

Naive Bayes

Conditional Probability:-

$P(A|B)$ = Probability of an event A given that event B has occurred.

$$= \frac{P(A \cap B)}{P(B)} ; P(B) \neq 0$$

Independent Events

A, B are said to be "independent" if :

$$\begin{cases} P(A|B) = P(A) \\ P(B|A) = P(B) \end{cases}$$

A : Getting value of 6 in die 1 throw ($D_1 = 6$).

B : Getting a value of 3 in die 2 throw ($D_2 = 3$).

Both the dice are independent of each other 😊

Mutually exclusive events:-

If $P(A|B) = P(B|A) = 0$, then A and B are said to be mutually exclusive.

A = 6 appears on die, D_1

B = 3 " " " " , D_1

Bayes Theorem

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

Proof:-

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

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

Naive Bayes Algorithm

$X = (x_1, x_2, x_3, x_4, \dots, x_n) \rightarrow n$ -dimensional vector of all column values,

It tries to calculate, $p(C_k | x_1, x_2, \dots, x_n)$

$C_k \rightarrow$ Classes or target variable

Now we know that,

$$P(C_k | X) = \frac{P(C_k) \cdot P(X | C_k)}{P(X)} \rightarrow \text{Constant for each } C_k = \frac{P(C_k, X)}{P(X)}$$

\rightarrow To determine this, only numerator is enough as denominator is constant.

$$P(X \cap C_k) = P(C_k, X) = P(C_k \text{ and } X)$$

$$P(A, B) = P(A|B) \cdot P(B)$$

$$\text{or } P(A, B, C) = P(A|B, C) P(B, C)$$

So, ~~Now~~ Using Chain rule of probability:-

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

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

$$= P(x_1 | x_2, x_3, \dots, x_n, C_k) \cdot P(x_2 | x_3, x_4, \dots, x_n, C_k) \cdot P(x_3, x_4, \dots, x_n, C_k)$$

$$\dots = P(x_1 | x_2, x_3, \dots, x_n, C_k) \cdot P(x_2 | x_3, x_4, \dots, x_n, C_k) \cdot \dots \cdot P(x_{n-1} | x_n, C_k) \cdot P(x_n | C_k) \cdot P(C_k)$$

Now, here the "Naive" assumption of "Conditional Independence" comes into picture.

$$P(A|B) = P(A) \Rightarrow \text{Independent events}$$

$$P(A|B, C) = P(A|C) \Rightarrow \text{Conditionally Independent.}$$

$$p(c_k | x_1, x_2, \dots, x_n) \propto$$

$$\text{So, } p(x_i | x_{i+1}, \dots, x_n, c_k) = p(x_i | c_k) \quad [\text{Using Naive assumption}]$$

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

↘ Multiplication

$$P(c_k | x_1, x_2, \dots, x_n) = \frac{1}{Z} p(c_k) \prod_{i=1}^n p(x_i | c_k), \quad Z = p(x)$$

Time Complexity :- $O(nd)$

Space Complexity :- $O(d * c) \Rightarrow$ To store $d * c$ likelihood probabilities

$d \Rightarrow$ no. of columns

$c \Rightarrow$ no. of classes

Naive Bayes is much better in terms of space complexity at runtime, i.e., $O(d * c)$

Naive Bayes for Text classification

It is used as a Baseline model ⁱⁿ ~~to~~ ^{its} classification problems.

Text ₁ →	0
Text ₂ →	1
	0
	1

Task = ?

$$P(y=0 | \text{text}_q) = ?$$

$$P(y=1 | \text{text}_q) = ?$$

After preprocessing of text like:- Stopwords removal, Stemming, n-grams

We get a particular sequence of words :- $\{w_1, w_2, w_3, \dots, w_n\}$

↓
Suppose we use Binary BOW on this.

