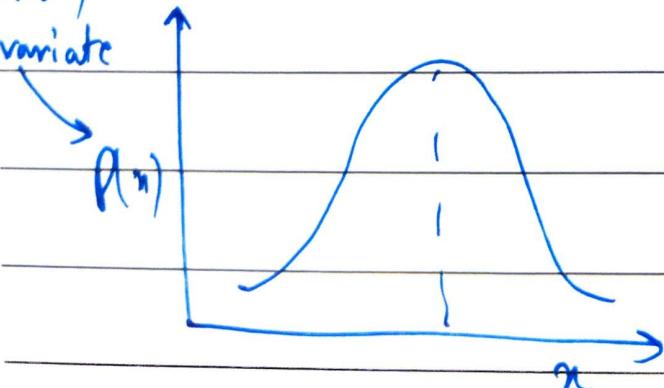
Conditional prob:

$$P(y=1|x) = \{ \text{Normal Distribution } (\mu, \sigma^2) \}$$

$\rightarrow$  some N D

for

Univariate/  
multivariate



Confusion Matrix:

		Predicted		Total (n=165)
		N	P	N → -ve
Actual	N	TN = 50	FP = 10	P → +ve
	P	FN = 5	TP = 100	

\* Accuracy =  $\frac{TP + TN}{n \text{ (total examples)}}$

\* Precision =  $\frac{TP}{TP + FP}$

= % of examples that ~~were~~ are really +ve from our prediction of +ve examples.

\* Recall =  $\frac{TP}{TP + FN}$

= % of examples that we predicted as +ve from all examples that were actually +ve.

\* Goal is to achieve high value of precision & recall but in reality if precision ↑ then recall ↓ & vice-versa.

\* To overcome this, we use:

F-measure, value between (0, 1)

Tells us the performance of model -

= ~~Harmonic mean~~ of Precision & Recall

$$\text{F-measure} = \frac{2 * \text{TP}}{2 * \text{TP} + \text{FP} + \text{FN}}$$

NB works better when you have:

- \* Small training data set
- \* Text classification
- \* Small features (dimensions)
- \* NB model size is low, won't overfit as it cannot represent complex behaviour.
- \* Use when data changes quickly
- \* Conditional independence of ~~the~~ features.