

Naive Bayes

①

- It is a classification Algorithm.
- It is a probability based Technique.
- k-nnth neighborhood based Technique.

Conditional probability

$$P(A|B) = P_r(A=a \mid B=b)$$

Variable variable.
↑ ↑
A: Random Variable
B: Random Variable
Value Value

Two dies $D_1 + D_2$

	D_2					
D_1	1	2	3	4	5	6
1	2	3	4	5	6	7
2	3	4	5	6	7	8
3	4	5	6	7	8	9
4	5	6	7	8	9	10
5	6	7	8	9	10	11
6	7	8	9	10	11	12

→ what is the probability that $D_1 + D_2 \leq 5$?

$$P(D_1 + D_2 \leq 5) = 10/36$$

→ what is the probability that $D_1 = 2$, given $D_1 + D_2 \leq 5$?

$$P(D_1 = 2 \mid D_1 + D_2 \leq 5) = 3/10$$

Condition on

Let us assume

$D_1 = 2$ (is an Event A)

$D_1 + D_2 \leq 5$ (is an Event B)

$P(A|B) = \frac{P(A \cap B)}{P(B)}$ (definition of Conditional probability)

$P(A \cap B)$ = Probability of Event A has happened and also Event B also has happened.

∴ This is Valid $P(B) \neq 0$

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{3/36}{10/36} = 3/10$$

Reference

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

22.2 Independent vs mutually exclusive events

Independent Events

A, B are said to be independent if $P(A|B) = P(A)$

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

A: getting Value of 6 in die 1's throw ($D_1=6$)

B: getting Value of 3 in die 2's throw ($D_2=3$)

$$P(D_1=6 | D_2=3) = P(D_1=6)$$

Here throwing D_1 and D_2 are two independent events. There is no point that when $D_2=3$, D_1 cannot not get 6

$$P(D_2=3 | D_1=6) = P(D=3)$$

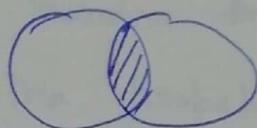
mutually Exclusive Events

If $P(A|B) = P(B|A) = 0$ Then $A \cup B$ are said to be mutually Exclusive Events

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

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

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



$$\text{Event A : } D_1 = 6$$

$$\text{Event B : } D_1 = 3$$

$$P(A|B) = P(D_1=6|D_1=3)$$

As we already got 3 probability of getting 6 is zero

$$\textcircled{1} \quad P(A|B) \leftarrow \frac{P(A \cap B)}{P(B)} = \frac{P(D_1=6 \cap D_1=3)}{P(D_1=3)} = \frac{0}{0.3} = 0$$

$$\textcircled{2} \quad P(A|B) \leftarrow \frac{P(A \cap B)}{P(B)} = \frac{P(D_1=6 \cap D_1=6)}{P(D_1=6)} = \frac{0.1}{0.6} = 0.1667$$

$$\textcircled{3} \quad P(A|B) \leftarrow \frac{P(A \cap B)}{P(B)} = \frac{P(D_1=6 \cap D_1=1)}{P(D_1=1)} = \frac{0.1}{0.4} = 0.25$$

$$\textcircled{4} \quad P(A|B) \leftarrow \frac{P(A \cap B)}{P(B)} = \frac{P(D_1=6 \cap D_1=2)}{P(D_1=2)} = \frac{0.1}{0.2} = 0.5$$

Bayes Theorem

→ Simple, elegant, beautiful, useful theorem

→ Early 1700's

Thm:

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

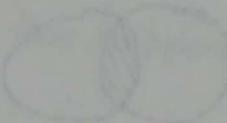
Posterior prob

Joint likelihood

Probability

Evidence

class.



Proof

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(A, B)}{P(B)} \stackrel{(1)}{=} \frac{P(A \cup B)}{P(B)}$$

$A \cap B = B \cap A \rightarrow$ Set Theory

$$P(A|B) = \frac{P(B \cap A)}{P(B)} = \frac{P(B, A)}{P(B)} \rightarrow \text{Eqn } (1)$$

$$P(B|A) = \frac{P(B \cap A)}{P(A)} \Leftarrow \text{By defn.}$$

$$P(A) \cdot P(B|A) = P(B \cap A) \rightarrow \text{Eqn } (2)$$

Replacing Eqn 2 with Eqn 1

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

(5)

Example

we have three machines and produce an output

$m_1 = 0.2$, $m_2 = 0.3$, $m_3 = 0.5$ and the fraction of defective items produced is this

$d_1 = 0.05$, $d_2 = 0.03$, $d_3 = 0.01$

if an item is chosen at random from the total o/p and is found to be defective what is the probability that is produced by m_3

B : Item is defective.