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

$$= P(y=1) \times P(w_1 | y=1) \times P(w_2 | y=1) \dots P(w_n | y=1)$$

$$\text{So, } P(y=1 | \text{text}) \propto P(y=1) \times \prod_{i=1}^n P(w_i | y=1)$$

Each term is easy to calculate.

$$P(w_i | y=1) \Rightarrow \frac{\text{No. of pts containing } w_i \text{ with } y=1}{\text{No. of pts with } y=1}$$

Laplace Smoothing / Additive Smoothing

During training, we have the following values:-

$$P(y=1), P(y=0), P(w_1|y=1), P(w_1|y=0), \dots, P(w_m|y=1), P(w_m|y=0)$$

Now, suppose in the test data we have a test t containing an unseen word.

$$P(y=1 | \text{test}_t) = P(y=1 | w_1, w_2, w_3, \dots, w')$$

$$= P(y=1) \times P(w_1|y=1) \times P(w_2|y=1) \times \dots \times P(w'|y=1)$$

↓
Unknown to model

We can't just simply ignore it as in that case $P(w'|y=1)=1$ and $P(w'|y=0)=0$ will be considered which is logically incorrect.

$$P(w'|y=1) = \frac{P(w', y=1)}{P(y=1)} = \frac{\# \text{ train pts where } w' \text{ occurs and } y=1}{\# \text{ train pts where } y=1}$$

$$= \frac{0}{n_1}$$

We also can't consider it 0 as the whole expression will then become 0.

So, what happens in Laplace Smoothing is that:-

$$P(w'|y=1) = \frac{0 + \alpha}{n_1 + K\alpha} = \frac{0 + \alpha}{100 + 10\alpha}$$

\swarrow \searrow
 $n_1 = 100$ $K = 2$ (possible values w' can take)

So, let's say if $\alpha = 1 \Rightarrow P(w'|y=1) = \frac{1}{102} \neq 0$

If α is large enough, $\alpha = 10,000$.

$$P(w|y=1) = \frac{10000}{20100} \approx \frac{1}{2}$$

Now here we are considering that $P(w|y=1) = P(w|y=0) \approx \frac{1}{2}$

The w has equal chance of occurring when $y=0, 1$.

Even with small α we are getting rid of multiplying with 0.

→ for all words

$$P(w_i|y=1) = \frac{(\# \text{ data pts with } w_i \text{ and } y=1) + \alpha}{(\# \text{ data pts } y=1) + \alpha K}$$

↓
Present in my
training data

Log Probabilities for numerical stability

$0.2 \times 0.004 \times 0.003 \times \dots$ what happens in Python is that for a it starts approximating after too many decimal digits which can lead to errors.

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

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

Instead of calculate the above expression, then log is calculated

$$\log \left[P(y=1|w_1, w_2, \dots, w_d) \right] = \log(P(y=1)) + \sum_{i=1}^d \log(P(w_i|y=1))$$

Bias - Variance Tradeoff

High Bias \rightarrow Underfitting
High Variance \rightarrow Overfitting

$\left\{ \begin{array}{l} \alpha \text{ is Laplace smoothing or the hyper-} \\ \text{parameter that determines Underfit} \\ \text{and overfit.} \end{array} \right\}$

Now, let's understand this with two extreme cases:-

Case I:- $\alpha = 0$

$$P(w_i | y=1) = \frac{\# \text{ Train data pts where } w_i \text{ occurs \& } y=1}{\# \text{ Train data pts with } y=1}$$

Suppose, $n=2000$ $\left\{ \begin{array}{l} 1000 \text{ +ve} \\ 1000 \text{ -ve} \end{array} \right.$ $= \frac{2}{1000}$ \rightarrow w_i occurs in only 2 tests

So, overfitting means that small changes in ~~training~~ training data results in model change dramatically.

Since w_i occurs only in 2 tests out of 2000 tests:-

\hookrightarrow Small change in $D_{\text{train}} \rightarrow$ Removing the 2 tests containing w_i

$$\text{Probability, } P(w_i | y=1) = \frac{2}{100} \text{ to } \frac{0}{100}$$

Case II:- α is \forall large, $\alpha = 10,000$

$$P(w_i | y=1) = \frac{2 + 10,000}{1000 + 20,000} \approx \frac{1}{2}$$

\downarrow \uparrow
on line 2 values $K=2$
on 1

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

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

If these terms are same

Then majority class will decide the test data

Q How to deal with this Problem?

