

In this chapter we will look at an algorithm called Naive Bayes.

32.1 Naive Bayes Algorithm

Probabilistic model [edit]

Abstractly, naive Bayes is a conditional probability model: given a problem instance to be classified, represented by a vector $\mathbf{x} = (x_1, \dots, x_n)$ representing some n features (independent variables), it assigns to this instance probabilities for each of K possible outcomes or classes C_k .^[7]

The problem with the above formulation is that if the number of features n is large or if a feature can take on a large number 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, the conditional probability can be decomposed as

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

In plain English, using Bayesian probability terminology, the above equation can be written as

$$\text{posterior} = \frac{\text{prior} \times \text{likelihood}}{\text{evidence}}$$

In practice, there is interest only in the numerator of that fraction, because the denominator does not depend on C and the values of the features x_i are given, so that the denominator is effectively constant. The numerator is equivalent to the joint probability model

$$p(C_k, x_1, \dots, x_n)$$

which can be rewritten as follows, using the chain rule for repeated applications of the definition of conditional probability:

Timestamp 0:57

Naive Bayes algorithm is based on the Bayes Theorem. Bayes theorem plays a very crucial rule in the world of statistics, in fact there is one entire branch called Bayesian Statistics.

Probabilistic model [edit]

Abstractly, naive Bayes is a conditional probability model: given a problem instance to be classified, represented by a vector $\mathbf{x} = (x_1, \dots, x_n)$ representing some n features (independent variables), it assigns to this instance probabilities for each of K possible outcomes or classes C_k .^[7]

The problem with the above formulation is that if the number of features n is large or if a feature can take on a large number 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, the conditional probability can be decomposed as

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

In plain English, using Bayesian probability terminology, the above equation can be written as

$$\text{posterior} = \frac{\text{prior} \times \text{likelihood}}{\text{evidence}}$$

In practice, there is interest only in the numerator of that fraction, because the denominator does not depend on C and the values of the features x_i are given, so that the denominator is effectively constant. The numerator is equivalent to the joint probability model

$$p(C_k, x_1, \dots, x_n)$$

which can be rewritten as follows, using the chain rule for repeated applications of the definition of conditional probability:

Timestamp 2:28

In this chapter we will discuss the theory behind this algorithm.

Say we have a datapoint x , we can represent it in vector form as $\langle x_1, x_2, \dots, x_n \rangle$, here x contains n features, each datapoint can belong to one of the K classes. So, we are dealing with a multi class classification problem, if $K = 2$ we have a binary classification problem.

Probabilistic model [edit]

Abstractly, naive Bayes is a conditional probability model: given a problem instance to be classified, represented by a vector $\mathbf{x} = (x_1, \dots, x_n)$ representing some 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 .^[7]

The problem with the above formulation is that if the number of features n is large or if a feature can take on a large number 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, the conditional probability can be decomposed as

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

In plain English, using Bayesian probability terminology, the above equation can be written as

$$\text{posterior} = \frac{\text{prior} \times \text{likelihood}}{\text{evidence}}$$

In practice, there is interest only in the numerator of that fraction, because the denominator does not depend on C and the values of the features x_i are given, so that the denominator is effectively constant. The numerator is equivalent to the joint probability model

$$p(C_k, x_1, \dots, x_n)$$

which can be rewritten as follows, using the chain rule for repeated applications of the definition of conditional probability:

Timestamp 4:25

Now we want to find out $p(C_k | x)$, it means that given a feature x , what is the probability that the feature x belongs to class(C) k , denoted by subscript. So, we can find this probability for all classes from 1 to k , and say that the point belongs to that class which has the maximum probability.

Now, using Bayes theorem the conditional probability can be written as

$$p(C_k | x) = (p(C_k) \times p(x | C_k)) / p(x)$$

In Bayesian terminology we have,

$$\text{Posterior} = (\text{prior} \times \text{likelihood}) / (\text{evidence})$$

Where,

$$\text{Posterior} = p(C_k | x)$$

$$\text{Prior} = p(C_k)$$

$$\text{Likelihood} = p(x | C_k)$$

$$\text{Evidence} = p(x)$$

$\underline{x \rightarrow \text{class}}$

$p(C_k | x)$

$p(C_i | x)$

$p(C_3 | x) \text{ - target}$

$p(C_k | x) = \frac{p(C_k) p(x | C_k)}{p(x)}$

$p(C_k) p(x | C_k) = p(x \cap C_k)$

The problem with the above formulation is that if the number of features n is large or if a feature can take on a large number 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, the conditional probability can be decomposed as

In plain English, using Bayesian probability terminology, the above equation can be written as

$\text{posterior} = \frac{\text{prior} \times \text{likelihood}}{\text{evidence}}$

In practice, there is interest only in the numerator of that fraction, because the denominator does not depend on C and the values of the features x_i are given, so that the denominator is effectively constant. The numerator is equivalent to the joint probability model

$$p(C_k, x_1, \dots, x_n)$$

which can be rewritten as follows, using the chain rule for repeated applications of the definition of conditional probability:

$$\begin{aligned} p(C_k, x_1, \dots, x_n) &= p(x_1, \dots, x_n, C_k) \\ &= p(x_1 | x_2, \dots, x_n, C_k)p(x_2, \dots, x_n, C_k) \\ &= p(x_1 | x_2, \dots, x_n, C_k)p(x_2 | x_3, \dots, x_n, C_k)p(x_3, \dots, x_n, C_k) \\ &= \dots \\ &= p(x_1 | x_2, \dots, x_n, C_k)p(x_2 | x_3, \dots, x_n, C_k) \dots p(x_{n-1} | x_n, C_k)p(x_n | C_k)p(C_k) \end{aligned}$$

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

Timestamp 6:18

So our main task is to calculate the probabilities for each class and determine which one is the maximum. Now, if we carefully notice the denominator, it is the same for all the classes as shown in the image, so for determining the maximum we don't need the denominator, we can only focus on the numerator.

The numerator part is $\Rightarrow p(C_k) \times p(x | C_k)$, this term calculates the probability of an event where both C_k and x occurs. So, we can simply write this as $p(x \cap C_k)$.

Timestamp 9:29

Now, $p(C_k | x) \propto p(x \cap C_k)$. The proportionality comes because we have ignored the denominator.

$p(x \cap C_k)$ can also be written as $p(C_k, x)$. So whichever class possesses the maximum value we can select that class as our prediction. Here we are calculating the probability of two events, this is often called a joint probability model.

$p(C_k, x)$ can be written as $p(C_k, x_1, x_2, \dots, x_n)$, here we have just written x in its components of n - features/dimensions.

The image shows a handwritten derivation of Bayes' Theorem. It starts with a formula for the joint probability of variables x_1, x_2, \dots, x_n and C_k , where x_1 is circled in blue and labeled 'A', and the remaining variables x_2, \dots, x_n, C_k are grouped under a brace labeled 'B'. This is equated to the product of the probability of x_1 given x_2, \dots, x_n, C_k and the probability of x_2, \dots, x_n, C_k . Below this, it is shown that $p(A, B) = p(A|B)p(B)$, which is then equated to $p(A|B) = \frac{p(A, B)}{p(B)}$. To the right, there is a note 'defn. of cond. prob.' and 'Bayes thm'.

Timestamp 12:21

Our task has boiled down to finding this $\Rightarrow p(C_k, x_1, x_2, \dots, x_n)$.

Here we will use something known as the chain rule of conditional probability, where we use the conditional probability recursively.

Conditional probability for an event A given B, is given by $\Rightarrow p(A | B) = p(A \cap B) / p(B)$. Also $p(A \cap B) = p(A | B) \times p(B)$

We will apply this simple formula recursively on this term $\Rightarrow p(C_k, x_1, x_2, \dots, x_n)$.

Now,

$$p(C_k, x_1, x_2, \dots, x_n) = p(x_1, x_2, \dots, x_n, C_k) \quad [\text{we can exchange terms because } A \cap B = B \cap A]$$

Now, we can break our elements into two components say $x_1 = A$ and $x_1, x_2, \dots, x_n, C_k = B$, and

apply conditional probability.

So, we have

$$p(x_1, x_2, \dots, x_n, C_k) = p(x_1 | x_2, \dots, x_n, C_k) \times p(x_2, \dots, x_n, C_k)$$

we want to compute

$p(C_k | \mathbf{x})$

\propto

values of the features x_i are given, so that the denominator is effectively constant. The numerator is equivalent to the joint probability model

$$p(C_k, x_1, \dots, x_n) = p(x_1, \dots, x_n, C_k)$$

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

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

$$= \dots$$

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

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

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

Thus, the joint model can be expressed as

$$p(C_k | x_1, \dots, x_n) \propto p(C_k, x_1, \dots, x_n)$$

$$\propto p(C_k) p(x_1 | C_k) p(x_2 | C_k) p(x_3 | C_k) \dots$$

$$\propto p(C_k) \prod_{i=1}^n p(x_i | C_k).$$

This means that under the above independence assumptions, the conditional distribution over the class variable C is:

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

Timestamp 15:45

We will continue applying conditional probability recursively.

This is what we get after the first expansion.

$$p(x_1, x_2, \dots, x_n, C_k) = p(x_1 | x_2, \dots, x_n, C_k) \times p(x_2, \dots, x_n, C_k)$$

Now we will keep the first part $\Rightarrow p(x_1 | x_2, \dots, x_n, C_k)$ as it is and expand the second one $\Rightarrow p(x_2, \dots, x_n, C_k)$, as shown in the image above.

Continuing like this, the final term that we get is,

$$= p(x_1 | x_2, \dots, x_n, C_k) \times p(x_2 | x_3, \dots, x_n, C_k) \dots p(x_{n-1} | x_n, C_k) \times p(x_n | C_k) \times p(C_k)$$

Now we need to solve for components of this probability terms from the data.

values of the features x_i are given, so that the denominator is effectively constant. The numerator is equivalent to the joint probability model

$$p(C_k, x_1, \dots, x_n) = p(x_1, \dots, x_n, C_k)$$

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

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

$$\dots$$

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

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

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

Thus, the joint model can be expressed as

$$p(C_k | x_1, \dots, x_n) \propto p(C_k, x_1, \dots, x_n)$$

$$\propto p(C_k) p(x_1 | C_k) p(x_2 | C_k) p(x_3 | C_k) \dots$$

$$\propto p(C_k) \prod_{i=1}^n p(x_i | C_k).$$

This means that under the above independence assumptions, the conditional distribution over the class variable C is:

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

Train
Test

Timestamp 18:00

Now, finding these individual terms is very difficult because we have to count those data points for which our condition matches exactly. For that we need lots of data, rather lots of combinations of data, which may not be always possible.

Naive Bayes

values of the features x_i are given, so that the denominator is effectively constant. The numerator is equivalent to the joint probability model

$$p(C_k, x_1, \dots, x_n)$$

which can be rewritten as follows, using the chain rule for repeated applications of the definition of conditional probability:

$$p(C_k, x_1, \dots, x_n) = p(x_1, \dots, x_n, C_k)$$

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

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

$$\dots$$

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

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

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

Thus, the joint model can be expressed as

$$p(C_k | x_1, \dots, x_n) \propto p(C_k, x_1, \dots, x_n)$$

$$\propto p(C_k) p(x_1 | C_k) p(x_2 | C_k) p(x_3 | C_k) \dots$$

$$\propto p(C_k) \prod_{i=1}^n p(x_i | C_k).$$

$p(A \cap B) = p(A)$

$\{ p(A \cap B) = p(A \cap C) \}$

$B \cap C$

This means that under the above independence assumptions, the conditional distribution over the class variable C is:

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

Timestamp 20:05

In order to overcome this problem of finding an exact match we make a naive assumption. That's why this is called Naive Bayes. The assumption is that of conditional independence of the features.

Two events are said to be independent if the occurrence of one doesn't affect the other.

$$p(A|B) = p(A)$$

Similarly conditional independence can be written as

$$p(A|B, C) = p(A|C),$$

where $p(A|B, C)$ is the probability of A given both B and C. Since the probability of A given C is the same as the probability of A given both B and C, this equality expresses that B contributes nothing to the certainty of A. In this case, A and B are said to be conditionally independent given C.

Timestamp 21:44

So our naive bayes assumption is that feature x_i is conditionally independent of every other feature x_j for $j \neq i$, given the class k .

Now this idea will be used for every term in our probability. In general,

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

Here x_i is conditionally independent of $x_{i+1}, x_{i+2}, \dots, x_n$ given C_k . This vastly simplifies our calculation.

values of the features x_i are given, so that the denominator is effectively constant. The numerator is equivalent to the joint probability model:

$$p(C_k, x_1, \dots, x_n) = p(x_1, \dots, x_n, C_k)$$

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

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

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

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

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

Thus, the joint model can be expressed as

$$\begin{aligned} p(C_k | x_1, \dots, x_n) &\propto p(C_k, x_1, \dots, x_n) \\ &\propto p(C_k) p(x_1 | C_k) p(x_2 | C_k) p(x_3 | C_k) \dots p(x_n | C_k) \\ &\propto p(C_k) \prod_{i=1}^n p(x_i | C_k). \end{aligned}$$

This means that under the above independence assumptions, the conditional distribution over the class variable C is:

Timestamp 23:00

Now the conditional independence can be applied to each of the terms. After following the steps as shown in the image we reach to our final probability, which is

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

The proportionality comes because we are ignoring the denominator. It is a product of individual terms along with the probability of the class.

$\propto p(C_k) p(x_1 | C_k) p(x_2 | C_k) p(x_3 | C_k) \dots$
 $\propto p(C_k) \prod_{i=1}^n p(x_i | C_k).$

This means that under the above independence assumptions, the conditional distribution over the class variable C is:

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

where the evidence $Z = p(x) = \sum_k p(C_k) p(x | C_k)$ is a scaling factor dependent only on x_1, \dots, x_n , that is, a constant if the values of the feature variables are known.

Constructing a classifier from the probability model [edit]

The discussion so far has derived the independent feature model, that is, the naive Bayes probability model. The naive Bayes classifier combines this model with a decision rule. One common rule is to pick the hypothesis that is most probable; this is known as the *maximum a posteriori* or MAP decision rule. The corresponding classifier, a Bayes classifier, is the function that assigns a class label $\hat{y} = C_k$ for some k as follows:

$$\checkmark \hat{y} = \underset{k \in \{1, \dots, K\}}{\operatorname{argmax}} p(C_k) \prod_{i=1}^n p(x_i | C_k).$$

Parameter estimation and event models [edit]

A class's prior may be calculated by assuming equiprobable classes (i.e., priors $\approx 1 / (\text{number of classes})$), or by calculating an

Timestamp 25:28

If we want an exact term, we can remove the proportionality and bring in the denominator as shown in the image.

Now for a classifier we will find the probability for each of the classes, and select that class which gives the maximum probability.

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

This rule is known as *maximum a posteriori* or MAP decision rule, as we are selecting the argmax of the maximum posterior.

https://en.wikipedia.org/wiki/Naive_Bayes_classifier

32.2 Toy example: Train and test Stages

<http://shatterline.com/blog/2013/09/12/not-so-naive-classification-with-the-naive-bayes-classifier/>

Naive Bayes

$y_i = \text{Yes or No}$

f_1, f_2, f_3, f_4 : Categorical feat

	Predictors				Response
	Outlook	Temperature	Humidity	Wind	Class Play-Yes
Day1	Sunny	Hot	High	Weak	No
Day2	Sunny	Hot	High	Strong	No
Day3	Overcast	Hot	High	Weak	Yes
Day4	Rain	Mild	High	Weak	Yes
Day5	Rain	Cool	Normal	Weak	Yes
Day6	Rain	Cool	Normal	Strong	No
Day7	Overcast	Cool	Normal	Strong	Yes
Day8	Sunny	Mild	High	Weak	No
Day9	Sunny	Cool	Normal	Weak	Yes
Day10	Rain	Mild	Normal	Weak	Yes
Day11	Sunny	Mild	Normal	Strong	Yes
Day12	Overcast	Mild	High	Strong	Yes
Day13	Overcast	Hot	Normal	Weak	Yes
Day14	Rain	Mild	High	Strong	No

The Learning Phase

In the learning phase, we compute the table of likelihoods (probabilities) from the training data. They are:

$P(\text{Outlook} = o | \text{Class}_i = b)$, where $o \in \{\text{Sunny, Overcast, Rainy}\}$ and $b \in \{\text{yes, no}\}$

$P(\text{Temperature} = t | \text{Class}_i = b)$, where $t \in \{\text{Hot, Mild, Cool}\}$ and $b \in \{\text{yes, no}\}$

$P(\text{Humidity} = h | \text{Class}_i = b)$, where $h \in \{\text{high, Normal}\}$ and $b \in \{\text{yes, no}\}$

$P(\text{Wind} = w | \text{Class}_i = b)$, where $w \in \{\text{Weak, Strong}\}$ and $b \in \{\text{yes, no}\}$

Timestamp 2:28

In this chapter we will look at a toy dataset and apply naive bayes algorithm on this dataset. It's a binary classification problem with target variable play, $y_i \in \{\text{Yes(Y), No(N)}\}$. Our dataset contains 4 features => Outlook(f_1), Temperature(f_2), Humidity(f_3), Wind(f_4). All these features are categorical features, naive bayes can also be used for continuous features but it is easy to visualize for categorical data. So we will stick to categorical data.

Outlook(f_1) can have 3 categories: sunny, overcast, rain. Temperature(f_2) can have 3 categories: hot, mild, cool. Humidity(f_3) can have 2 categories: high, normal. Wind(f_4) can have 2 categories: strong, weak.

Here we are trying to predict whether it is suitable to play tennis outside given various weather parameters.

Windy ∈ [Weak, Strong]

The class label is the variable, Play and takes the values yes or no.

Play ∈ {Yes, No}

We read-in training data below that has been collected over 14 days.

	Predictors				Response
	Outlook	Temperature	Humidity	Wind	Class
					Play=Yes Play=No
Day1	Sunny	Hot	High	Weak	No
Day2	Sunny	Hot	High	Strong	No
Day3	Overcast	Hot	High	Weak	Yes
Day4	Rain	Mild	High	Weak	Yes
Day5	Rain	Cool	Normal	Weak	Yes
Day6	Rain	Cool	Normal	Strong	No
Day7	Overcast	Cool	Normal	Strong	Yes
Day8	Sunny	Mild	High	Weak	No
Day9	Sunny	Cool	Normal	Weak	Yes
Day10	Rain	Mild	Normal	Weak	Yes

MAP

$p(\text{class} | \text{play} = \text{yes}) | x_q$

$\checkmark p(\text{play} = \text{no}) | x_q$

$x_q \rightarrow y_q = (\text{play} = \text{yes})$

'if'

Timestamp 4:32

So, here we need to calculate two probabilities, $p(\text{Play} = \text{Yes} | x_q)$, it says what is the probability that we should outside play outside given weather params x_q . Similarly $p(\text{Play} = \text{No} | x_q)$, what is the probability that we should not play outside given weather params?

