

CH - Naive Bayes.

Naive Bayes is a classification algorithm which is based on probability. Before understanding Naive Bayes, we need to be familiar with some terms.

i) Conditional probability :- $P(A|B)$:- Probability that $A = a$ given that B is b .

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad \left\{ \begin{array}{l} P(B) \neq 0 \\ \end{array} \right.$$

ii) Independent events :- Two events are said to be independent if $P(A|B) = P(A)$ similarly $P(B|A) = P(B)$

iii) Mutually exclusive events :- If $P(A|B) = P(B|A) = 0$ then $A \& B$ are said to be mutually exclusive.

$$P(A|B) = P(B|A) = 0$$

~~$P(A \cap B)$~~ - ~~$P(A)$~~ \Rightarrow ~~$P(A)$~~ = ~~$P(B)$~~

likelyhood

iv) Bayes' Theorem :- $P(A|B) = \frac{P(B|A) P(A)}{P(B)}$, $P(B) \neq 0$

Posterior Prior
 evidence

Proof : L.H.S. $P(A|B) = \frac{P(A \cap B)}{P(B)}$

$$= \frac{P(B \cap A)}{P(B)} \quad \text{--- (i)}$$

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

$$P(A) + P(B|A) = P(B \cap A) \quad \text{--- (ii)}$$

put value of $P(B \cap A)$ in (i)

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} = \text{R.H.S}$$

Now let's understand Naive Bayes algorithm.

Naive Bayes is a conditional probability model: given a problem instance to be classified, represented as vector $\mathbf{x} = (x_1, x_2, \dots, x_n)$ representing n features (independent variables), it assigns to this instance probabilities

$$P(C_k | x_1, \dots, x_n)$$

for each of K possible outcomes or classes C_k .

The problem with above formula is that if the number of features n is large or if a feature can take on a large no. of values, then basing such a model on probability tables is infeasible. We therefore reformulate the model to make it more tractable.

Using Bayes theorem,

$$P(C_k | \mathbf{x}) = \frac{P(\mathbf{x}|C_k) P(C_k)}{P(\mathbf{x})}$$

Numerator: $P(\mathbf{x}|C_k) P(C_k) = P(\mathbf{x} \cap C_k)$

Suppose we have 4 classes: we calculate

$$\begin{aligned} P(C_1 | \mathbf{x}) &= P(\mathbf{x}|C_1) * P(C_1) / P(\mathbf{x}) \\ P(C_2 | \mathbf{x}) &= P(\mathbf{x}|C_2) * P(C_2) / P(\mathbf{x}) \\ P(C_3 | \mathbf{x}) &= P(\mathbf{x}|C_3) * P(C_3) / P(\mathbf{x}) \\ P(C_4 | \mathbf{x}) &= P(\mathbf{x}|C_4) * P(C_4) / P(\mathbf{x}) \end{aligned} \quad \left. \begin{array}{l} \text{max class will choose} \\ \text{as class label of } \mathbf{x} \end{array} \right\}$$

This is common and output is not dependent on the.

So, let's focus only on numerator for now.

$$\therefore P(C_k | \mathbf{x}) \propto P(C_k, \mathbf{x}), \quad \{P(C_k \cap \mathbf{x})\}$$

This is called a joint probability model.

$$P(C_k, x_1, x_2, x_3, \dots, x_n) \Rightarrow P(C_k, x_1, x_2, x_3, \dots, x_n)$$

$$P(C_k, x_1, x_2, x_3, \dots, x_n) = P(x_1, x_2, x_3, \dots, x_n, C_k)$$

$$= P(x_1 | x_2, x_3, \dots, x_n, C_k) P(x_2, x_3, \dots, x_n, C_k)$$

$$= P(x_1 | x_2, x_3, \dots, x_n, C_k) P(x_2 | x_3, \dots, x_n, C_k) P(x_3, \dots, x_n, C_k)$$

$$= P(x_1 | x_2, x_3, \dots, x_n, C_k) \cdot \dots \cdot P(x_n | C_k) * P(C_k)$$

Applying Probability chain rule.

Now the "naive" conditional independence assumptions come into play: assume that each feature x_i is conditionally independent of every other feature x_j for $j \neq i$, given the category C . This means that,

$$p(x_i | x_{i+1}, \dots, x_n, C) = p(x_i | C)$$

Thus the joint model can be expressed as,

$$\begin{aligned} p(C_k | x_1, \dots, x_n) &\propto p(C_k, x_1, \dots, x_n) \\ &\propto p(C_k) p(x_1 | C_k) p(x_2 | C_k) \dots \\ &\propto p(C_k) \prod_{i=1}^n p(x_i | C_k) \end{aligned}$$

That means under above assumptions, the conditional distribution over the class variable C is:

$$p(C_k | x_1, \dots, x_n) = \frac{p(C_k) \cdot \prod_{i=1}^n p(x_i | C_k)}{p(X)}$$

→ Constructing a classifier from the above model.