→ Upsampling or downsampling

↳ $n_1 \approx n_2$ $P(y=1) = P(y=0) = \frac{1}{2}$

→ drop $P(y=1)$ & $P(y=0)$
i.e., $P(y=1) = P(y=0) = 1$

→ Modified NB to account for class imbalance
↳ Not often used

They won't affect the comparison after becoming equal

One more problem with imbalanced data is:-

Minority class gets impacted more with same α as compared to majority class.

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

Let $\alpha = 10$

So, $P(w_i | y=0) = \frac{12}{120} = 10\%$

Significant change of 8%.

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

Let $\alpha = 10$

$= \frac{28}{920} = 3.04\%$

Change of 1.04%

Hence, Same α impacts both class differently.

Handling Outliers

During training time, either outliers should be ignored or they can be handled by Laplace smoothing.

Missing Values :-

In case of text data there are no missing features.

But in case of categorical data, missing values can themselves be treated as a separate category.

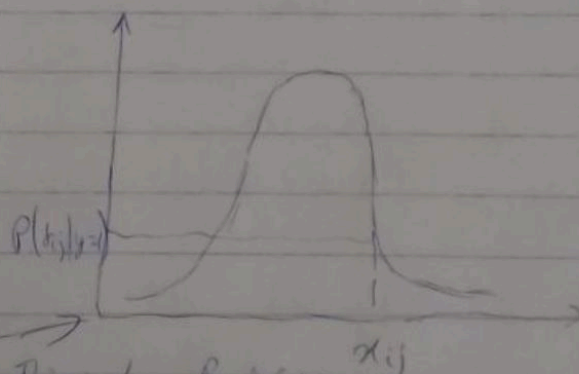
Handling Numerical Features

$$P(\underset{\substack{\uparrow \\ \text{numerical}}}{f} | y=0) \text{ or } P(f | y=1)$$

So, how do we calculate this probability.

	f_1	f_2	...	f_j	...	f_d	y
x_i	x_{i1}	x_{i2}	x_{id}	1
							1
							1
							0
							0
							0

So, for a particular feature f_j , it is considered as Gaussian distribution and its PDF is drawn.

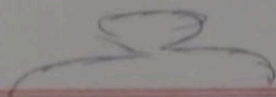


For the whole column where $y=1$ its (μ_j, σ_j) is calculated

To plot this

This is how Probability of a particular feature is calculated.

$$p(w/y=1)$$



This is called Gaussian NB, w is Gaussian.

When $w \in [0,1]$, it is Bernoulli NB.

When w can have multiple values, it is Multinomial NB.

Similarity matrix or Distance Matrix

Naive Bayes can't be used when we are provided with

Large Dimensionality :-

It can definitely handle data with large dimensionality but log-probabilities should be used to bring Numerical Stability.

Best and Worst Cases of Naive Bayes

1) ~~If~~ In case of NB, if conditional independence ~~is~~ :-

True :- Performs really well.

False :- It degrades and ~~deteriorates~~ deteriorates

In real world, some features are dependent, in those cases too Naive Bayes performs reasonably well.

2) Text classification $\begin{cases} \rightarrow \text{Email Spam} \\ \rightarrow \text{Review Polarity} \end{cases} \left\{ \begin{array}{l} \text{High Dimensionality} \\ \text{data} \end{array} \right\}$

NB can be used as a baseline model.

3) Categorical Features \rightarrow NB performs good

Real Valued features \rightarrow Only sometimes do good.

(As in real world ~~from~~ we don't get Gaussian Distribution)

4) ~~Interpretable~~ It is interpretable, feature importance can be easily found.

Runtime Complexity \rightarrow Low

Training Complexity \rightarrow Low

Run time Space Complexity \rightarrow Low