Out of both these terms whichever term possesses the highest value we declare that to be our class label, MAP rule.

The screenshot shows a Google Doc with a table of weather data and handwritten mathematical annotations.

Table Data:

Day	Predictors				Response
	Outlook	Temperature	Humidity	Wind	
Day1	Sunny	Hot	High	Weak	No
Day2	Sunny	Hot	High	Strong	No
Day3	Overcast	Hot	High	Weak	Yes
Day4	Rain	Mild	High	Weak	Yes
Day5	Rain	Cool	Normal	Weak	Yes
Day6	Rain	Cool	Normal	Strong	No
Day7	Overcast	Cool	Normal	Strong	Yes
Day8	Sunny	Mild	High	Weak	No
Day9	Sunny	Cool	Normal	Weak	Yes
Day10	Rain	Mild	Normal	Weak	Yes
Day11	Sunny	Mild	Normal	Strong	Yes
Day12	Overcast	Mild	High	Strong	Yes
Day13	Overcast	Hot	Normal	Weak	Yes
Day14	Rain	Mild	High	Strong	No

Handwritten Annotations:

- $p(C = \text{No}) = 5/14$
- f_1, f_2, f_3, f_4 (labeled above the first four columns)
- $p(C|f_1, f_2, f_3, f_4) \propto p(f_1|C)p(f_2|C)p(f_3|C)p(f_4|C)p(C)$
- $p(C|f_1) * p(f_2|f_1, C) * p(f_3|f_1, f_2, C) * p(f_4|f_1, f_2, f_3, C)$ (with arrows pointing from each term to its respective column)
- $p(C) = 9/14$

Timestamp 8:02

As per Naive Bayes,

$$p(C | f_1, f_2, f_3, f_4) \propto p(f_1|C)p(f_2|C)p(f_3|C)p(f_4|C)p(C)$$

We have written the above term directly from Naive Bayes where f_i 's are our features and C is our class.

$p(C = \text{Yes}) = 9/14$, we have 9 yes out of 14 data points.

$p(C = \text{No}) = 5/14$, we have 5 no out of 14 data points.

This are called prior probabilities.

The screenshot shows a Google Docs slide with handwritten annotations. At the top left, there is a yellow circle with the word 'Train'. To the right, there is a large orange circle containing the handwritten text $p(f_i | c)$ at the top and $p(\text{outlook} = \text{sunny} | c = \text{Yes}) = 2/9$ below it. Below these circles is a table for 'Outlook' with columns for 'Play=Yes' and 'Play=No', and rows for 'Sunny', 'Overcast', and 'Rain'. The total count for 'Play=Yes' is circled as 9, and for 'Play=No' as 5. The probability $2/9$ is circled in the 'Play=Yes' row under 'Sunny'. Another table for 'Temperature' is partially visible below it.

$P(\text{Outlook} = o \text{Class Play} = \text{Yes/No})$	Frequency		Probability in Class	
Outlook =	Play=Yes	Play=No	Play=Yes	Play=No
→ Sunny	2	3	2/9	3/5
→ Overcast	4	0	4/9	0/5
→ Rain	3	2	3/9	2/5
	total(9)	total(5)		

$P(\text{Temperature} = t \text{Class Play} = \text{Yes/No})$	Frequency		Probability in Class	
Temperature =	Play=Yes	Play=No	Play=Yes	Play=No

Timestamp 12:33

Now we need to calculate the feature probabilities(likelihood), for calculating these we will use the train data. This phase is called the learning/training phase.

Let's calculate one of the terms => $p(\text{Outlook} = \text{sunny} | C = \text{Yes})$

Now here we need to find the probability of an event where outlook = sunny given $C = \text{yes}$. So, at first we will count all the data points which have class Yes and out of those data points we need to count those features where outlook = sunny.

So, after counting we got $p(\text{Outlook} = \text{sunny} | C = \text{Yes}) = 2/9$, Out of 9 data points which have class label = yes, 2 of them had an outlook = sunny.

Similarly $p(\text{Outlook} = \text{sunny} | C = \text{No}) = 3/5$, Out of 5 data points which have class label = No, 3 of them had an outlook = sunny.

We need to calculate the likelihood of each of the features, and store it in some data structure. As shown in the image above. Please follow the shatterline link to get a complete list of all the probabilities.

The screenshot shows a Google Docs slide with a table of weather data and handwritten annotations explaining the Naive Bayes classification process.

Table Data:

	Predictors				Response
	Outlook	Temperature	Humidity	Wind	Class
Day1	Sunny	Hot	High	Weak	No
Day2	Sunny	Hot	High	Strong	No
Day3	Overcast	Hot	High	Weak	Yes
Day4	Rain	Mild	High	Weak	Yes
Day5	Rain	Cool	Normal	Weak	Yes
Day6	Rain	Cool	Normal	Strong	No
Day7	Overcast	Cool	Normal	Strong	Yes
Day8	Sunny	Mild	High	Weak	No
Day9	Sunny	Cool	Normal	Weak	Yes
Day10	Rain	Mild	Normal	Weak	Yes
Day11	Sunny	Mild	Normal	Strong	Yes
Day12	Overcast	Mild	High	Strong	Yes
Day13	Overcast	Hot	Normal	Weak	Yes
Day14	Rain	Mild	High	Strong	No

Annotations:

- $P(F|C)$ is written above the first column of the table.
- n is circled and has an arrow pointing to the bottom of the table.
- c is circled and has an arrow pointing to the 'Class' column.
- Training ph: is written above the table.
- $\rightarrow \text{likelihood prob}$ is written next to the table.
- $\rightarrow p(\text{class})$ is written next to the table.
- $O(nd)$ is written twice in boxes with arrows pointing to the table.

Timestamp 15:19

Let's digress for a minute and calculate the complexities.

Time Complexity

Say we have n data points in our dataset and each data point has d dimension. Also say we have c classes. Then if we use a brute force algorithm our time complexity is $O(ndc)$.

This is because for calculating the individual probabilities we need to go through each feature i.e. d , n times because we have n data points. This we need to do for c classes.

If we use some optimization, time complexity can be reduced to $O(n)$. if d is small. Otherwise it is $O(nd)$.

Contents - Google Docs

Shatterline Blog - Not-so-naive-classification-with-the-naive-bayes-classifier/

Naive Bayes classifier - Wikipedia

Chakru Srikan...

shatterline.com/blog/2013/09/12/not-so-naive-classification-with-the-naive-bayes-classifier/

The screenshot shows a table of 14 data points (Day1 to Day14) with columns for Outlook, Temperature, Humidity, Wind, Class Play=Yes, and Class Play=No. Handwritten notes on the right side of the screen say "Time: - O(n)" and "Space: - O(d*c)".

	Outlook	Temperature	Humidity	Wind	Class Play=Yes	Class Play=No
Day1	Sunny	Hot	High	Weak	No	
Day2	Sunny	Hot	High	Strong	No	
Day3	Overcast	Hot	High	Weak	Yes	
Day4	Rain	Mild	High	Weak	Yes	
Day5	Rain	Cool	Normal	Weak	Yes	
Day6	Rain	Cool	Normal	Strong	No	
Day7	Overcast	Cool	Normal	Strong	Yes	
Day8	Sunny	Mild	High	Weak	No	
Day9	Sunny	Cool	Normal	Weak	Yes	
Day10	Rain	Mild	Normal	Weak	Yes	
Day11	Sunny	Mild	Normal	Strong	Yes	
Day12	Overcast	Mild	High	Strong	Yes	
Day13	Overcast	Hot	Normal	Weak	Yes	
Day14	Rain	Mild	High	Strong	No	

The Learning Phase

In the learning phase, we build tables for each feature (Outlook, Temperature, Humidity, Wind) that will have to be classified (i.e., are we going to play tennis under the conditions specified by x').

With the MAP rule, we compute the posterior probabilities. This is easily done by looking up the tables we built in the learning phase.

$P(\text{ClassPlay}=\text{Yes}|x') = [P(\text{Sunny}|\text{ClassPlay}=\text{Yes}) \times P(\text{Cool}|\text{ClassPlay}=\text{Yes}) \times P(\text{High}|\text{ClassPlay}=\text{Yes}) \times P(\text{Strong}|\text{ClassPlay}=\text{Yes})] \times P(\text{ClassPlay}=\text{Yes})$

$= 2/9 \times 3/9 \times 3/9 \times 3/9 \times 9/14 = 0.0053$

$P(\text{ClassPlay}=\text{No}|x') = [P(\text{Sunny}|\text{ClassPlay}=\text{No}) \times P(\text{Cool}|\text{ClassPlay}=\text{No}) \times P(\text{High}|\text{ClassPlay}=\text{No}) \times P(\text{Strong}|\text{ClassPlay}=\text{No})] \times P(\text{ClassPlay}=\text{No})$

$= 3/5 \times 1/5 \times 4/5 \times 3/5 \times 5/14 = 0.0205$

Since $P(\text{ClassPlay}=\text{Yes}|x')$ less than $P(\text{ClassPlay}=\text{No}|x')$, we classify the new instance x' to be Class Play=No.

Timestamp 17:00

Space Complexity is $O(dc)$, because we need to store the probabilities of each feature i.e. d, c times i.e. no of classes. Also we need to store class probabilities.

Contents - Google Docs

Shatterline Blog - Not-so-naive-classification-with-the-naive-bayes-classifier/

Naive Bayes classifier - Wikipedia

Chakru Srikan...

shatterline.com/blog/2013/09/12/not-so-naive-classification-with-the-naive-bayes-classifier/

Handwritten notes on the right side of the screen include $x' = (\text{Sunny}, \text{cool}, \text{high}, \text{Strong})$, $y_{\text{in}} = ?$, and $P(C=\text{Yes}|x')$.

Temperature=Cool, Humidity=High, Wind=Strong) that will have to be classified (i.e., are we going to play tennis under the conditions specified by x').

With the MAP rule, we compute the posterior probabilities. This is easily done by looking up the tables we built in the learning phase.

$P(\text{ClassPlay}=\text{Yes}|x') = [P(\text{Sunny}|\text{ClassPlay}=\text{Yes}) \times P(\text{Cool}|\text{ClassPlay}=\text{Yes}) \times P(\text{High}|\text{ClassPlay}=\text{Yes}) \times P(\text{Strong}|\text{ClassPlay}=\text{Yes})] \times P(\text{ClassPlay}=\text{Yes})$

$= 2/9 \times 3/9 \times 3/9 \times 3/9 \times 9/14 = 0.0053$

$P(\text{ClassPlay}=\text{No}|x') = [P(\text{Sunny}|\text{ClassPlay}=\text{No}) \times P(\text{Cool}|\text{ClassPlay}=\text{No}) \times P(\text{High}|\text{ClassPlay}=\text{No}) \times P(\text{Strong}|\text{ClassPlay}=\text{No})] \times P(\text{ClassPlay}=\text{No})$

$= 3/5 \times 1/5 \times 4/5 \times 3/5 \times 5/14 = 0.0205$

Since $P(\text{ClassPlay}=\text{Yes}|x')$ less than $P(\text{ClassPlay}=\text{No}|x')$, we classify the new instance x' to be Class Play=No.

Timestamp 18:30

Now, given a data point during test time we need to classify the data point. Say we have a data point $x_q = (\text{sunny, cool, high, strong})$, here we are given values for the 4 features. We need to calculate two posteriors, because we have two classes.

$p(\text{class} = \text{Yes} | x_q)$ and $p(\text{class} = \text{No} | x_q)$.

We can write these posteriors in terms of their features as shown in the image.

The screenshot shows a Google Doc with handwritten annotations. The text discusses classifying a new instance x' (Cool, High, Wind=Strong) using the MAP rule. It shows the calculation of posteriors for ClassPlay=Yes and ClassPlay=No, referencing a training table. A handwritten Python dictionary is shown, mapping feature combinations to probabilities.

Temperature=Cool, Humidity=High, Wind=Strong) that will have to be classified (i.e., are we going to play tennis under the conditions specified by x').

With the MAP rule, we compute the posterior probabilities. This is easily done by looking up the tables we built in the learning phase.

$$P(\text{ClassPlay}=\text{Yes}|x') = [P(\text{Sunny}|\text{ClassPlay}=\text{Yes}) \times P(\text{Cool}|\text{ClassPlay}=\text{Yes}) \times P(\text{High}|\text{ClassPlay}=\text{Yes}) \times P(\text{Strong}|\text{ClassPlay}=\text{Yes})] \times P(\text{ClassPlay}=\text{Yes})$$

$$= 2/9 \times 3/9 \times 3/9 \times 3/9 \times 9/14 = 0.0053$$

$$P(\text{ClassPlay}=\text{No}|x') = [P(\text{Sunny}|\text{ClassPlay}=\text{No}) \times P(\text{Cool}|\text{ClassPlay}=\text{No}) \times P(\text{High}|\text{ClassPlay}=\text{No}) \times P(\text{Strong}|\text{ClassPlay}=\text{No})] \times P(\text{ClassPlay}=\text{No})$$

$$= 3/5 \times 1/5 \times 4/5 \times 3/5 \times 5/14 = 0.0205$$

Since $P(\text{ClassPlay}=\text{Yes}|x')$ less than $P(\text{ClassPlay}=\text{No}|x')$, we classify the new instance x' to be No.

Dictionary - Python

Sunny	Cool	High	Wind	Probability
Yes	Yes	Yes	Yes	2/9
Yes	Yes	Yes	No	3/9
Yes	No	Yes	Yes	1/9
No	Yes	No	Yes	4/5
No	No	Yes	Yes	3/5

dict get

Timestamp 20:02

For individual terms we need to use the probabilities that we have calculated during the training phase. For example, one of the feature is $p(\text{Outlook} = \text{Sunny} | \text{class} = \text{Yes})$, so during training we can build a dictionary with feature combinations as key and the probabilities as value. We can have a key $\langle \text{Sunny}, \text{Yes} \rangle$ and its probability i.e. $2/9$ as its value. We can use python's inbuilt dictionary methods to implement these.

Temperature=Cool, Humidity=High, Wind=Strong) that will have to be classified (i.e., are we going to play tennis under the conditions specified by x').

With the MAP rule, we compute the posterior probabilities. This is easily done by looking up the tables we built in the learning phase.

$$P(\text{ClassPlay}=\text{Yes}|x') \propto [P(\text{Sunny}|\text{ClassPlay}=\text{Yes}) \times P(\text{Cool}|\text{ClassPlay}=\text{Yes}) \times P(\text{High}|\text{ClassPlay}=\text{Yes}) \times P(\text{Strong}|\text{ClassPlay}=\text{Yes})] \times P(\text{ClassPlay}=\text{Yes})$$

$$\propto 2/9 \times 3/9 \times 3/9 \times 3/9 \times 9/14 = 0.0053$$

$$P(\text{ClassPlay}=\text{No}|x') = [P(\text{Sunny}|\text{ClassPlay}=\text{No}) \times P(\text{Cool}|\text{ClassPlay}=\text{No}) \times P(\text{High}|\text{ClassPlay}=\text{No}) \times P(\text{Strong}|\text{ClassPlay}=\text{No})] \times P(\text{ClassPlay}=\text{No})$$

$$= 3/5 \times 1/5 \times 4/5 \times 3/5 \times 5/14 = 0.0205$$

Since $P(\text{ClassPlay}=\text{Yes}|x')$ less than $P(\text{Classplay}=\text{No}|x')$, we classify the new instance x' to be No.

Timestamp 21:24

After putting all the likelihood values and priors in we get,

$$p(\text{class} = \text{Yes} | x_q) \approx 0.0053$$

$$p(\text{class} = \text{No} | x_q) \approx 0.0205$$

So as per MAP rule we can say that our class instance is No.

Classification Phase

Let's say, we get a new instance of the weather condition, $x'=(\text{Outlook}=\text{Sunny}, \text{Temperature}=\text{Cool}, \text{Humidity}=\text{High}, \text{Wind}=\text{Strong})$ that will have to be classified (i.e., are we going to play tennis under the conditions specified by x').

With the MAP rule, we compute the posterior probabilities. This is easily done by looking up the tables we built in the learning phase.

$$P(\text{ClassPlay}=\text{Yes}|x') = [P(\text{Sunny}|\text{ClassPlay}=\text{Yes}) \times P(\text{Cool}|\text{ClassPlay}=\text{Yes}) \times P(\text{High}|\text{ClassPlay}=\text{Yes}) \times P(\text{Strong}|\text{ClassPlay}=\text{Yes})] \times P(\text{ClassPlay}=\text{Yes})$$

$$= 2/9 \times 3/9 \times 3/9 \times 3/9 \times 9/14 = 0.0053$$

$$P(\text{ClassPlay}=\text{No}|x') = [P(\text{Sunny}|\text{ClassPlay}=\text{No}) \times P(\text{Cool}|\text{ClassPlay}=\text{No}) \times P(\text{High}|\text{ClassPlay}=\text{No}) \times P(\text{Strong}|\text{ClassPlay}=\text{No})] \times P(\text{ClassPlay}=\text{No})$$

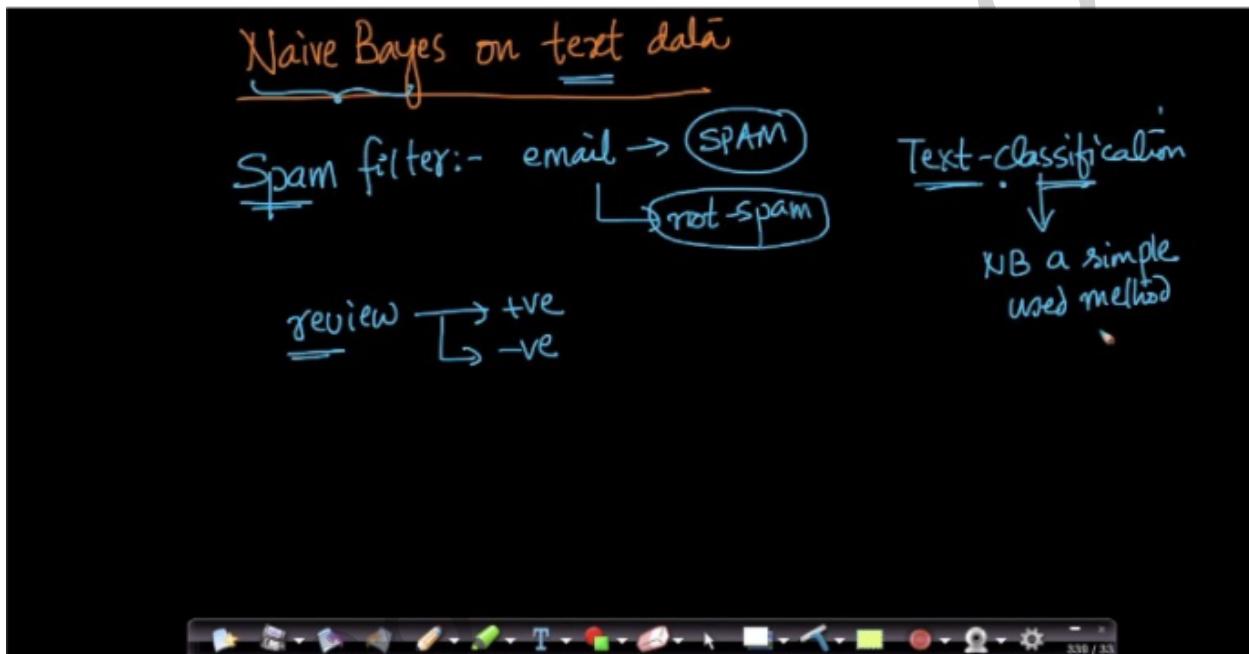
Timestamp 24:17

Now let's look at complexity during test time.