$P(B|A_1)$: Probability of an item is defective given it is made by m_1

$P(B|A_2)$: " made by m_2

$P(B|A_3)$: " by m_3

$$P(B|A_1) = 0.05 \quad P(A_1) = 0.2$$

$$P(B|A_2) = 0.03 \quad P(A_2) = 0.3$$

$$P(B|A_3) = 0.01 \quad P(A_3) = 0.5$$

$$P(B) = P(B|A_1)P(A_1) + P(B|A_2)P(A_2) + P(B|A_3)P(A_3)$$

$$= (0.05)(0.2) + (0.03)(0.3) + (0.01)(0.5)$$

$$P(B) = 0.024$$

$P(B) = 2.4\%$ of all the products are defective.

$$P(A_3|B) = \frac{P(B|A_3)P(A_3)}{P(B)} = \frac{0.01 \times 0.5}{0.024} = 5/24$$

↓
What is the probability, given an item is defective what is probability that it is made

1. A total of 46% of the voters in a certain city classify themselves as independent, whereas 30% classify themselves as liberals & 24% say they are conservative.

In a recent local election 35% of the independents, 62% of the liberals and 58% of the conservatives voted.

If a voter is chosen at random. Given that this person is in the local election, what probability that he or she is an independent is _____

$P(J)$: Probability of person being an independent voter

$$P(J) = 46\%$$

$$P(L) = 30\%$$

$$P(C) = 24\%$$

\Rightarrow And also given that 35% of independents are voters

$$P(V|J) = 0.35, P(V|L) = 0.62, P(V|C) = 0.58$$

$$P(J|V) = \frac{P(V|J) P(J)}{P(V)}$$

$$= \frac{0.35 \times 0.46}{P(V|J) P(J) + P(V|L) P(L) + P(V|C) P(C)}$$

$$P(V|J) P(J) + P(V|L) P(L) + P(V|C) P(C)$$

$$P(J|V) = 0.33$$

Naïve Bayes Algorithm

(7)

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

→ Abstractly, Naïve Bayes is a Conditional probability

~~Naïve~~ means simplistic or unsophisticated

I have a datapoint x with n features

$$x = (x_1, x_2, \dots, x_n)$$

If have a K classes Each datapoint represent a K

$$x \rightarrow 1, 2, 3, \dots, K$$

$$K \rightarrow C_K$$

$P(C_K|x)$ = probability of class label given a point x ,

Random variable has n features and hence n random variable.

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

$$\text{posterior} = \frac{\text{prior} \times \text{likelihood}}{\text{Evidence.}}$$

Numerator

$$P(C_K) P(x|C_K) = P(x \cap C_K)$$

I have calculated

$$P(C_1|x), P(C_2|x), \boxed{P(C_3|x)}, \dots, P(C_K|x)$$

last arrow with the largest

$P(x)$

→ The denominator is same of all conditions, hence we can ignore them for a while.

→ Let's focus only on numerators.

$$P(C_k|x) = P(C_k) \quad P(x|C_k) = P(C_k \cap x) \propto P(C_k, x)$$

$$P(C_k|x) \propto P(C_k, x)$$

proportional
to

The numerator is called Joint probability model, C_k & x are jointly happening.

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

→ There is something called as chain rule for conditional probability.

$$A \cap B = B \cap A$$

$$P(C_k, x_1, \dots, x_n) = P(x_1, \dots, x_n, C_k) \quad \text{Conditional probability, not bayes theorem}$$

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

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

= ...

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

$$P(x_n | C_k) P(C_k)$$

We want to compute $P(C_k|x)$ for each column and this is equal to $P(C_k|x)$

Let take one first Condition.

class label

obere
fit
lears

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

$$P(w=180 | \text{haircolor} = b, \text{hairlength} = 5\text{cm}, \text{eyecolor} = br, c_k = \text{abies})$$

This is very hard find these in the real data.

what Naive Bayes does make Naive assumption of Conditional Independence.

$$P(A|B) = P(A) \quad \underbrace{\text{As } A \text{ and } B \text{ are two independent Events}}_{\text{Simple independence}}$$

$$P(A|B, c) = P(A|c) \quad \begin{matrix} \text{Conditionally} \\ \text{i.e., } A \text{ is independent of } B \text{ given } c \\ \text{or} \\ B \cap c \end{matrix}$$

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

x_i is independent of $x_{i+1}, x_{i+2}, \dots, x_n$ given c_k

Naive + Conditionally independence of Each feature pair

$X = (x_1, x_2, \dots, x_n)$ are features

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

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

$$\propto p(x_1 | c_k) p(x_2 | c_k) p(x_3 | c_k) \dots p(x_n | c_k) p(c_k)$$

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

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 depending 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

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 (MAP) decision rule.