Naive Baye's classifier combines this model with a decision rule. One common rule is to pick the hypothesis that is most probable. This is known as maximum a posteriori or MAP decision rule.

$$\hat{y} = \max_{k \in \{1, \dots, K\}} p(C_k) \cdot \prod_{i=1}^n p(x_i | C_k)$$

Predicted class label.

lets take an example:-

	outlook	Temperature	Humidity	Wind	Response: (Play = YES/no)
Day 1	Sunny	Hot	High	weak	No
Day 2	Sunny	Hot	High	Strong	No
Day 3	Overcast	Hot	High	Weak	Yes
Day 4	Rain	Mild	High	weak	Yes
Day 5	Rain	Cool	Normal	weak	Yes
Day 6	Rain	Cool	Normal	Strong	No
Day 7	Overcast	Cool	Normal	Strong	Yes
Day 8	Sunny	Mild	High	Weak	No
Day 9	Sunny	Cool	Normal	Weak	Yes
Day 10	Rain	Mild	Normal	Weak	Yes
Day 11	Sunny	Mild	Normal	Strong	Yes
Day 12	Overcast	Mild	High	Strong	Yes
Day 13	Overcast	Hot	Normal	Weak	Yes
Day 14	Rain	Mild	High	Strong	No

The learning phase :- In learning phase , we compute the table of likelihoods from the training data : They are:

$P(\text{outlook} = o | \text{class play} = b)$ where $o \in \{\text{sunny, overcast, rainy}\}$, $b \in \{\text{yes, no}\}$

$P(\text{Temp} = t | \text{class play} = b)$ where $t \in \{\text{Hot, Mild, Cool}\}$

$P(\text{Humidity} = h | \text{class play} = b)$ where $h \in \{\text{High, Normal}\}$

$P(\text{Wind} = w | \text{class play} = b)$ where $w \in \{\text{Weak, Strong}\}$

	Yes	No	$P(\text{Out})$	$P(\text{No})$
$\cdot P(o c)$	Sunny	2	3	$2/9$
	Overcast	4	0	$4/9$
	Rain	3	2	$3/9$

	Y	N	$P(Y o)$	$P(T N)$
Hot	2	2	$2/9$	$2/5$
Mild	4	2	$4/9$	$2/5$
Cold	3	1	$3/9$	$1/5$

$P(H/C)$

	Y	No	$P(H/Y)$	$P(H/N)$
High	3	4	3/9	4/5
Normal	6	1	6/9	1/5

$P(W/C)$

	Y	N	$P(S/Y)$	$P(S/N)$	$P(W/Y)$	$P(W/N)$
Strong	3	3	3/9	3/5		
Weak	6	2	6/9	2/5		

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

$$P(\text{clou} = \text{No}) = 5/14.$$

→ Classification phase: Let's say we get a new instance of the weather condition, $x' = (\text{sunny}, \text{cool}, \text{high}, \text{strong})$ we have to clarify it.

$$\hat{y} = \max \left(\left\{ P(\text{Yes}) * P(\text{sunny}/\text{Y}) * P(\text{cool}/\text{Y}) * P(\text{high}/\text{Y}) * P(\text{strong}/\text{Y}) \right\} / \left\{ P(\text{No}) * P(\text{sunny}/\text{N}) * P(\text{cool}/\text{N}) * P(\text{high}/\text{N}) * P(\text{strong}/\text{N}) \right\} \right)$$

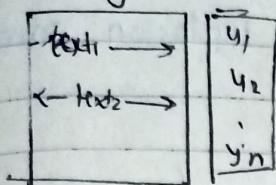
$$\begin{aligned} \hat{y} &= \max \left\{ \frac{9}{14} * \frac{2}{3} * \frac{1}{2} * \frac{4}{5} * \frac{2}{3} \right\} / \left\{ \frac{5}{14} * \frac{3}{5} * \frac{1}{3} * \frac{1}{5} * \frac{2}{3} \right\} \\ \hat{y} &= \max \left\{ 0.0053, 0.0205 \right\} \quad \rightarrow \max : \quad \hat{y} = \text{No} \end{aligned}$$

Time complexity: $O(n * d)$] → Training phase
 Space complexity: $O(d * c)$]

Time comp: $O(d * c) \rightarrow$ No. of features
 Space: \rightarrow No. of classes] → Test time

* Naive Bayes is much more memory efficient at run time.