Time Complexity = $O(d*c)$, where d is the dimension of the data, c is the no of classes. This is because we need to look up d times i.e. for each feature, for each class i.e. c times. We are also considering our lookup time to be constant.

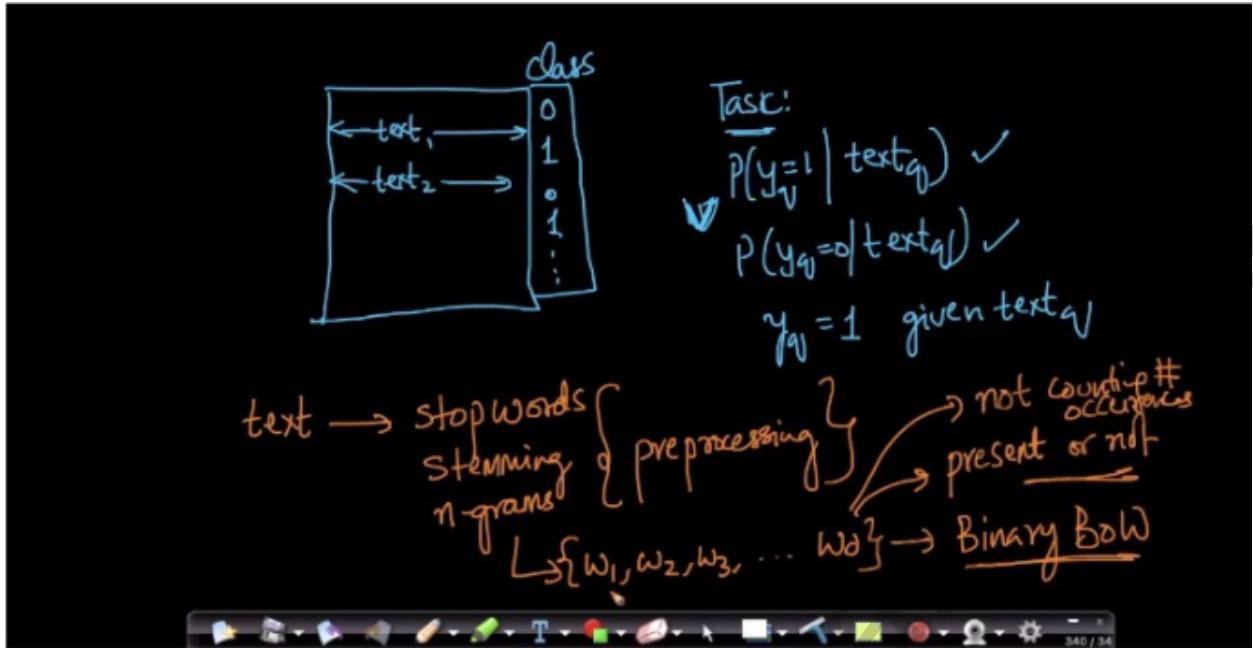
If we compare NB with KNN, we see that NB is extremely space efficient at runtime because in KNN we need all the data points in memory, whereas NB requires only the look up tables to be in the memory.

32.3 Naive Bayes on Text Data



Timestamp 1:31

NB is extremely useful when it comes to text data. One of the earliest spam filters was built using NB. A spam filter classifies a mail into $\{\text{Spam, Not Spam}\}$. NB is very efficient in classifying text data into various classes.



Timestamp 4:36

Suppose we are given a dataset where we have textual data; each text belonging to a certain class. Let's keep it to binary classification; it can be easily extended to multi-class classification problems. So our class labels are $\epsilon\{0, 1\}$. As shown in the image above.

Now we want to calculate,

$p(y_q = 1 | text_q)$ and $p(y_q = 0 | text_q)$ i.e. probability that the class is 0 or 1, given text $text_q$.

For handling text data we need to preprocess it, we can remove stopwords, do stemming, calculate n-grams.etc. All these steps are application dependent and not concrete.

After doing the preprocessing we need to convert this text data into numerical vectors. In NB we typically use binary bag of words(BOW). Binary BOW contains cells for each unique word in the corpus. Each cell can have a value of 0 or 1, representing whether that particular word is present in the current text or not. We will consider only those cells where its value is 1, i.e. the words present in the text.

S

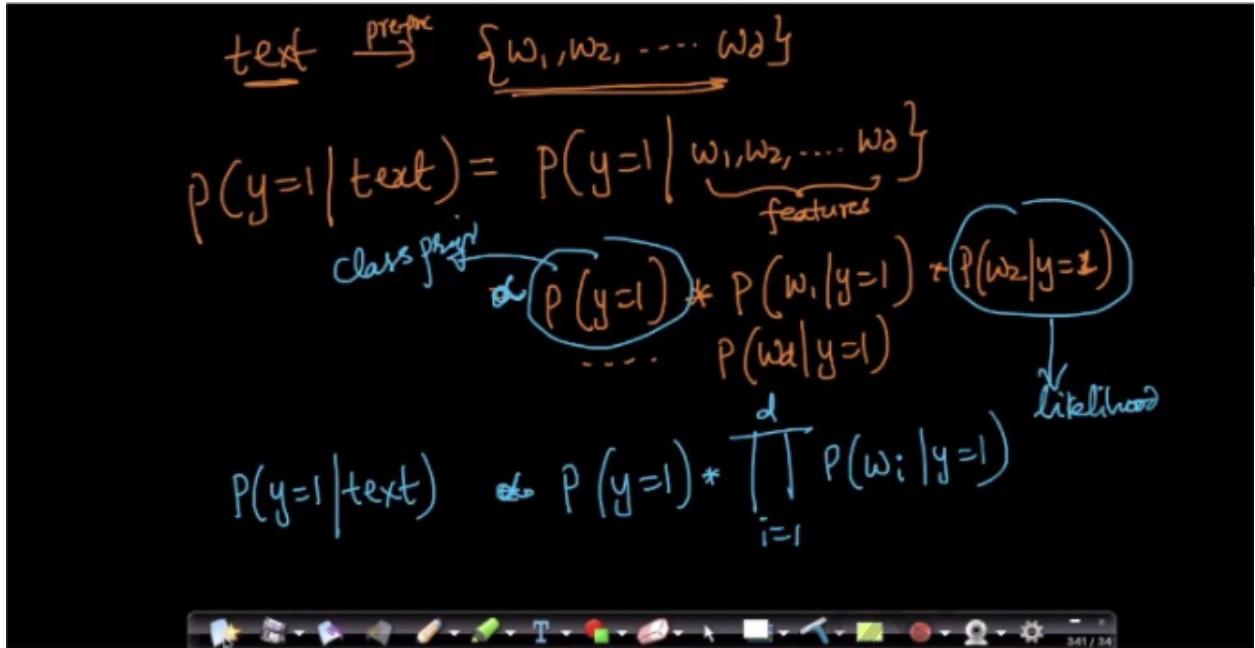
$$\underline{\text{text}} \xrightarrow{\text{preproc}} \underline{\{w_1, w_2, \dots, w_d\}}$$

$$p(y=1 | \text{text}) = p(y=1 | \underbrace{\{w_1, w_2, \dots, w_d\}}_{\text{features}})$$

class prior

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

likelihood

$$p(y=1 | \text{text}) \Leftrightarrow p(y=1) * \prod_{i=1}^d p(w_i | y=1)$$


Timestamp 7:00

Say we have a text text_q , after preprocessing we get $\Rightarrow \{w_1, w_2, \dots, w_d\}$. These individual words are features.

Now we need to calculate, $p(y = 1 | \text{text}_q)$

So,

$$p(y = 1 | \text{text}_q) \propto p(y = 1) \prod_{i=1}^d p(w_i | y = 1)$$

The above term can be derived as shown in the image above.

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

$p(y=1) = \frac{\# \text{ Train pts with } y=1}{\text{Total } \# \text{ Train pts}}$
 $p(y=0) = \frac{\# \text{ Train pts with } y=0}{\text{Total } \# \text{ Train pts}}$

Timestamp 9:10

Similarly,

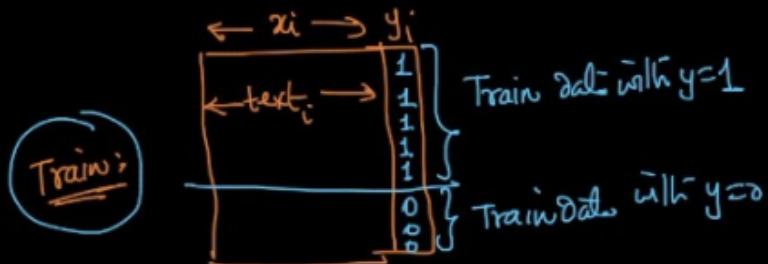
$$p(y=0 | \text{text}_q) \propto p(y=0) \prod_{i=1}^d p(w_i | y=0)$$

The prior probabilities i.e. $p(y=1)$ and $p(y=0)$ can be calculated as,

$$p(y=1) = (\# \text{train points with } y=1) / (\text{total } \# \text{ of points})$$

$$p(y=0) = (\# \text{train points with } y=0) / (\text{total } \# \text{ of points})$$

$$P(w_i | y=0) = \frac{\# \text{Train data pts with } w_i \text{ & } y=0}{\# \text{Train data pts with } y=0}$$



Timestamp 12:19

Now we can calculate the likelihoods, it can be calculated as

$$p(w_i | y=0) = (\# \text{Train data pts which contains } w_i \text{ and } y=0) / (\# \text{Train data pts with } y=0)$$

Similarly it can be done for $y=1$.

Now we have all the pieces for calculating the probabilities.

Text-classification problems

{ Spam-detection } → NB is a very good baseline

{ polarity of a review }

benchmark

① NB → acc 98%	④ DL → 98.5%
② LR →	③ GBDT →

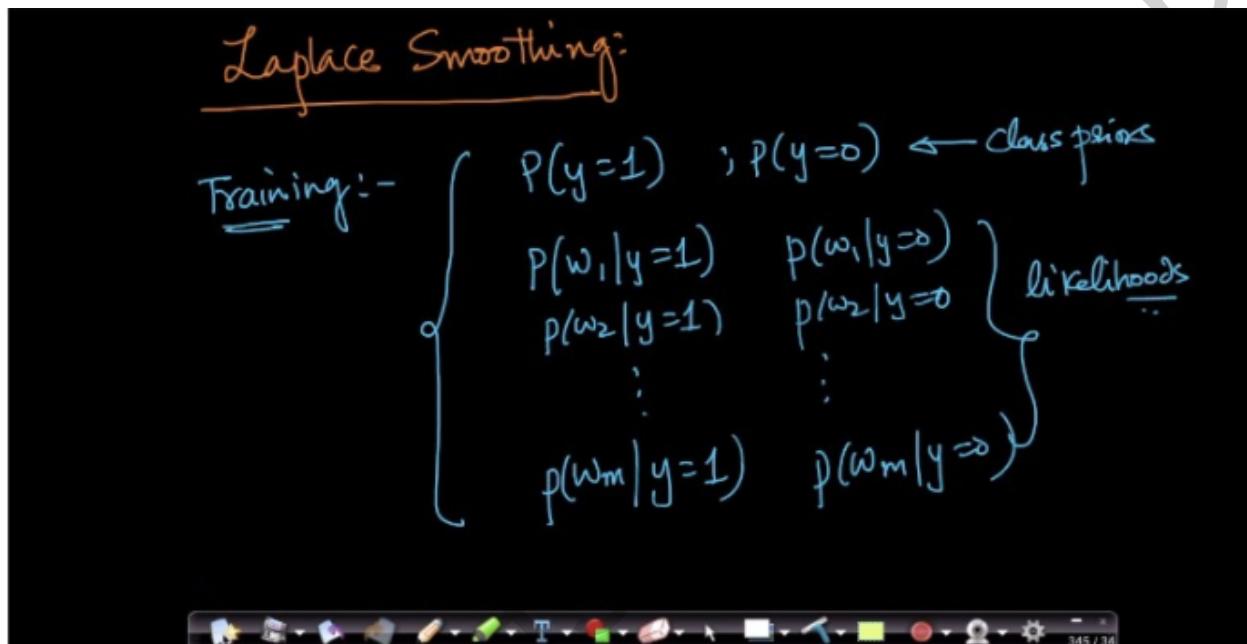
Timestamp 14:46

Naive Bayes serves as a very good baseline mode for problems like text classification, polarity of a review.etc. A baseline model is like a benchmark. We can train other models for text classification and compare with NB.

32.4 Laplace/Additive Smoothing

Laplace Smoothing:-

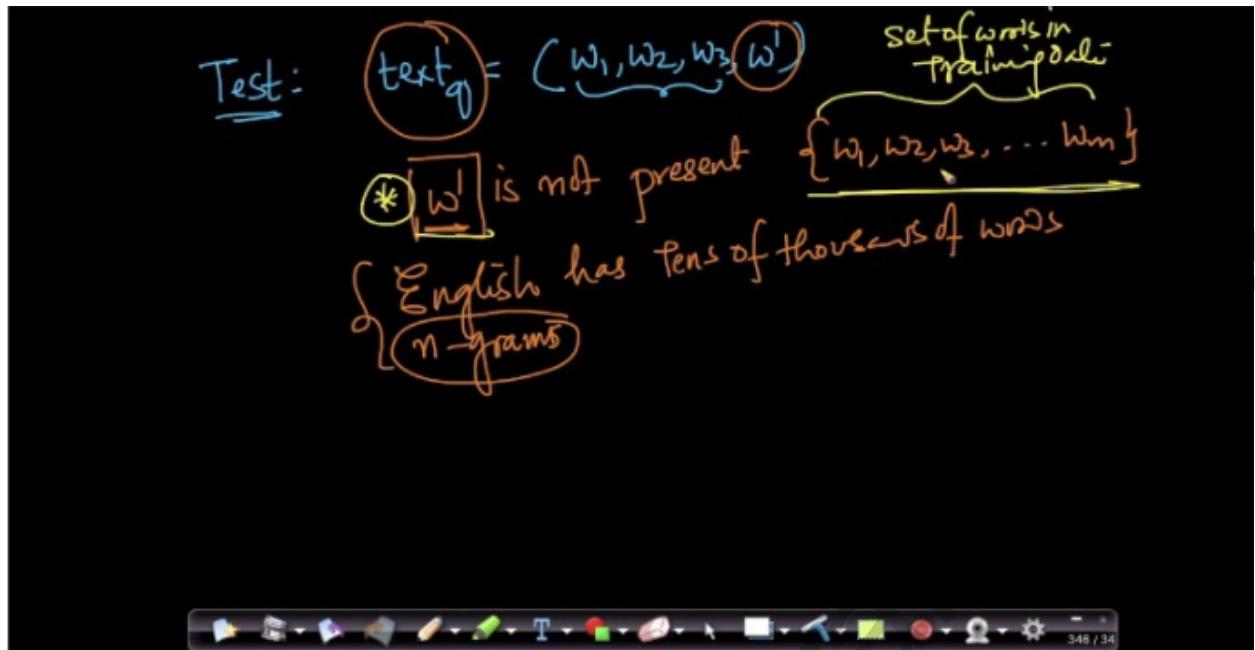
Training :-

$$\left\{ \begin{array}{ll} p(y=1) & ; p(y=0) \leftarrow \text{class priors} \\ p(w_1|y=1) & p(w_1|y=0) \\ p(w_2|y=1) & p(w_2|y=0) \\ \vdots & \vdots \\ p(w_m|y=1) & p(w_m|y=0) \end{array} \right\} \text{likelihoods}$$


Timestamp 1:14

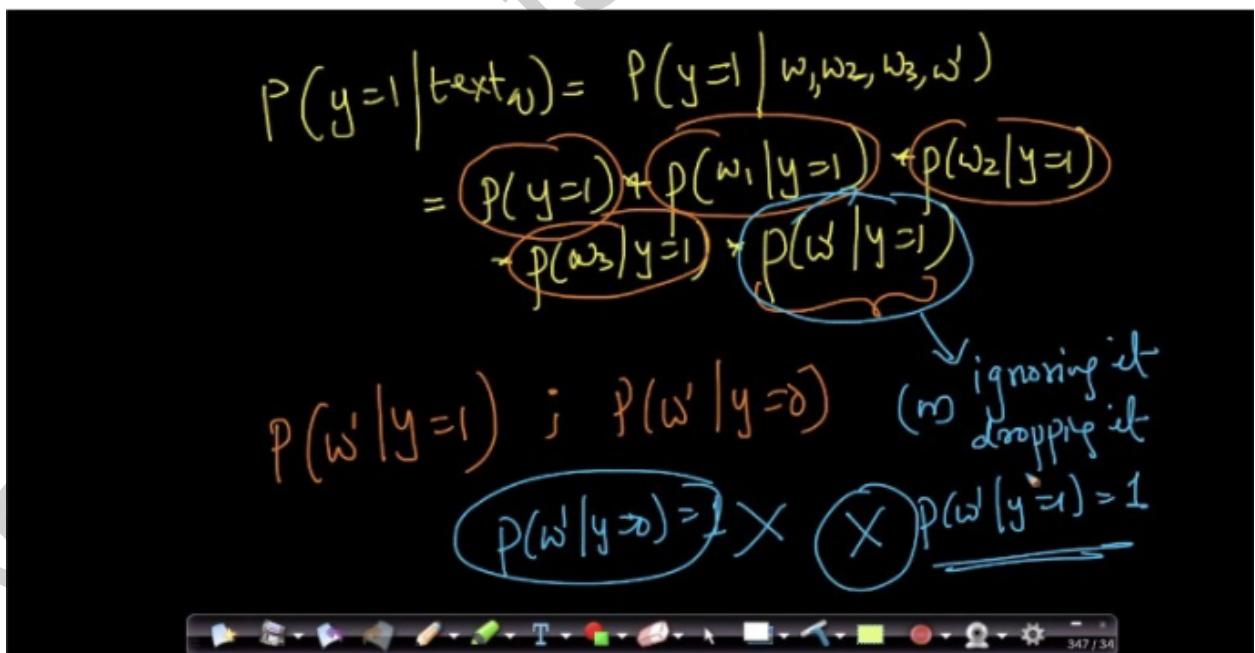
Here we will discuss something known as laplace smoothing/additive.

During the training phase we compute all the class priors and the likelihood of the features as shown in the image above. Now during train say we have a corpus of words $\{w_1, w_2, w_3, \dots, w_m\}$



Timestamp 3:24

Now during test time suppose we get a text $\text{text}_q = \{w_1, w_2, w_3, w^1\}$, it contains word w^1 , now this word is not present in our train dataset. English language contains lots of words and it is not always possible to get all the words in our corpus. Also if we are using n-grams, it may happen that some combination of words may not be present in our train data.



Timestamp 5:30

Now in order to calculate the probability of this text_q , we need to compute the likelihoods of individual words. One of the words that this text_q contains is w^l , which is not present in our train dataset.

So, the question is how to deal with this situation and what value should be used for $p(y = 1 | w^l)$ or $p(y = 0 | w^l)$.

One solution is to ignore or drop these unknown words. But this is not correct, because ignoring means that we are replacing it with 1. Meaning, $p(y = 1 | w^l) = 1$, which is incorrect.

$$\left\{ \begin{aligned} P(w^l | y=1) &= \frac{P(w^l, y=1)}{P(y=1)} \\ &= \frac{\# \text{train pts st } w^l \text{ occurs in } y=1}{n_1 \rightarrow \# \text{pts where } y=1} \\ &= \frac{0}{n_1} = 0 \end{aligned} \right.$$

Timestamp 12:17

Another solution is to calculate the actual probability, here n_1 is the no of data points where $y = 1$.

We have to calculate the term $p(w^l | y = 1)$, here the numerator denotes no. of terms where word w^l occurs and $y = 1$, which is 0, because there are no texts in the train dataset which contains the word w^l .

If we include this term in our calculation of probability the whole term becomes 0, which is not desirable.

Laplace smoothing (not Laplacian smoothing)

Additive smoothing

$$p(w^i | y=1) = \frac{o + \alpha}{n_1 + \alpha k}$$

$\alpha = 1$ Typically

$\begin{cases} 1 \rightarrow \text{present} \\ 0 \rightarrow \text{not present} \end{cases}$

$k = \# \text{ distinct values}$

w^i can take

Timestamp 10:23

Now in order to deal with this problem we use something known as laplace/additive smoothing. Remember not laplacian smoothing, which is used in image processing.

In this technique we add constant terms to both numerator and denominator. In numerator we add a α whose value is typically 1(not always), and in the denominator we add $\alpha \times k$, where k is the number of distinct values a feature can take.

In our example $k = 2$ as a word w_i can be either present or not.

So our probability for w^i becomes,

$$p(w^i | y=1) = (0 + \alpha) / (n_1 + \alpha \times k)$$

$$p(w^i | y=1) = \frac{0 + \alpha}{100 + 2\alpha}$$

Let $n_1 = 100$

$p(w^i = 0.00)$

$k=2$ because $w^i = 0 \text{ or } 1$

case 1:- $\alpha = 1 = \frac{1}{102} \neq 0$

$0 \leq p(y=1 | text_i) \leq p(w^i | y=1) \neq 0$

Timestamp 12:28

Now, we have gotten around the problem of getting 0 values as we are adding constants. Let's take a case where $\alpha = 1$, $n_1 = 100$ and in our example we had $k = 2$.

So, $p(w^i | y=1) = (0 + 1) / (100 + 1 \times 2) = 1/102 \neq 0$.

Now, the question arises why this is the correct way to deal with this problem, we could have simply replaced our probability for words which don't occur in our train data with some small value(ε).

Case 2 :- $\alpha = 10000$

$$P(\omega^1 | y=1) = \frac{0 + 10000}{100 + 20000} = \frac{10000}{20100} \approx \frac{1}{2}$$

$n_1 = 100$

$\checkmark \omega^1 = 0 \rightarrow \text{two possibilities}$

$\checkmark \omega^1 = 1$

$$\boxed{P(\omega^1 | y=1) = P(\omega^1 | y=0) = \frac{1}{2}}$$

Timestamp 15:35

We will look at another case where α is large, say = 10000. In this case we get the probabilities as 1/2 for both the classes as shown in the image above which is much better than ignoring it. As we are saying that we don't know about this unknown word and how it affects class values so let's take equal probabilities for both the classes. Introducing large α 's means introducing high bias.

Laplace Smoothing

$P(\omega_i | y=1) = \frac{\# \text{data pts with } \omega_i \text{ by } y=1 + \alpha}{\# \text{data pts } y=1 + \alpha k}$

ω_i *for all words*

present in my training data

Timestamp 17:08

So we will introduce the laplace smoothing for all the words in our training dataset, not just for unknown words. So, for every term we are adding something to the likelihoods, that's the reason it is called additive smoothing.

$$p(w_i | y=1) = \frac{2 + \alpha}{50 + \alpha}$$

$$\alpha = 1 \Rightarrow \frac{2+1}{50+1} = \frac{3}{54}$$

$$\alpha = 10 \Rightarrow \frac{2+10}{50+10} = \frac{12}{70}$$

$$\alpha = 100 \Rightarrow \frac{2+100}{50+100} = \frac{102}{250}$$

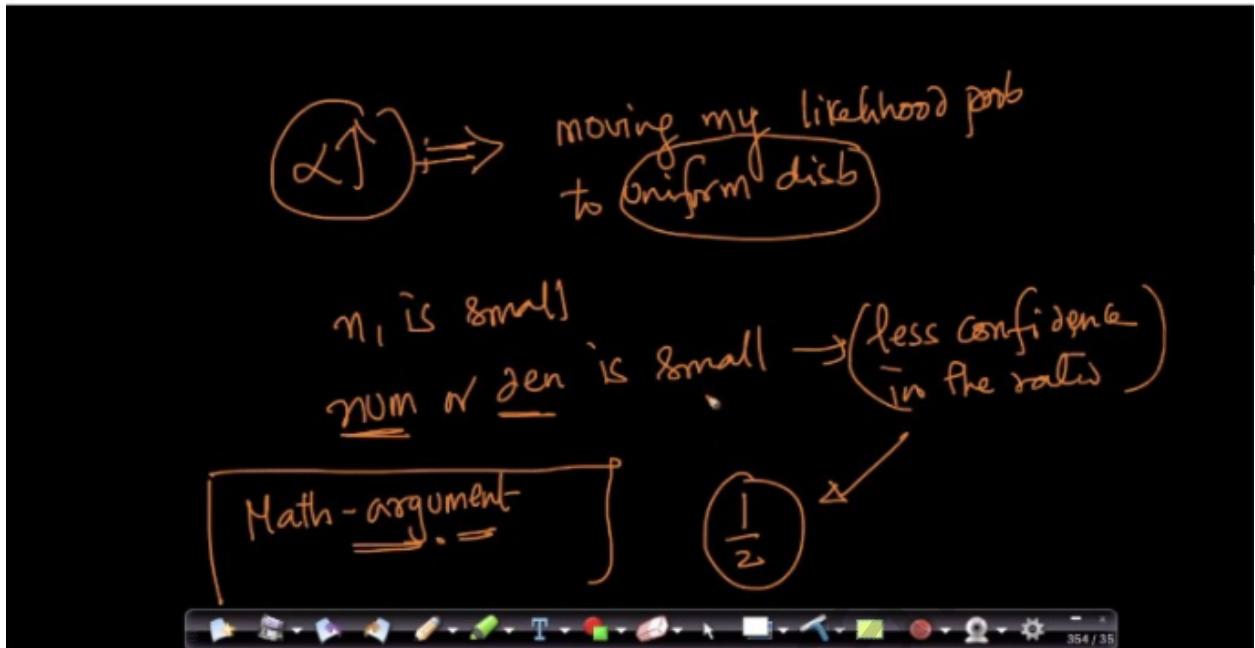
$$\alpha = 1000 \Rightarrow \frac{2+1000}{50+2000} = \frac{1002}{2050} = \frac{1}{2}$$

Timestamp 17:26

Here we will try to answer why it is called smoothing. For that let's take a word w_i , and calculate the probability $p(w_i | y = 1)$,

Say $p(w_i | y = 1) = 2/50$.

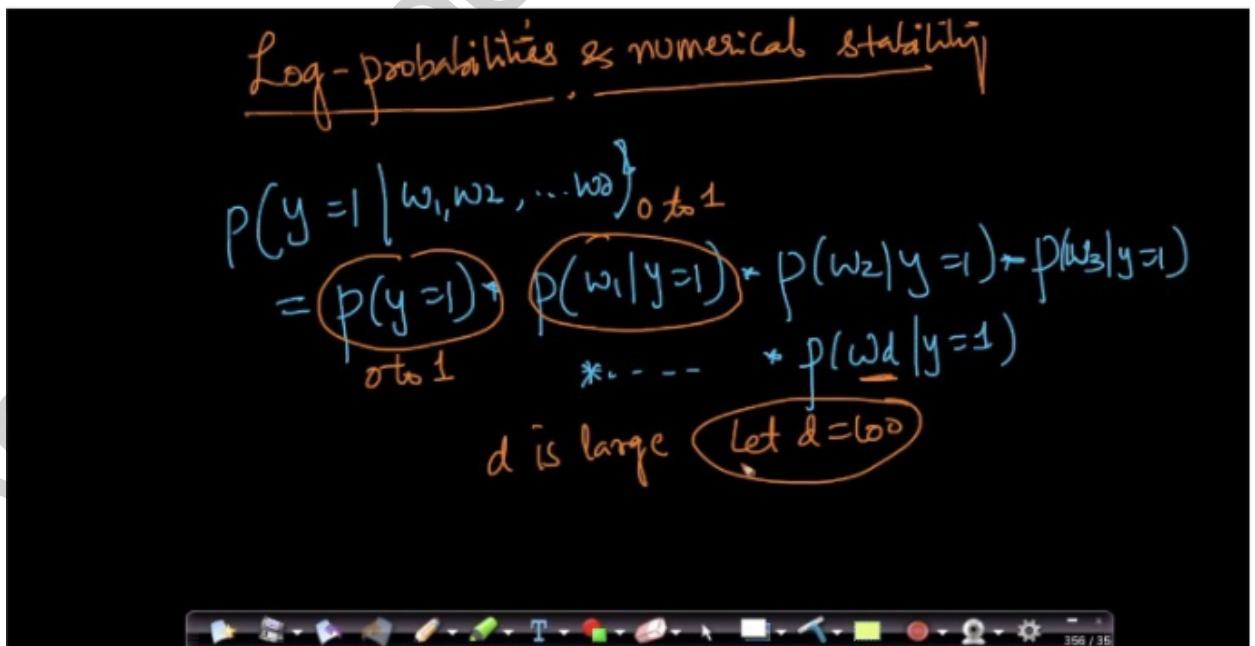
Now we will introduce laplace smoothing and calculate the probability values for different α 's. We have calculated various probabilities based on different values α as it can be seen in the image above. We can observe from the values that we got that as the value of α 's are increasing the probability values are moving towards 1/2.



Timestamp 21:38

As the value of α increases we are moving our likelihood probabilities to uniform distribution. Remember in uniform distribution all the values have equal probabilities. So it causes a smoothing of the probabilities values. In our case we had $k = 2$, that's why we were getting $1/2$. We can also use laplace smoothing if we have very few points of a certain class. The mathematical proof of laplace smoothing is beyond the scope of this course.

32.5 Log-probabilities for numerical stability



Timestamp 0:57

In order to calculate the probability we need to multiply the likelihoods of different features. Say we have a probability term $p(y = 1 | w_1, w_2, \dots, w_d)$ as shown in the above image, where d is the number of features, which say is very large like 100. Also keep in mind that all the terms are probability terms so their value lie between 0 & 1.

The slide shows a mathematical calculation and handwritten notes:

Handwritten notes:

- 100 numbers all of which 0 to 1
- ↓ numerical stability
- numerical underflow
- python double precision
- 16 significant digits
- (rounding)

Calculation:

$$\{ 0.2 \times 0.1 \times 0.2 \times 0.1 \} = 0.0004$$

Timestamp 2:40

Suppose we take four probability values 0.2, 0.1, 0.2, 0.1. If we multiply these values we get 0.0004. This is a very small number and if we multiply 100 such numbers it can lead to issues of numerical stability. In python if we have a variable of double precision it can have upto 16 significant digits. If we get more significant digits than this, it starts rounding the number. This causes numerical underflow.

Log-probabilities:

$$\log \left(\frac{P(y=1 | w_1, w_2, \dots, w_d)}{P(y=0 | w_1, w_2, \dots, w_d)} \right) = P(y=1) + \prod_{i=1}^d P(w_i | y=1) - P(y=0) + \prod_{i=1}^d P(w_i | y=0)$$

$x \uparrow ; \log x \uparrow$
monotonic fn

Timestamp 5:25

In order to deal with this numerical instability we use log probabilities. Log probabilities are nothing but we use log on either side of the equation. We can do this because $\log x$ is a monotonic function, as x increases $\log x$ also increases. If we take $\log(0.0004) = -3.39$, which is better than dealing with high precision numbers. As this negative no. has lesser significant digits.

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

$\log(a+b) = \log a + \log b$
mult
add

Timestamp 6:53

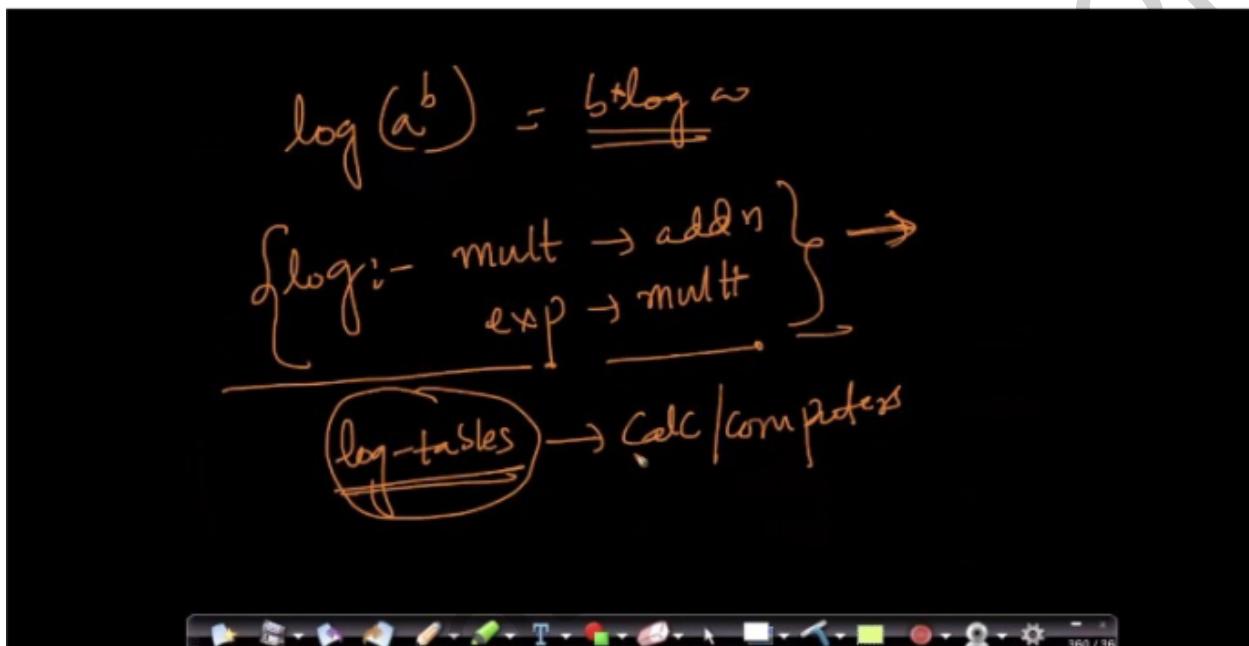
So, our probability equation is,

$$p(y=1 | w_1, w_2, w_3, \dots, w_d) = p(y=1) \prod_{i=1}^d p(w_i | y=1)$$

Taking log on both sides,

$$\log(p(y=1 | w_1, w_2, w_3, \dots, w_d)) = \log(p(y=1)) + \sum_{i=1}^d \log(p(w_i | y=1))$$

The product term is converted into a sum term because $\log(a*b) = \log(a) + \log(b)$. So we can now add log values without worrying about numerical instability.



Timestamp 9:40

Logs are extremely useful because it converts multiplication into addition

$$\log(a*b) = \log(a) + \log(b),$$

It also converts exponentiation into multiplication

$$\log(a^b) = b * \log(a)$$

In earlier times we used to use log-tables for calculating these log - values.

Log-probabilities help us prevent numerical instability because it's easier to deal with negative numbers than numbers of high precision.

32.6 Bias and Variance Tradeoff

Bias - Variance tradeoff

{ Naive Bayes :- high bias \rightarrow underfitting }
 high variance \rightarrow overfitting }

↳ in laplace smoothing
 ↳ high bias
 ↳ high variance

Timestamp 0:50

In this chapter we will deal with the bias - variance tradeoff for NB.

Recall that high bias means underfitting and high variance means overfitting. In NB only α determines whether we have high bias or high variance.

Case 1:- $\alpha = 0$

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

$$(n = 2000 \text{ pts} \rightarrow 1000 \text{ +ve}) = \frac{\# \text{Train pts with } y=1}{1000} \rightarrow \text{only 2 times} \Rightarrow \text{rare}$$

{ words that are rare } $\uparrow P(w_i | y=1)$ \rightarrow overfitting

Timestamp 3:55

Let's look at various cases indicating various values of α ,

Case 1: $\alpha = 0$

Here we have considered the value of α to be 0. Let's say we want compute the likelihood of w_i ; in our dataset we have 2000 data points, where 1000 = positive and 1000 = negative.

$$p(w_i | y = 1) = (\# \text{Train data pts which contains } w_i \text{ and } y = 1) / (\# \text{Train data pts with } y = 1)$$

Say we got this probability value as = 2/1000. We can see that this particular word w_i is very rare as it occurs in only two positive documents and also note that we don't have any smoothing as α is set to 0. Now even for this word which is very rare we will have a probability value. If we change our model data i.e. train data slightly and remove these two texts where w_i occurs, these probability changes from 2/1000 to 0/1000. This means slight change in data is changing our model considerably indicating high variance. High variance means overfitting.

So lower α values means overfitting.