* Naive Bayes on text data:-



text $\rightarrow \{ p_{\text{processing}} \} \rightarrow \{ w_1, w_2, \dots, w_d \}$

$$P(y=1 | \text{text}) \propto P(y=1 | w_1, w_2, \dots, w_d)$$

$$\propto P(y=1) * P(w_1 | y=1) * P(w_2 | y=1) * \dots * P(w_d | y=1)$$

$$\propto P(y=1) * \prod_{i=1}^d P(w_i | y=1)$$

Similarly $P(y=0 | \text{text}) \propto P(y=0) * \prod_{i=1}^d P(w_i | y=0)$

If $P(y=1 | \text{text}) > P(y=0 | \text{text}) \quad \hat{y} = 1$
else $\hat{y} = 0$

→ Laplace Smoothing :- At the end of training a Bayes classifier, we have.

$$\begin{aligned} & P(y=1), P(y=0) \\ & P(w_1 | y=1), P(w_2 | y=1), \dots, P(w_d | y=1) \\ & P(w_1 | y=0), P(w_2 | y=0), \dots, P(w_d | y=0) \end{aligned} \quad \left. \begin{array}{l} \text{likelihood} \\ \text{prior} \end{array} \right\}$$

Now at Test stage :- $\text{text}_q = (w_1, w_2, w_3, w')$

There are odd words

So, what we will do, if we get new word in test data.

$$\begin{aligned} P(y=1 | \text{text}_q) &= P(y=1 | w_1, w_2, w_3, w') \\ &= P(y=1) * P(w_1 | y=1) * P(w_2 | y=1) * P(w_3 | y=1) \\ &\quad * P(w' | y=1) \end{aligned}$$

How to find this

$$\begin{aligned} P(w' | y=1) &= \frac{P(w', y=1)}{P(y=1)} : \frac{\text{No. of data points s.t } w' \text{ occurs } \& y=1}{\text{No. of points where } y=1} \\ &= \frac{0}{\text{No. of points where } y=1} \rightarrow 0, \text{ because } w' \text{ occurs no data point which can be } w' \text{ because } w' \text{ occurs now} \\ &= 0 \end{aligned}$$

But we cannot put 0 - because we will get $P(y=1 | \text{text}_q) = P(y=0 | \text{text}_q) = 0$

To remove this we ~~you~~ use

Laplace smoothing (it is Laplace not Laplacian smoothing)

it says: $P(w^i | y=1) = \frac{\alpha}{n + \alpha k}$

$\hookrightarrow k = \text{No. of distinct values } w^i \text{ can take.}$

(In our case, w^i can be present or not present)

hence w^i can have 2 values $\therefore k = 2$)

α can be any value, typically $\alpha = 1$

$$\therefore P(w^i | y=1) = \frac{1}{n_i + 2}$$

• Come let a case where $\alpha = 10000$

$$P(w^i | y=1) = \frac{10000}{n_i + 20000} \quad (\text{let } n_i = 10000 = \frac{10000}{20000} \approx 0.5)$$

when α is large we get a balanced probability.

So, for each word whether it is present in training and

$$P(w^i | y=1) = \frac{\#\text{datapoint with } w^i \text{ and } y=1 + \alpha}{\#\text{data points } y=1 + \alpha k}$$

as $\alpha \uparrow \rightarrow$ moving my likelihood prob to conformable.

One more problem is:

$$\text{we know } P(y=1 | w_1, w_2, \dots, w_d)$$

$$= P(y=1) * P(w_1 | y=1) * P(w_2 | y=1) \dots P(w_d | y=1)$$

All these values are b/w 0 to 1 and

hence if α is very large this number becomes very small. It is called numerical stability.

In lang like python, C etc. we can have numerical underflow (i.e. we cannot represent them)

so, solution is.. take log of the probabilities.

$$p(y=1|w_1, w_2, \dots, w_d) \propto p(y=1) \prod_{i=1}^d p(w_i|y=1)$$

$$p(y=0|w_1, w_2, \dots, w_d) \propto p(y=0) \prod_{i=1}^d p(w_i|y=0)$$

Instead of comparing these probabilities, we can compare log of these probabilities.

$$\log \{ p(y=1|w_1, w_2, \dots, w_d) \} \Rightarrow \log(p(y=1)) + \log(p(w_1|y=1)) + \log(p(w_2|y=1)) \\ \cdot = \log(p(y=1)) + \sum_{i=1}^d \log(p(w_i|y=1))$$

) Bias-variance tradeoff in Naive Bayes:-

In Naive Bayes, only one parameter α i.e. α (Laplace smoothing) decides underfitting & overfitting.

Case 1: $\alpha = 0$

$$P(w_i|y=1) = \frac{\# \text{ train datapoint with } w_i \text{ occur } \& y=1}{\# \text{ train points with } y=1}$$

Let w_i occur only 2 times and there are 1000 train datapoints.

$$P(w_i|y=1) = \frac{2}{1000}$$

Hence even for a rare word we are assigning some prob. hence we are overfitting in a way. Now suppose $\alpha = 0$

$$\text{Then } P(w_i|y=1) = \frac{2+0}{1000+2*0} = \frac{2}{1000}$$

If we remove these 2 words then probability become 0 that means even for a small change model changes drastically i.e. model have high variance when α is small hence $\alpha=0$ have high variance problem i.e. overfitting

Case 2: α is very large $\approx \alpha = 10000$

$$P(w_i|y=1) = \frac{2 + 10000}{1000 + 2 * 10000} = \frac{10002}{21000} \approx 0.5$$

when α is very large, $P(w_i|y=1) \approx P(w_i|y=0) \approx 0.5$
hence it is underfitting.

hence if α is high, every word have approx prob as 0.5
hence this is a case of underfitting. When $\alpha \neq 1$
every time predicted class is majority class.
Hence we need a balanced α .

How to find the right α . similar to K-NN, finding
K we used cross-validation, same here we
use cross validation to find α .

→ Feature importance & interpretability:-

In KNN we used feature selection to find feature
importance. Now let see this in NB.

In NB, we have likelihood probability

$$\forall w_i, p(w_i | y_0 = 1)$$

$$\forall w_i, p(w_i | y_0 = 0)$$

If we sort the words on the basis of $p(w_i | y_0 = 1)$ in decrescent
order so, top words occur very often because their probability
is very high. We can say that, words which have
high value of probability are the important features
in determining that a data point belongs to corresponding
class.

Similarly we can sort on basis of $p(w_i | y_0 = 0)$.

Feature importance in Naive Bayes is simple to calculate.

Now, what about Interpretability:-

Let we have y_0 and declare that $y_0 = 1$, we can
say that $y_0 = 1$ because we have ~~we have~~ $p(w_i | y_0 = 1)$
is high. We can make statements like this i.e.
we can give reasoning for our prediction. Hence NB
is an interpretable ~~used~~ model.

→ Imbalanced data:- what happen to NB when we have imbalance data set.

$$n \begin{cases} \rightarrow n_1 & +ve \\ \rightarrow n_2 & -ve \end{cases}$$

$$n_1 >> n_2 \quad \left\{ \begin{array}{l} \text{let } n_1 = 90\% \\ n_2 = 10\% \end{array} \right.$$

$$\frac{n_1}{n} = 0.9 \quad \frac{n_2}{n} = 0.1$$

$$P(y=1 | w_1, w_2, \dots, w_d) = P(y=1) \cdot \prod_{i=1}^d P(w_i | y=1)$$

$$P(y=0 | w_1, w_2, \dots, w_d) = P(y=0) \prod_{i=1}^d P(w_i | y=0)$$

if these values
are equal then $y=1$

∴ (1) & because of class priors \rightarrow majority / dominating class has an advantage.

~~How~~ what is solution?

solⁿ: ① use the standard tech of upsampling or down sampling.

11 simply drop the class prior terms.

$$P(y=1 | w_1, w_2, \dots, w_l) = \prod_{i=1}^l P(w_i | y=1)$$

$$P(y=0 | w_1, w_2, \dots, w_d) = \prod_{i=1}^d P(w_i | y=0)$$

- ⑪ modified NB formulation to account for imbalanced data.
But they are not often used.

There is another problem with imbalanced data.

$$n \xrightarrow{n_1 + v_e \quad (90\%)} P(w_i | y=1) \Rightarrow \frac{9}{9+1} (0 \text{ to } 900)$$

$$\xrightarrow{n_2 - v_e \quad (10\%)} P(w_i | y=0) \Rightarrow \frac{1}{9+1} (0 \text{ to } 100)$$

This have more possibilities hence it tends to be larger.

And hence if we give a , then d impacts more of minority class and less impact on majority class.