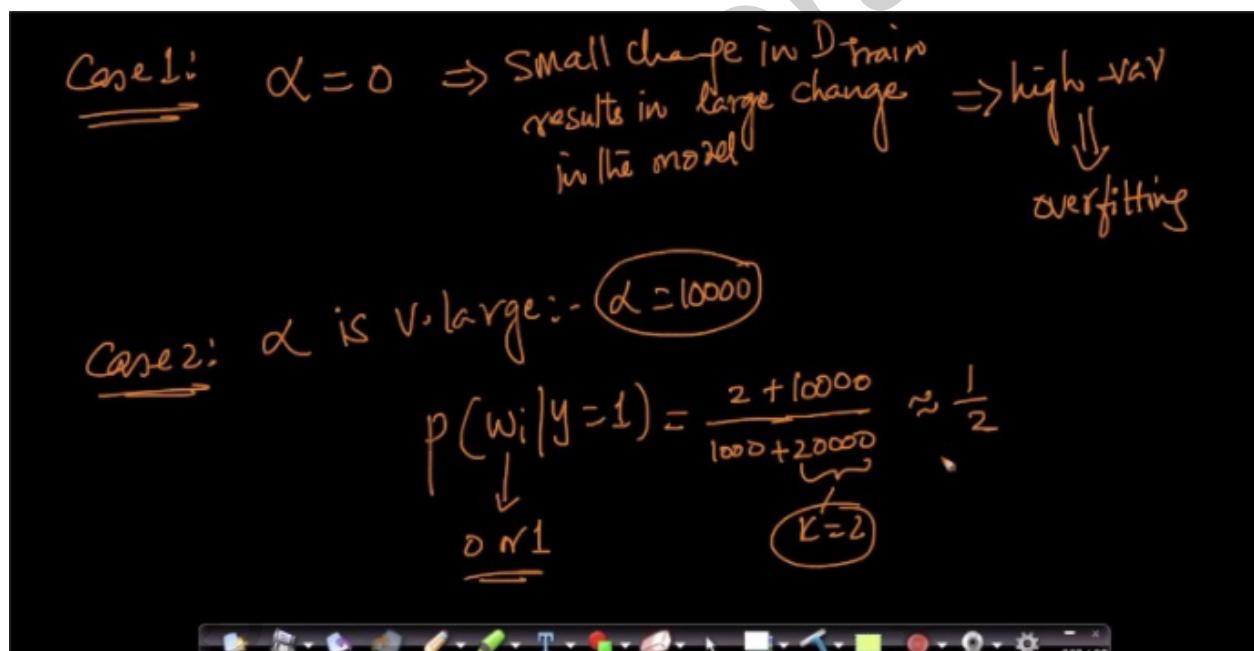
Case 1: $\alpha = 0 \Rightarrow$ small change in D_{train}
results in large change in the model \Rightarrow high variance
overfitting

Case 2: α is very large: - ($\alpha = 10000$)

$$p(w_i | y = 1) = \frac{2 + 10000}{1000 + 20000} \approx \frac{1}{2}$$

\downarrow
0 or 1

$K=2$



Timestamp 7:14

Case 2: α is very large say 10000

Again we calculate the probability of the word w_i , but this time with high smoothing. Also in our example $k = 2$.

So,

$$p(w_i | y = 1) = (2 + 10000) / (1000 + 20000) \approx 1/2$$

$$\text{For } \alpha \gg 1 \quad p(w_i | y=1) \approx p(w_i | y=0) \approx \frac{1}{2}$$

$$p(y_{\text{obs}} | \underline{x}_{\text{obs}}) \approx p(y_w=0 | \underline{x}_{\text{obs}}) \left(\approx \frac{1}{2} \right)$$

(Underfitting)

K-NN; $K = n$

$x_q \rightarrow$ Same class label (+ve)

$n_1 > n_2$

$n_1: +ve$

$n_2: -ve$

Timestamp 11:13

For every class we will get the same probability. As our α term dominates the probability values and less importance is given to actual probability values. Our actual prediction depends only on the class priors. So, whichever class is present in the majority our model will always predict that class. This means we are underfitting.

So higher α values means underfitting.

We can compare this situation with KNN, where we can set our neighbours $K = n(\# \text{data points})$, then for every query we will get the same prediction. The prediction will be the class, which is present in majority in the train dataset. Leading to underfitting.

(Q) How to find the right α

→ KNN; right $K \rightarrow$ using Simple CV
10-fold CV

→ right α : → Simple CV
10-fold CV

Timestamp 12:45

Now how to find the right α , we can find this using cross validation, like we did in KNN. α is called a hyperparameter. Similarly the value of k in KNN is also a hyperparameter.

32.7 Feature Importance and Interpretability

Feature importance & interpretability

(NB) - feature imp

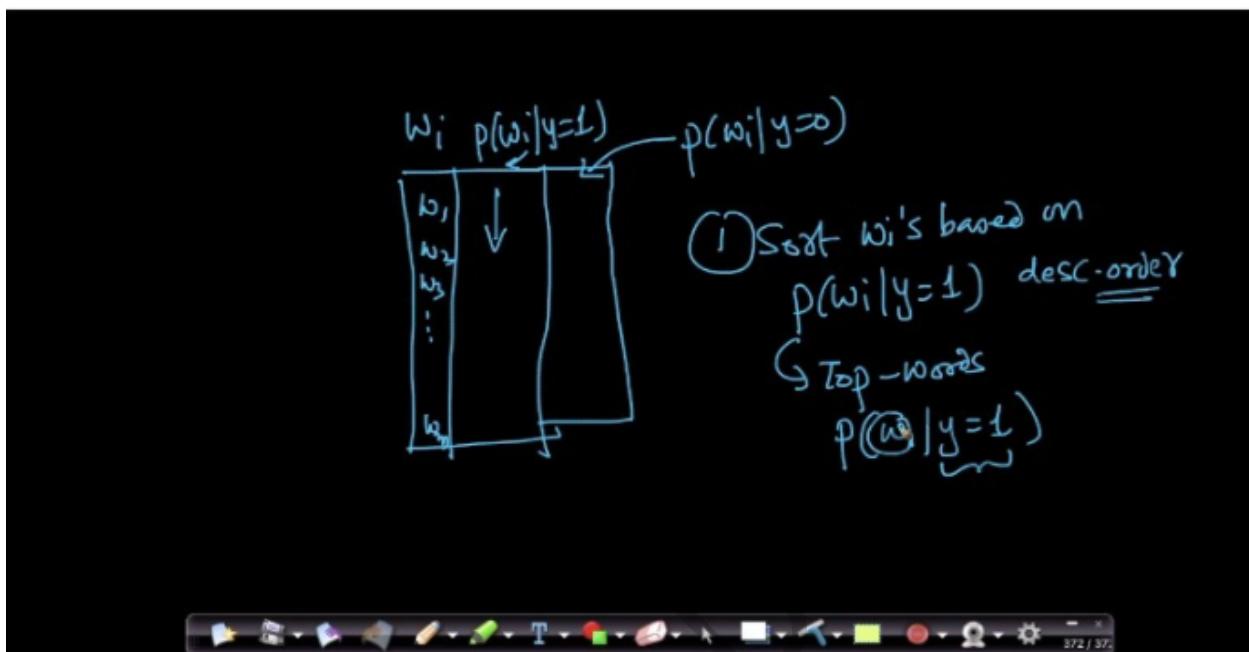
Likelihood:

$w_i, p(w_i | y=1)$

for all $w_i, p(w_i | y=0)$

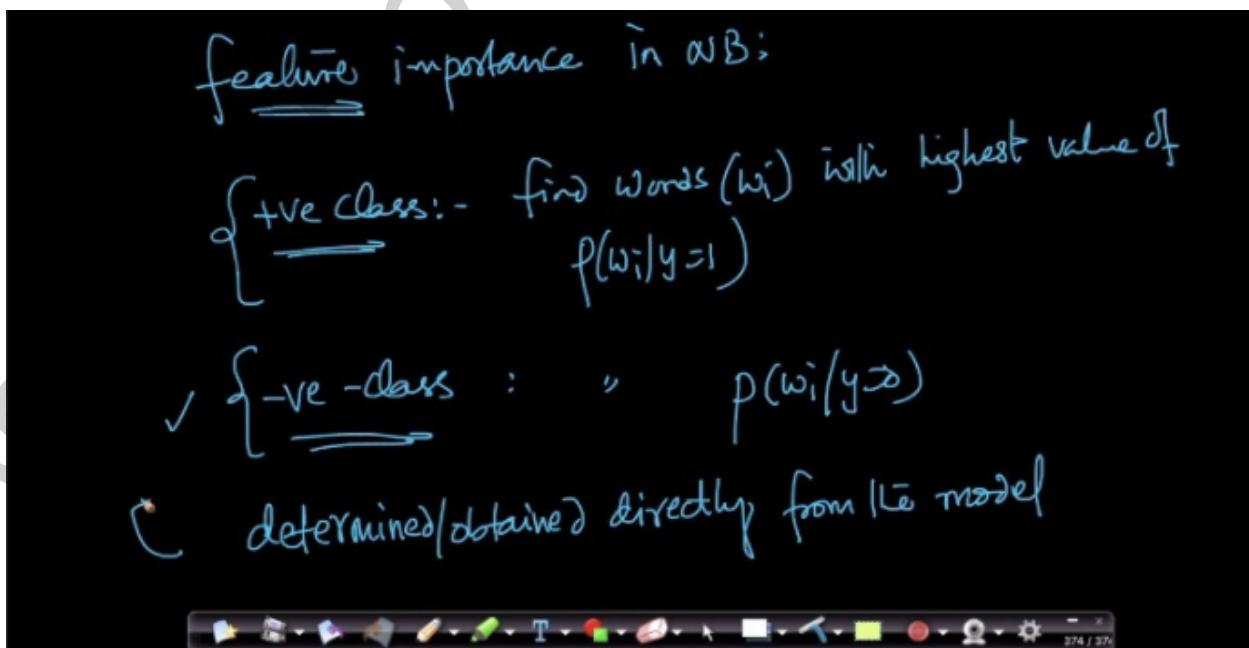
Timestamp 1:07

In this chapter we will try to find the feature importance and interpretability for a prediction. In order to find these we need to calculate the likelihood for all words.



Timestamp 2:35

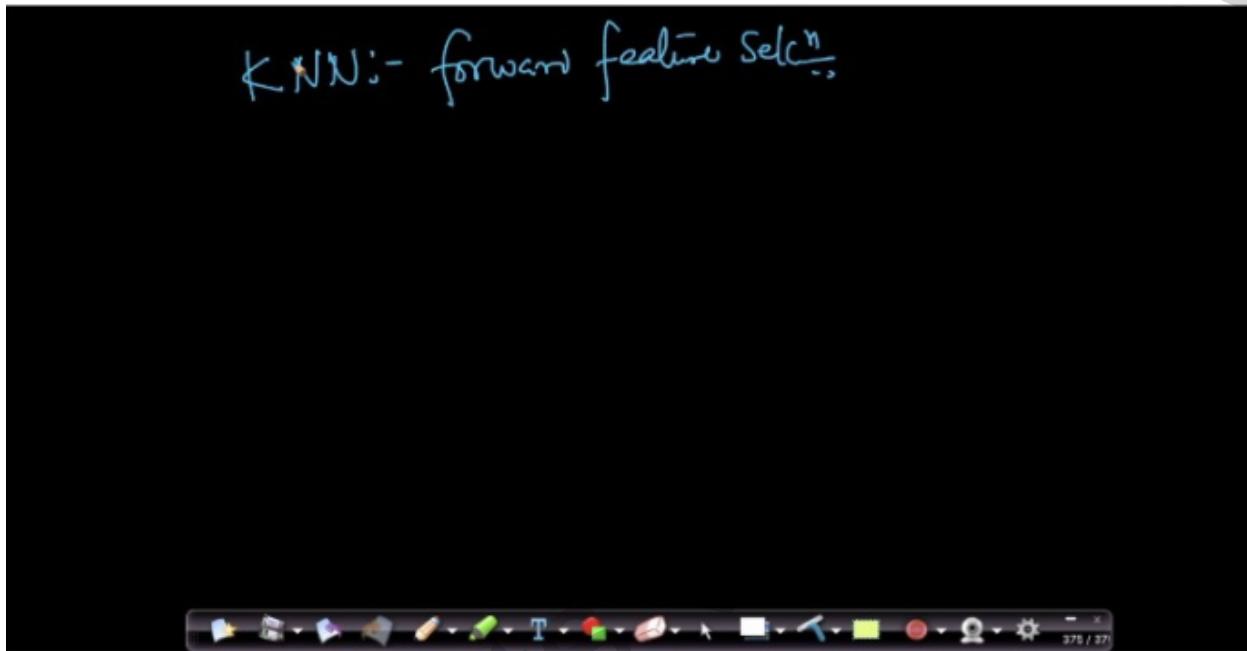
We will calculate the likelihood i.e. $p(w_i | y = 1)$ and $p(w_i | y = 0)$ for each word and store it into a table. Now we need to sort this table as per their probability values in descending order. Words which occur on top of this sorted table i.e. highest probabilities are the most frequent words in a particular class.



Timestamp 5:29

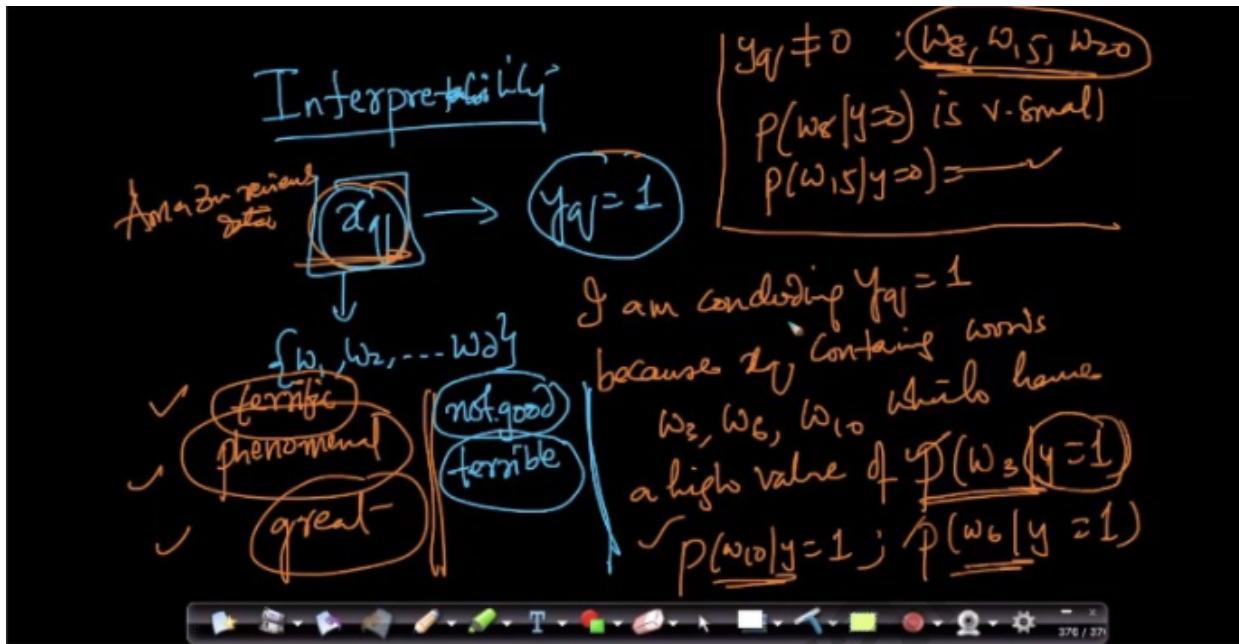
So features/words which possess the highest probability values are the most important features for a particular class. As when we compute the total probability, features/words which have the highest likelihood have the highest impact in total probability value.

Feature importance in case of NB can be directly computed from the model.



Timestamp 5:44

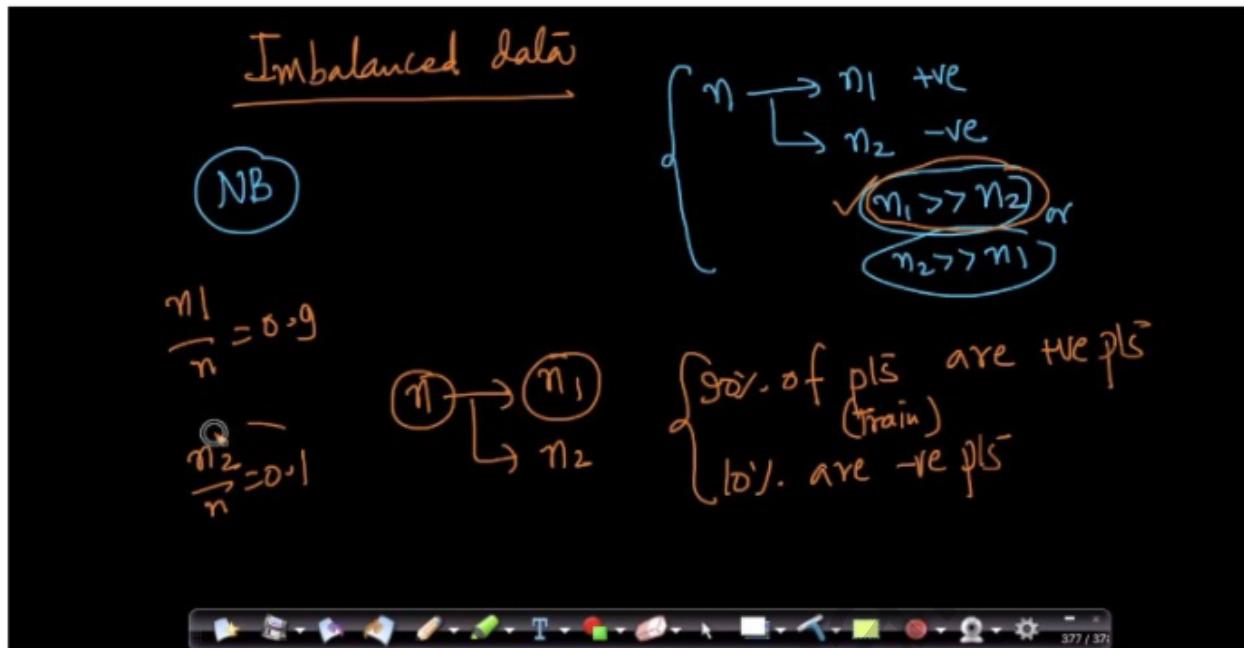
This is unlike KNN where we used techniques like forward and backward feature selection for determining the feature importance.



Timestamp 9:40

Now let's talk about interpretability, suppose we have a query point x_q , containing words $\{w_1, w_2, w_3, \dots, w_d\}$, suppose we got $y_q = 1$. Now we can look at the individual probability terms and find those words which have the highest values, say we got $\{w_3, w_6, w_{10}\}$. From this we can conclude that we got the class label $y_q = 1$, because of the high likelihood of these words. These words determined our class label.

32.8 Imbalanced data



Timestamp 1:48

In this chapter we will look at the effect of an imbalanced dataset on Naive Bayes. Let's quickly define imbalanced dataset, say we have dataset containing n data points, where n_1 is the number of positive points and n_2 is the number of negative points. Now if n_1 is significantly greater than n_2 i.e. $n_1 \gg n_2$ or n_2 is significantly greater than n_1 i.e. $n_1 \ll n_2$, we call such a dataset as an imbalance dataset.

For our case let's assume that $n_1 \gg n_2$, say 90% of the data are positive and the remaining 10% are negative.

$$\text{So, } n_1/n = 0.9 \text{ and } n_2/n = 0.1$$

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

(let) $\frac{n_1}{n} = p(y=1)$

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

(let) $\frac{n_2}{n} = p(y=0)$

imbalanced data
① class priors → majority / dominating class has an advantage

Timestamp 3:35

If we calculate the probabilities for each of the classes, we need to calculate two terms: class priors and likelihoods. Suppose the likelihood is same for both the classes then we are left with priors. Prior for the majority class will dominate over the prior for the minority class. In our case the positive class will dominate over the minority class.

So, in case of an imbalanced dataset, the majority class has an undue advantage.

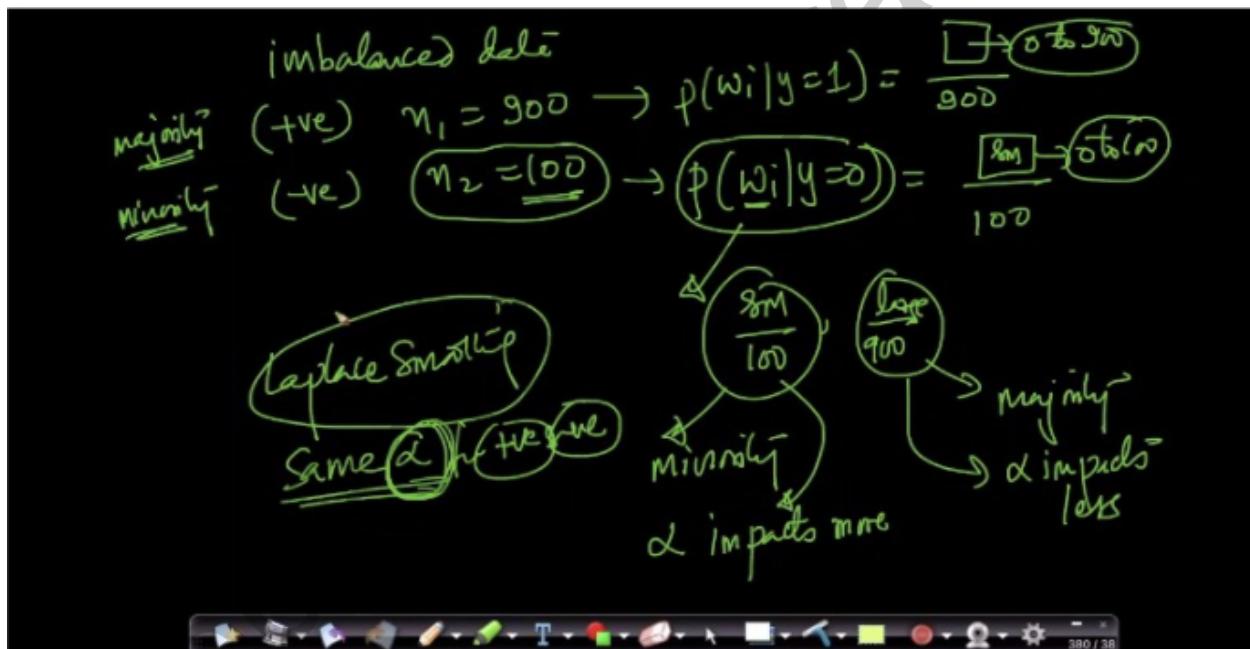
- solution:
- ① upsampling or down sampling
 - $n_1 \approx n_2$ $p(y=1) = p(y=0) = \frac{1}{2}$
 - ② drop $p(y=1) \neq p(y=0)$
 $p(y=1) = p(y=0) > 1$
 - ③ modified NB to account for class imbalance
 ↳ not often used

Timestamp 6:18

Here we will look at ways to handle imbalance,

1. Upsampling or Downsampling - We can either upsample the minority class or downsample the majority class to make both these class equal, i.e. $n_1 = n_2$. If we do this we are effectively making our class priors = 1/2
2. We can simply drop class priors. Dropping class priors effectively means setting them to 1. If we observe carefully both point 1 and 2 are equivalent as in both the cases we are setting the class priors to the same value. Since we only compare two probabilities it doesn't matter what priors we use as long as they are the same.
3. Modified NB, we can modify our NB to account for class imbalance, however it is seldom used.

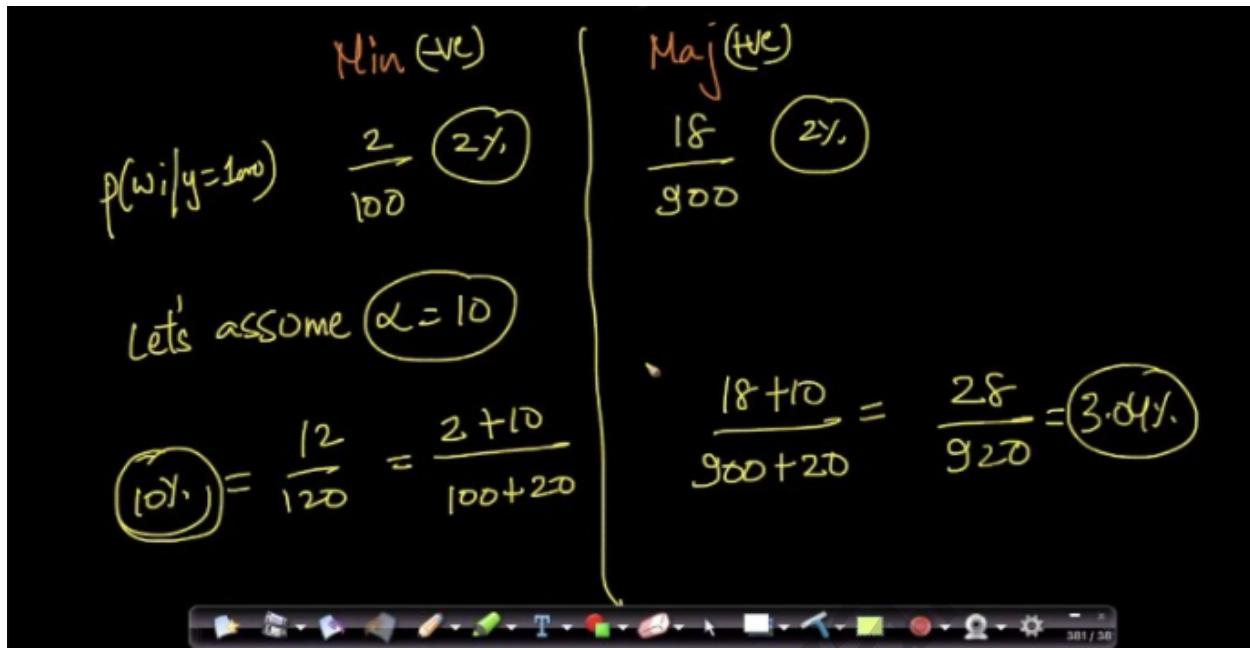
Mostly point 1 and 2 are used in practice.



Timestamp 8:56

Here we will look at the effect of laplace smoothing in an imbalance dataset.

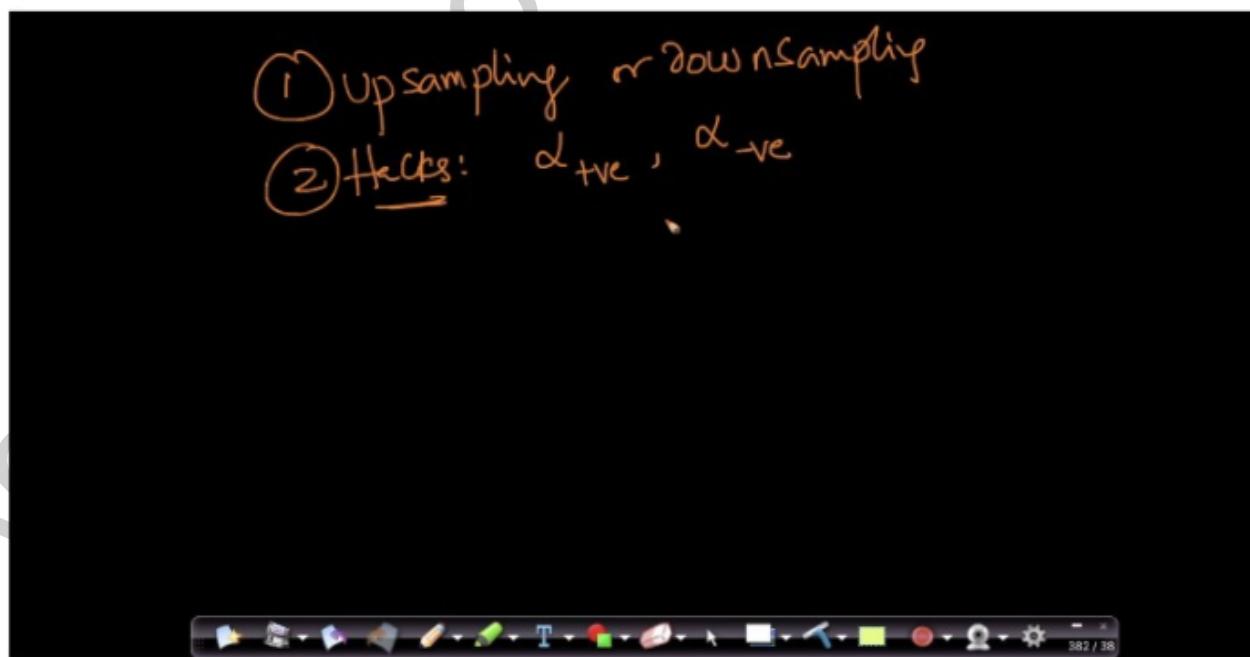
Let's say we have a dataset having 900 positive points and 100 negative points, here the positive class is the majority class. Now let's say we calculate the likelihood of the word w_i for both these classes as shown in the image. If we apply laplace smoothing on both these two terms the minority class is impacted more than the majority class.



Timestamp 11:15

Let's see it with an example.

As we can see from the image above we calculated the probabilities for both the classes without using any smoothing. For both the majority and minority class we got the probability % to be 2%. Now after applying laplace smoothing with $\alpha = 2$, the probability % for the minority class shot up to 10% from 2% whereas for the majority class it went to 3.04%, showing that laplace smoothing affects the minority class more.



Timestamp 12:24

One of the hacks we can use to handle the impact of laplace smoothing is to use different α 's for different classes. However this technique is not grounded in theory. It's better to follow techniques like upsampling/downsampling to handle class imbalance.

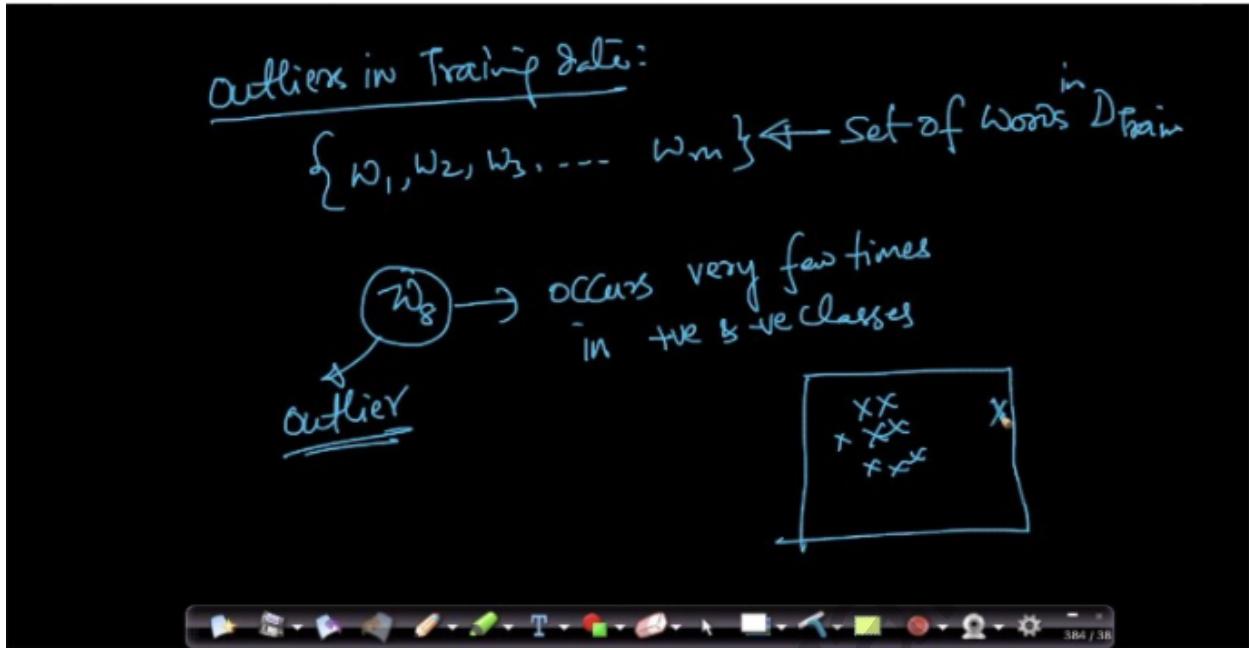
32.9 Outliers

The image shows handwritten notes on a black background. At the top left, there is a circle containing the text "NB". To its right, the word "Outliers :-" is written above the text "Text-classfn:-". Below this, another circle contains the text "outlier at test time". To the right of this circle, the formula $x_q = w_1, w_2, w_3, \dots, w_n$ is written, with a note below it stating " $w_i \notin \{w_1, w_2, \dots, w_m\}$ ". To the right of the formula, a large arrow points to the right with the handwritten note "Set of words that I have seen in Train". Below these notes, the formula for Laplace smoothing is written: $P(w/y) = \frac{O + \alpha}{n_1 + 2\alpha}$. At the bottom of the notes, the text "✓ Laplace Smoothing ✓" is written with two checkmarks. A decorative scroll graphic is visible on the right side of the slide.

Timestamp 1:42

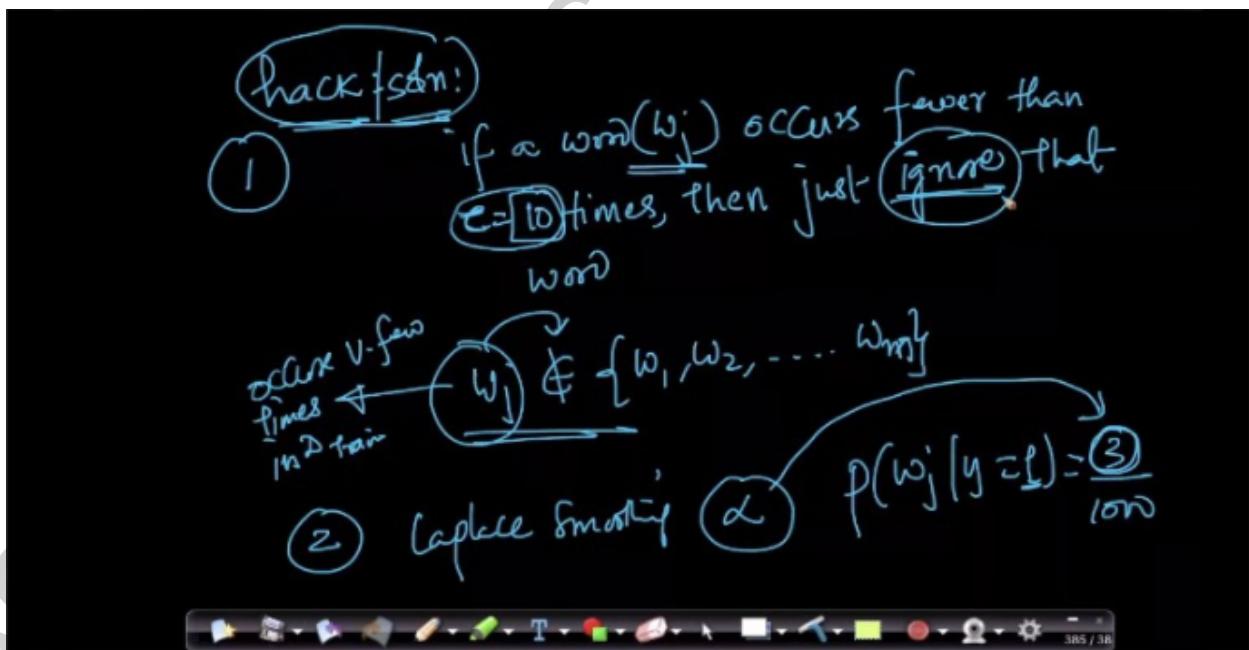
In this chapter we will see how to deal with an outlier.

Suppose we are doing text classification and during test time we see a word w^1 , which was not present during train time. If we calculate the likelihood of this value we will get 0 making the whole probability 0. So in order to deal with this situation we use laplace smoothing which we learnt in earlier chapters.



Timestamp 2:54

There can be outliers even in the training data. An outlier in training data can be defined as a word/feature which occurs very few times in training data. It is such a point around which there are no points.



Timestamp 5:06

We can deal with these problems in training data by simply ignoring the features/words. We can set a threshold and say if a particular feature occurs below a certain threshold we will ignore such words.

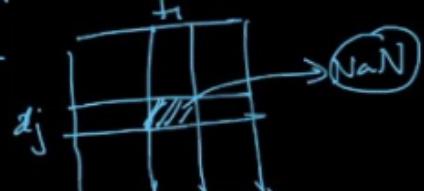
Another solution is to use laplace smoothing with an appropriate value of α .