Ex:- Let $P(w_i | y=0) = \frac{2}{100} = 2\%$

$$P(w_i | y=1) = \frac{18}{900} = 2\%$$

g = 10

$$\frac{2+10}{12 \times 2} = 10\%$$

-土3%

$$\frac{2+10}{920} = 3.09\%$$

→ Outlier How are outliers handled in naive bayes.
 If we get an outlier at test time, Laplace smoothing takes care of it.
 So, what happens to outliers in training data.
 $w_1, w_2, w_3, \dots, w_m$ → set of words in Dtrain

Let w_8 occurs very few times in the +ve class.
 Then we can say w_8 is outlier.
 Then soln is: if a word (w_j) occurs fewer than threshold then just ignore that word, and drop that word.

Second soln is to use Laplace smoothing in training data also.

→ Missing values

- (i) Case - 1: Text-data :- There is no case of missing data in Text-data
- (ii) Case - 2: categorical feature :- If a categorical feature is missing, then we can assume a another category NaN or null.
- (iii) Case - 3: Numerical features :- If we have numerical feature is missing we can use imputation as discussed previously.

★ NB with Numerical features:-

Till now we have seen NB on categorical feature and Binary (text) feature. Now suppose we have features with numerical data.

x_i	x_{i1}	x_{i2}	\vdots	x_{id}	y_i
					-
					+

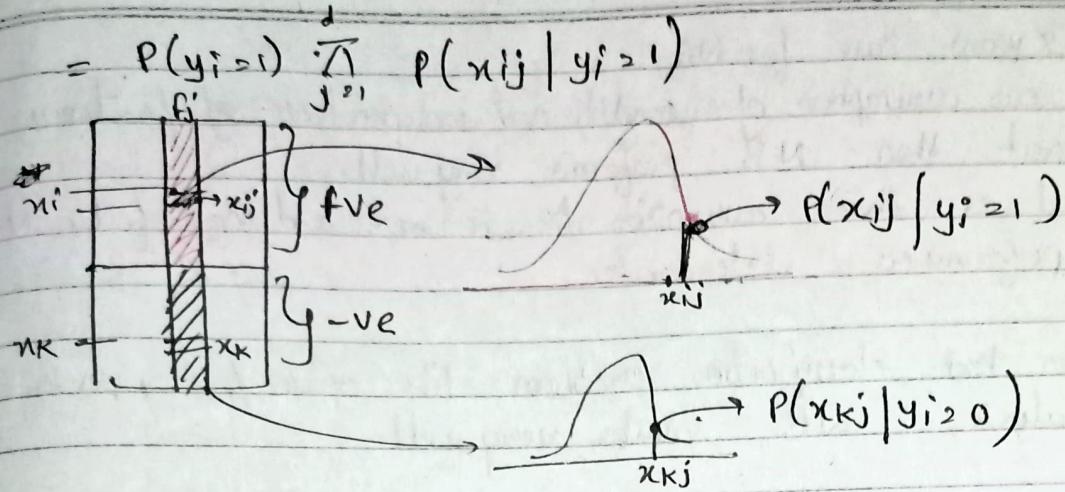
$\sum n_1$
+ve class
 $\sum n_2$
-ve class

Let f_j is Real valued.

$$P(y_i=1 | x_{i1}, x_{i2}, x_{i3}, \dots, x_{id})$$

$$\propto \frac{P(y_i=1)}{\sum_{j=1}^d P(x_{ij} | y_i=1)}$$

class prior



Hence, we can use PDF to get $P(x_{ij} | y_i = k)$
 ✓ PDF of feature j only on the data points of class K

people usually assume f_j in $D' \rightarrow$ Gaussian distn

1 " " f_j in $D' \rightarrow$ "
 $\hookrightarrow D^{-ve}$

And such a model with above 2 assumptions, is called Gaussian Naive based.

→ multiclass classification:- Naive Bayes can do multiclass classification.

→ Similarity matrix or distance matrix:-

Can NB do classification if instead of data point distance or similarity matrix is given. In general.

NB cannot use distance or similarity matrix. because NB is a probability based model where we need actual data values.

→ Large dimensionality:- NB is extensively used in Text classification which is a higher dimensional data. Hence NB can handle large dimensional data. only catch is we need to use log-probabilities.

→ Best & Worst case for NB:

- ① If our assumption of conditional independence of features is correct then NB performs very well.
But as this assumption become more and more false NB performance deteriorates.

- ② For text classification problem like spam filter, review analysis, NB works very well.

- ③ NB also perform well for categorical & Binary values.
- ④ for real values and numerical values - feature NB is not often used.
- ⑤ NB is very interpretable, we can easily get feature importance.
- ⑥ Run time complexity is low, training time complexity is low
Run time space complexity is also low.
- ⑦ ** We can easily overfit, if we don't do laplace smoothing and get right α .