32.10 Missing values

Missing-Values :-

① Text-data :- Text :- $\{w_1, w_2, \dots, w_n\}$
↳ no case of missing data

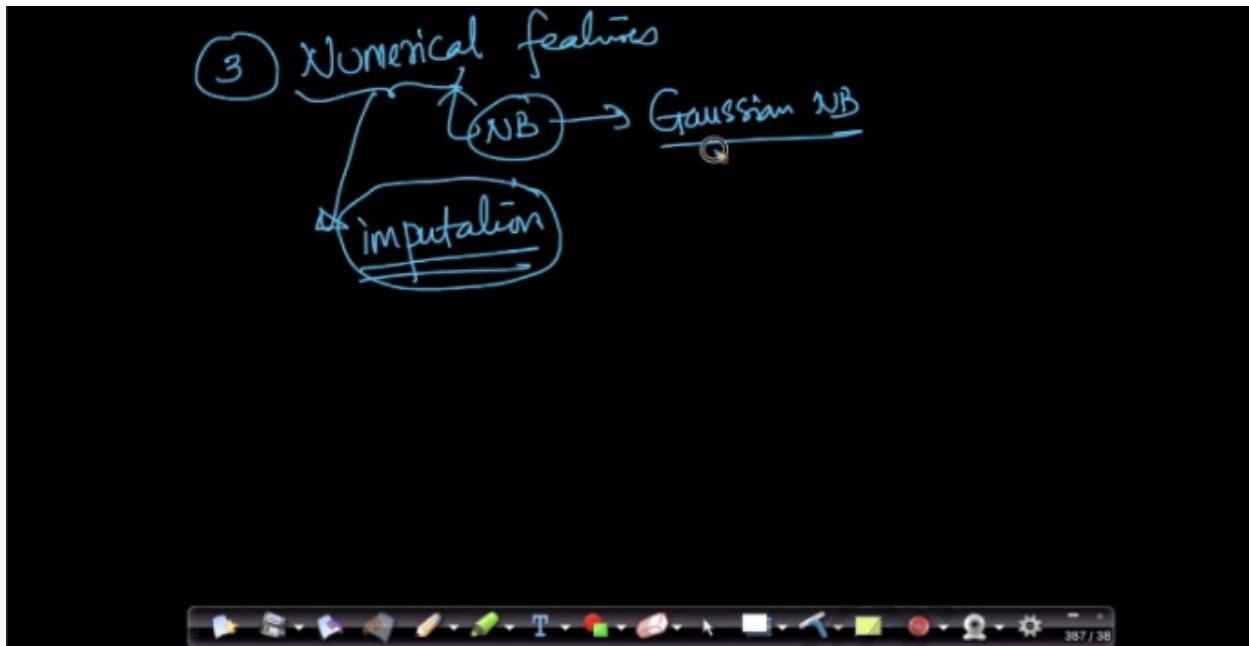
② Categorical features :-
 $f_i \in \{a_1, a_2, a_3\}$
→ consider NaN as a category
 $f_i \in \{a_1, a_2, a_3, \text{NaN}\}$



Timestamp 1:56

In this chapter we will talk about dealing with missing values.

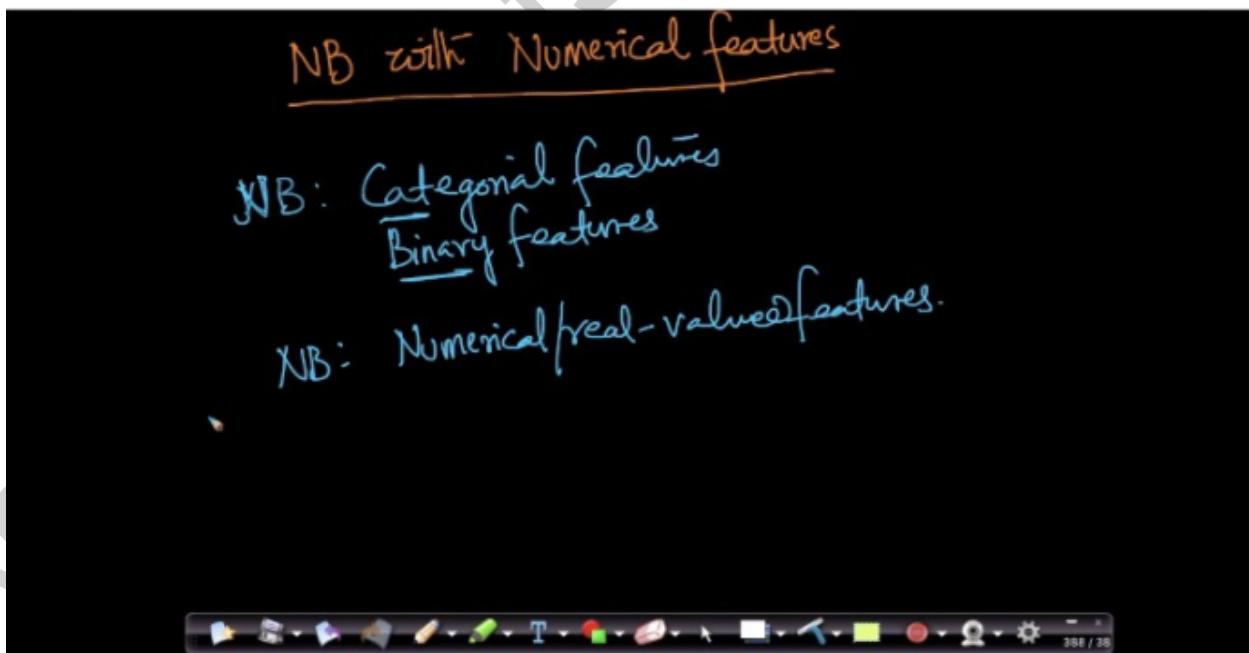
1. Text data - If we are dealing with text data there is no case of missing data, either the word is present or not present.
2. Categorical features - Say we have feature f_i , having categories $\{a_1, a_2, a_3\}$. Now suppose in one of the rows this feature value is missing. So we can consider the missing value NaN, itself as a feature. So, our feature set becomes $\{a_1, a_2, a_3, \text{NaN}\}$. This is one of the hacks that one can use. Also it is not limited to Naive Bayes but can be used anywhere.



Timestamp 2:55

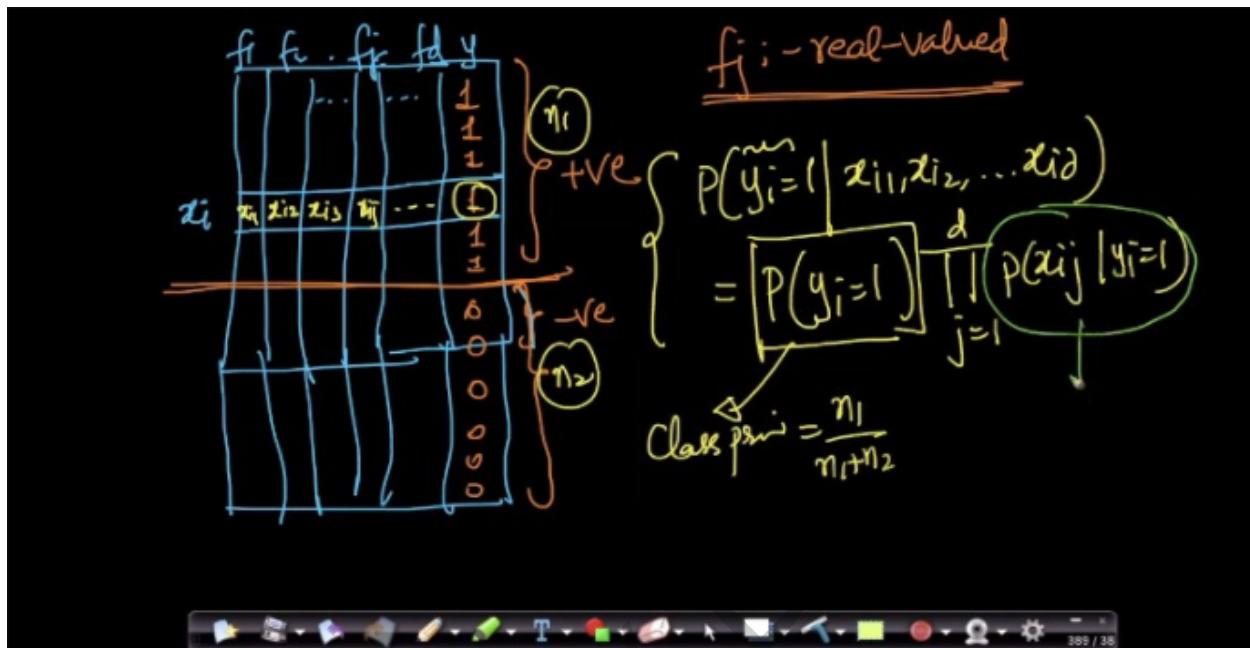
3. Numerical Features - We can handle missing numerical features by using various imputations like mean, median.etc.

32.11 Handling Numerical Features(Gaussian NB)



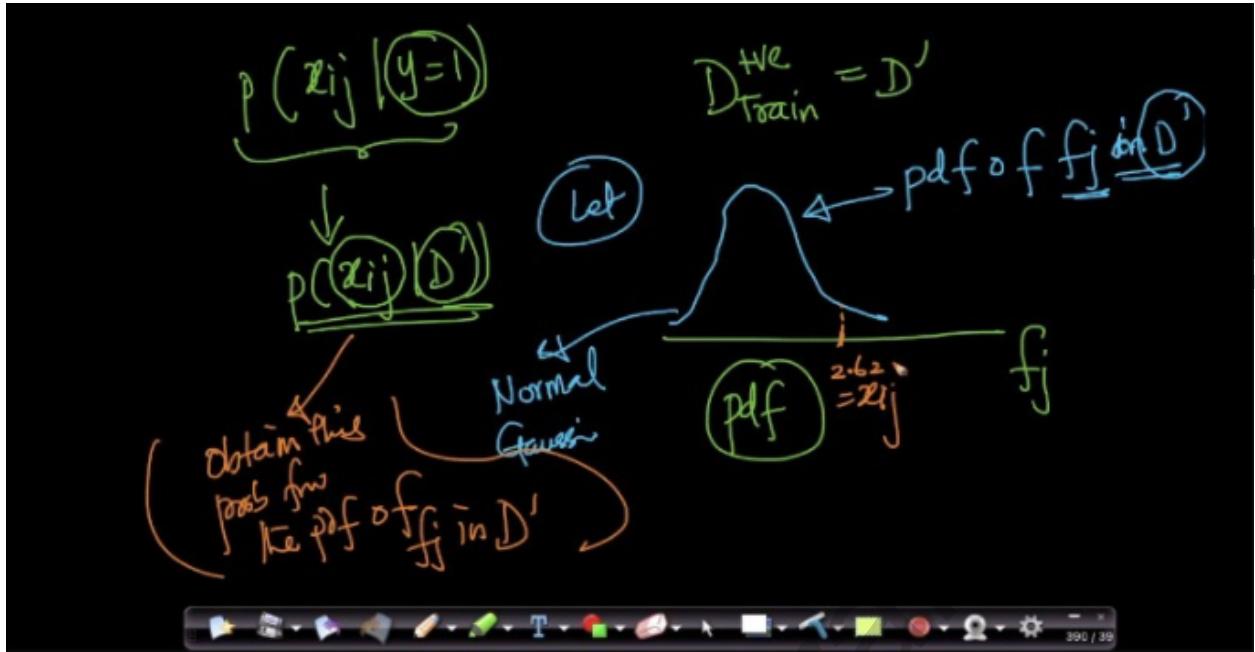
Timestamp 0:36

In this chapter we will learn how to deal with numerical(real valued) features in NB. Till now we have only worked with categorical and binary(words) features. In order to deal with numerical features we will use something known as Gaussian NB.



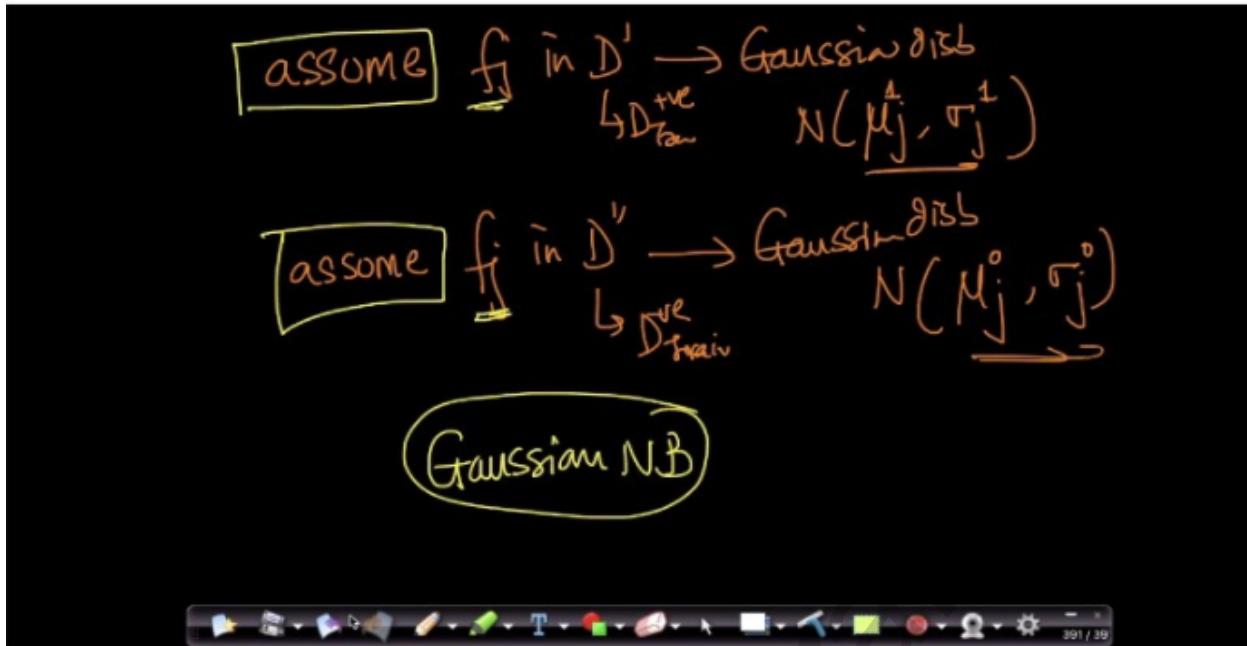
Timestamp 3:56

Consider a dataset where we have d features out of which say feature f_j is real valued. Now let's divide our dataset into groups one for all the positive labels(n_1) and another for all the negative labels(n_2). Consider a positive data point x_i , and it's j th feature i.e. x_{ij} , likelihood for this term can be written as $p(x_{ij} | y_i = 1)$. Now how should we calculate this term?



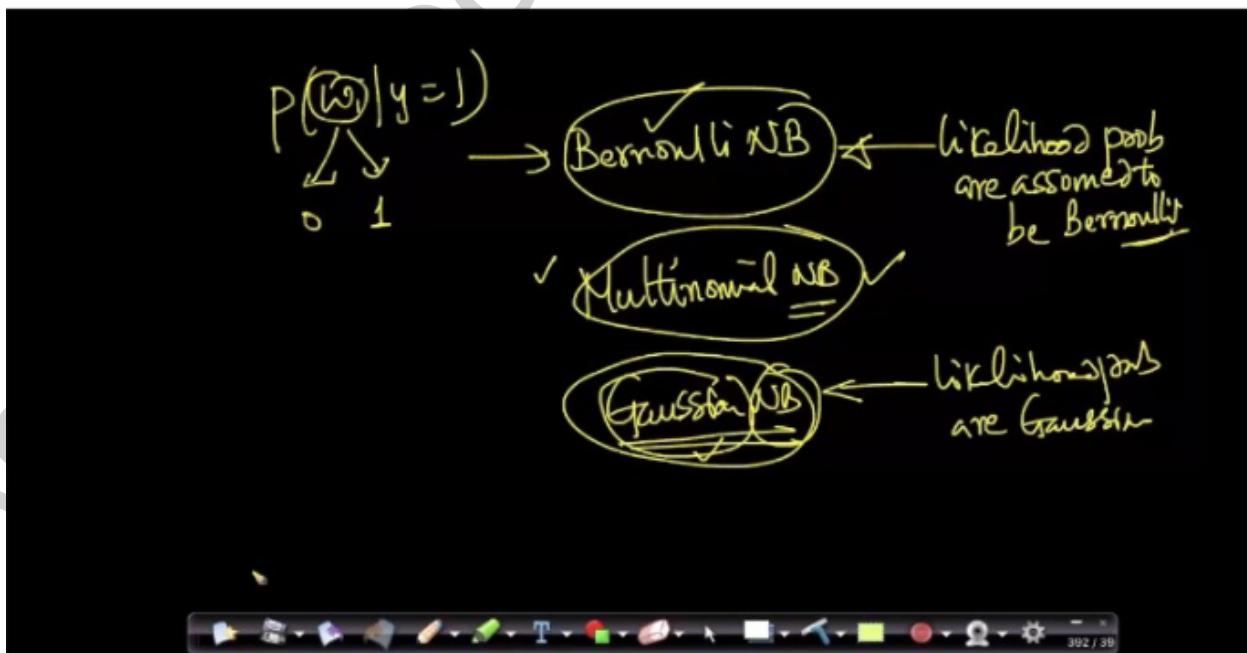
Timestamp 7:16

In order to calculate this term $p(x_{ij} | y_i = 1)$, where x_{ij} corresponds to feature j of datapoint x_i , we will first take only the positive data points from the whole dataset, let's call this D^1 . Now from all the data points of D^1 , we take feature f_j , and draw a pdf corresponding to its values as shown in the image above. Now, in order to get a probability for any value say x_{ij} , we can calculate the pdf value from the pdf curve as shown in the image above.



Timestamp 10:08

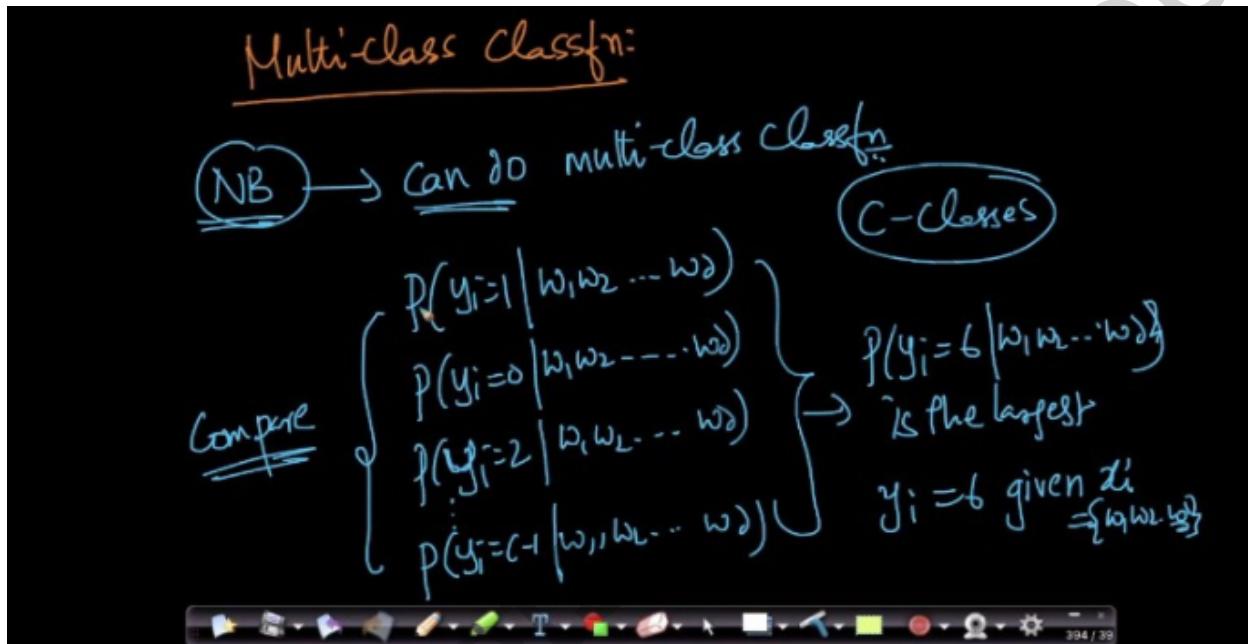
When dealing with numerical features we assume that these values follow a Gaussian Distribution. For the positive dataset D^1 , we assume that the real valued feature f_j containing only the values from D^1 follows a Gaussian Distribution with mean μ_j^1 and standard deviation σ_j^1 , computed only from the values of the D^1 . Similar steps can be followed for the negative dataset, as shown in the image above. This is called a Gaussian NB.



Timestamp 12:08

While dealing with words w_i , it could take two values either 0 or 1. It follows Bernoulli Distribution. It is also called Bernoulli NB. So, we can figure out what distribution a feature is following and based on the pdf/pmf of these distributions and we can calculate likelihood probabilities.

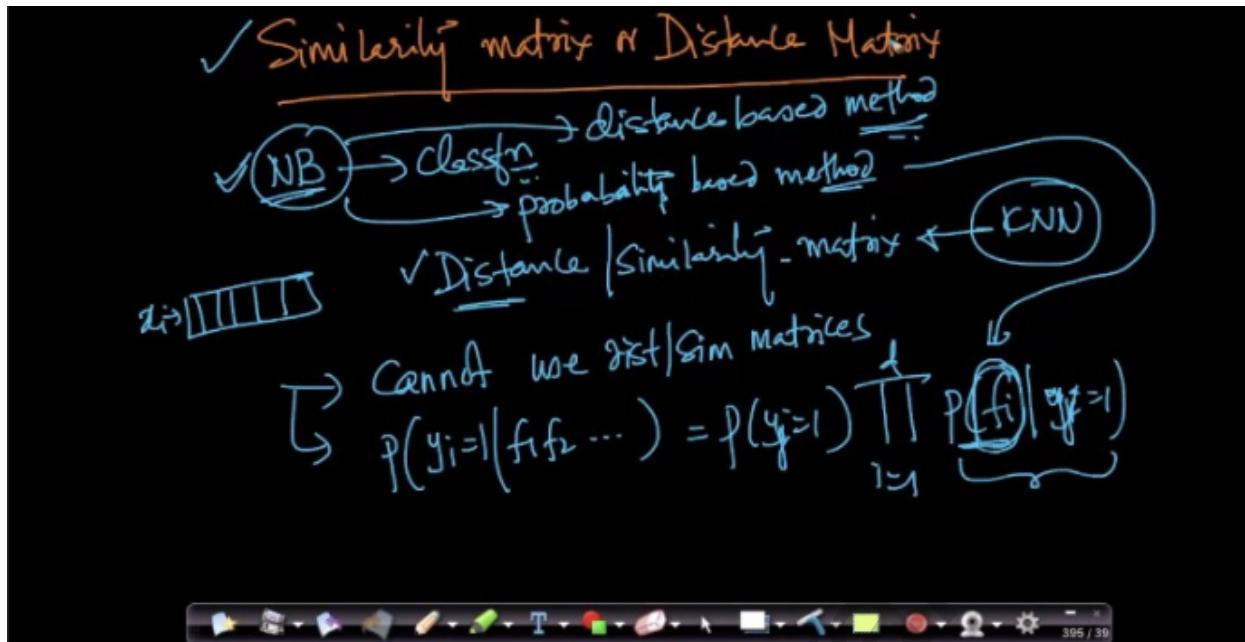
32.12 Multiclass classification



Timestamp 1:44

Naive Bayes can easily handle Multi Class Classification. We just need to calculate the probability for each of the classes, compare them and then select the class which has the highest probability.

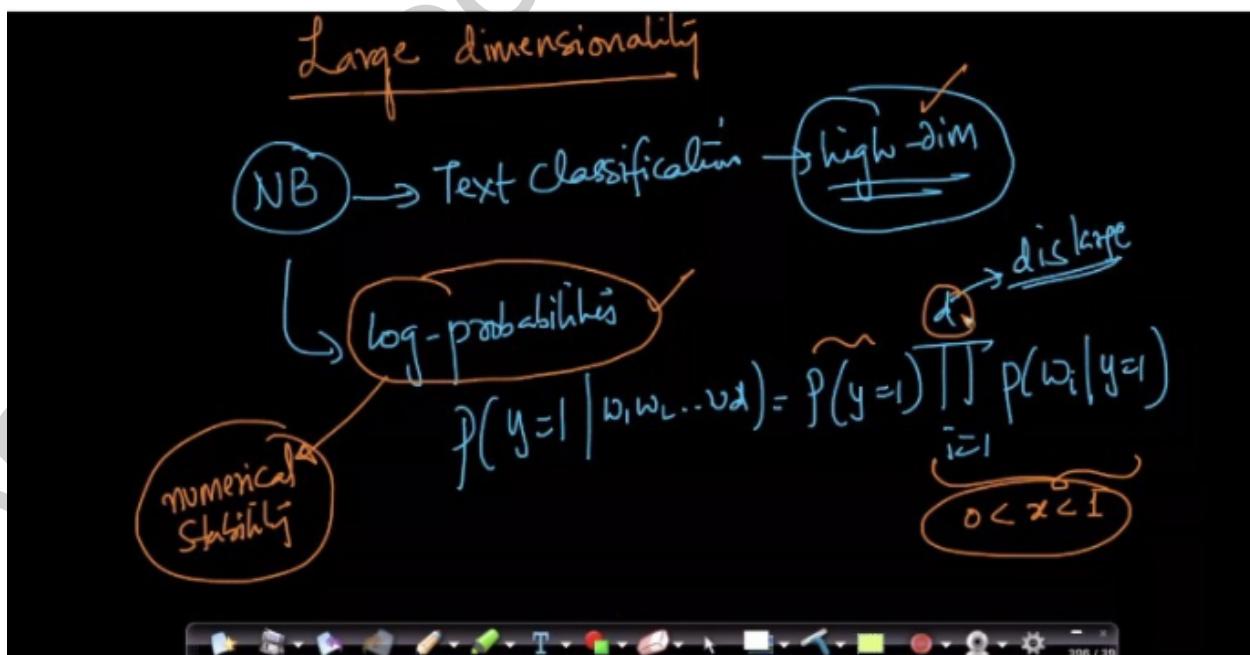
32.13 Similarity or Distance Matrix



Timestamp 2:15

Similarity or Distance matrix is usually used when it is difficult to represent a data point in the form of a vector, like we saw in KNN while doing pharma example. Naive Bayes cannot handle similarity or distance matrix in general, because it is not a distance based method rather a probability based method here we do not calculate the distance between two data points, rather we find out their likelihood probabilities.

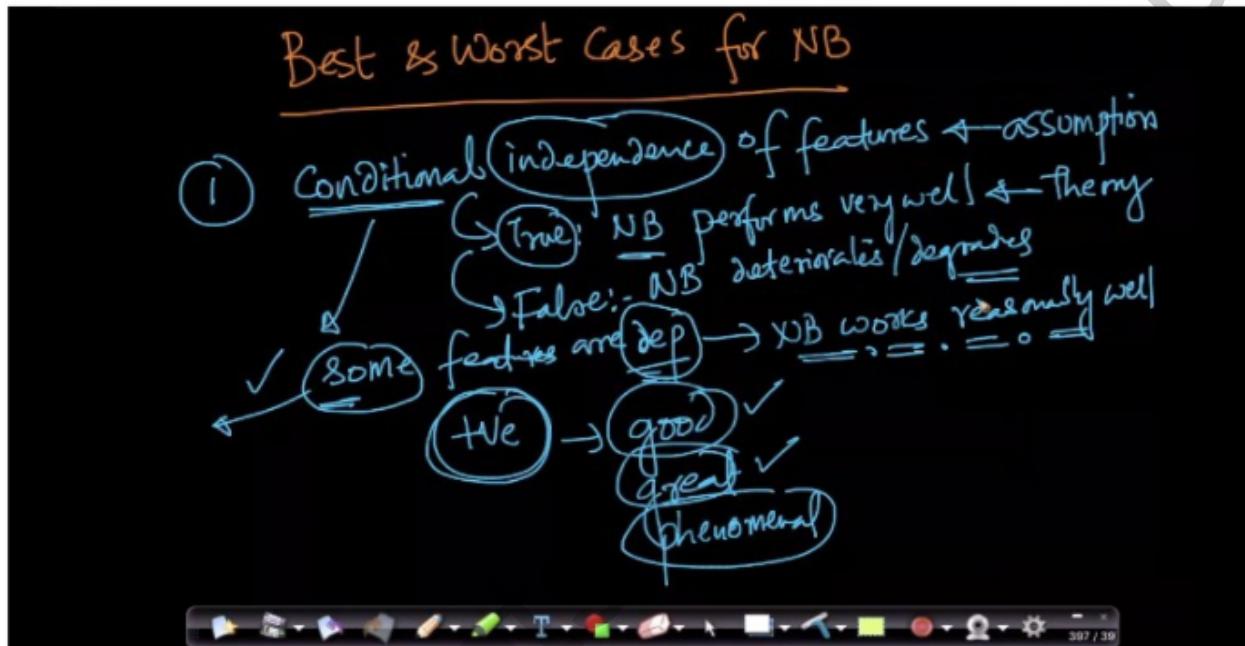
32.14 Large Dimensionality



Timestamp 1:40

Naive Bayes can easily handle data of large dimensions. While dealing with text classification we had text data which itself is of very high dimension. Only thing to keep in mind is to use log-probabilities to avoid numerical instability.

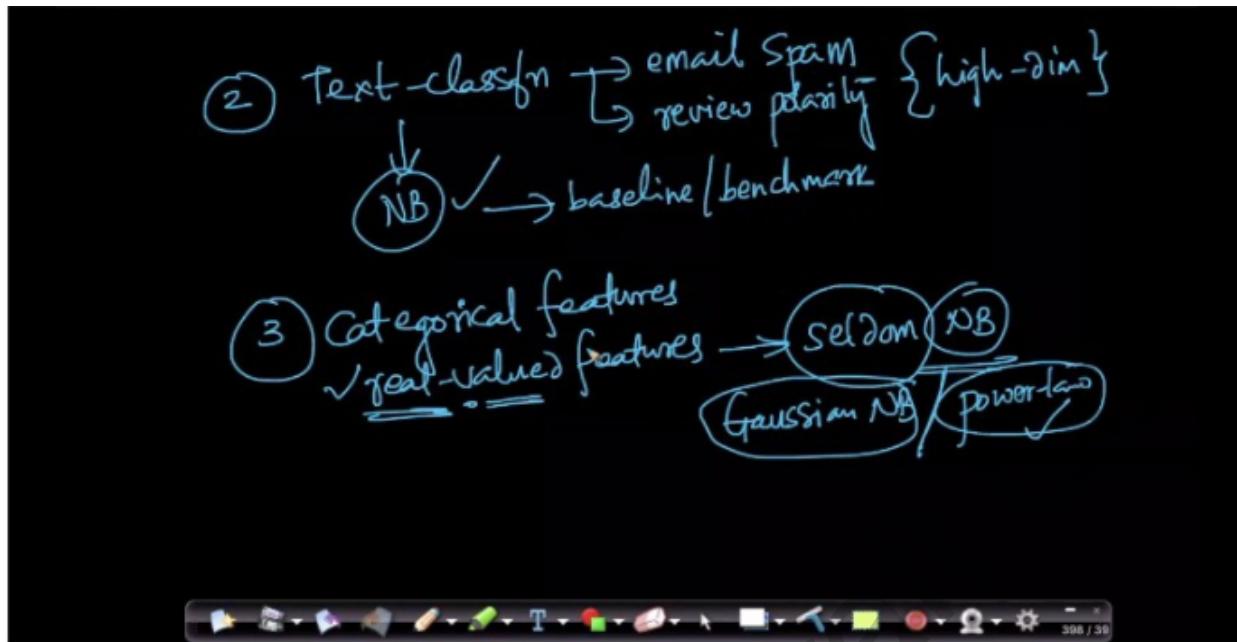
32.15 Best and worst cases



Timestamp 2:50

Now let's look at the cases where NB performs best and worst.

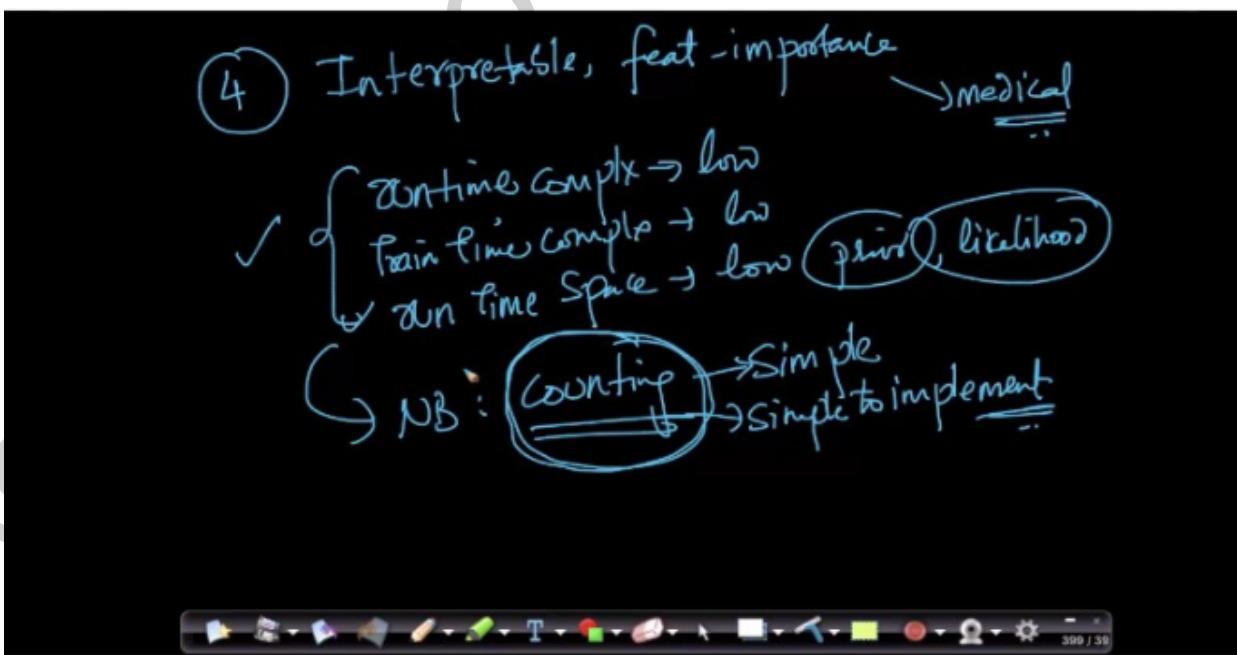
1. If the features are conditionally independent (the fundamental assumption of NB), then our model performs very well. Even if some of the features are slightly correlated then also our model performs reasonably well. Consider a positive case where words like good, great, phenomenal exist, even though these words are related our model will perform reasonably well in this situation. There are some mathematical proofs justifying this.



Timestamp 4:37

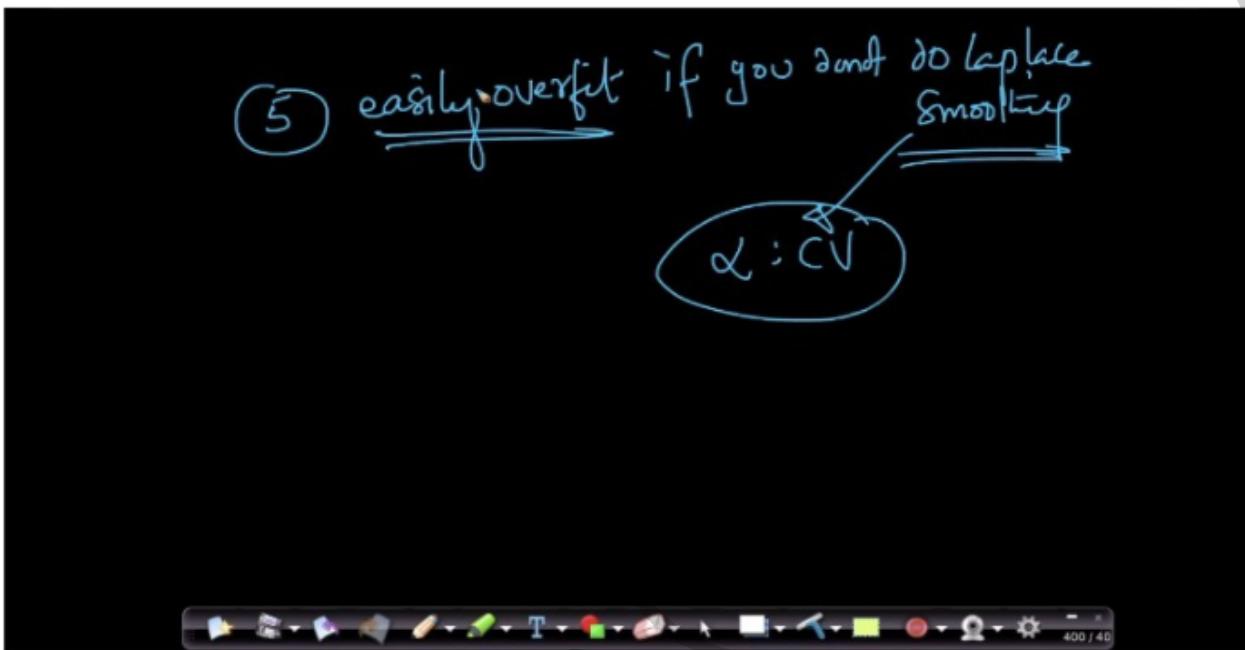
2. NB is extremely useful in dealing with high dimensional data as we saw during text classification. It serves as a good benchmark for problems like spam detection, polarity of a review, etc.

3. NB performs very well with categorical features. For real valued features we can use Gaussian NB. However it is seldom used at least in internet domains.



Timestamp 6:27

4. NB is highly interpretable. Feature importance can also be easily calculated. This makes it suitable for problems where interpretability is important like medical purposes. Also the time and space complexities for NB is very low when compared to KNN, which makes it useful in low latency solutions. It is also very easy to implement.



Timestamp 7:00

5. One of the disadvantages of NB is that it can easily overfit. In order to tackle this we need to do laplace smoothing. We can find the correct value of α using cross-validation.

32.16 Code example

We can see the scikit learn's implementation of Naive Bayes here =>

https://scikit-learn.org/stable/modules/naive_bayes.html

AppliedRoots (Draft Copy)