The corresponding classifier, a Bayes classifier is the function Bay assigns a class label $\hat{y} = c_k$ for some k as follows:

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

Toy Example

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

	Outlook	Temperature	Humidity	Wind	Play
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

$$P(C \mid \underbrace{f_1, f_2, f_3, f_4}_{\text{X}}) \propto P(f_1 \mid C) * P(f_2 \mid C) * P(f_3 \mid C) * P(f_4 \mid C) * P(C)$$

(Likelihoods)

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

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

In the training phase

→ we done frequency counting

f ₁	P(outlook = 0 class play = Yes/No)	Frequency		Probability in class	
		Play = Yes	Play = No	Play = Yes	Play = No
Sunny		2	3	2/9	3/5
overcast		4	0	4/9	0/5
Rain		<u>3</u>	<u>2</u>	<u>3/9</u>	<u>2/5</u>
		<u>9</u>	<u>5</u>		

$$P(\text{outlook} = \text{sunny} \mid C=\text{Yes}) = \frac{2/14}{9/14} = 2/9$$

$$P(\text{outlook} = \text{sunny} \mid C=\text{No}) = \frac{3/14}{5/14} = 3/5$$

Temperature	Frequency		Probability in class	
Hot	2	2	2/9	2/5
mild	4	2	4/9	2/5
cool	<u>3</u>	<u>1</u>	<u>3/9</u>	<u>1/5</u>
	<u>9</u>	<u>5</u>		

Humidity = Play = yes Play = no play = yes play = no

High 3 4 3/9 4/5

Normal 6 1 6/9 1/5

Wind

Strong 3 3 3/9 3/5

Weak 6 2 6/9 2/5

we also calculate $P(\text{class play} = \text{yes})$ & $P(\text{class play} = \text{no})$

$P(\text{class play} = \text{yes})$ & $P(\text{class play} = \text{no})$	Frequency	Probability in class
Play	Play = yes Play = no	Play = yes play = no

9 5 $9/14$ $5/14$

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 condition specified by x')

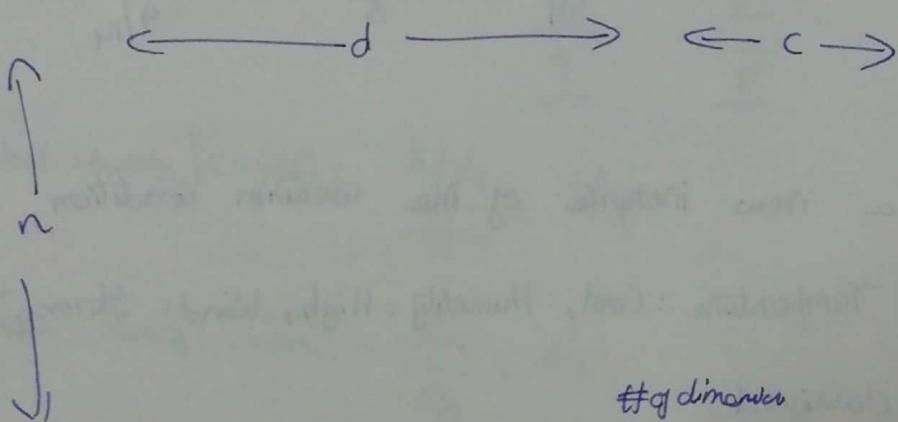
With MAP rule, we compute the posterior probabilities. This is easily done by looking up the tables we built in the learning phase

$$\begin{aligned}
 P(\text{class play}=\text{yes} | x) &\propto [P(\text{Sunny} | \text{play}=\text{yes}) * P(\text{Cool} | \text{play}=\text{yes}) * \\
 &\quad P(\text{High} | \text{play}=\text{yes}) * P(\text{Strong} | \text{class play}=\text{yes})] * \\
 &\quad P(\text{class play}=\text{yes}) \\
 &= \frac{2}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{9}{14} \\
 &= 0.0053
 \end{aligned}$$

$$\begin{aligned}
 P(\text{class play}=\text{no} | x) &\propto [P(\text{Sunny} | \text{no}) * P(\text{Cool} | \text{no}) * \\
 &\quad P(\text{High} | \text{no}) * P(\text{Strong} | \text{no})] * P(\text{no}) \\
 &= \frac{3}{5} \times \frac{1}{5} \times \frac{4}{5} \times \frac{3}{5} \times \frac{5}{14} \\
 &= 0.0255
 \end{aligned}$$

→ In training phase we will calculate all the likelihood probabilities, and the class probabilities

→ $P(\text{class})$



Time taken to calculate is $O(d * c)$

of attrs # of classes.

$(d) \quad (c)$
→ $P(f_i | c)$

The space complexity is $O(d * c)$

For Training

→ Time Complexity for Naive Bayes is $O(n)$
 Space " " " " $O(d * c)$ } d, c should be much smaller than n
 \hookrightarrow dictionary

$$P(\text{Yes} | x') \approx 0.0053$$

$$P(\text{No} | x') \approx 0.0205$$

Since n is large $x' : y' = \text{No}$

~~BB~~

For Test :-

If i have d dimensions, we have to do d -lookups.

→ Time Complexity $O(d * c)$

If d is small & c is small, the training time is less, space is also less.

→ Compared to KNN, NB is much better in terms of space complexity @ runtime.

$$O(d * c) \rightarrow \text{NB}$$

$$O(n * d) \rightarrow \text{KNN}$$

→ NB is much more memory efficient at runtime

NB Bernoulli \rightarrow Binary classification

NB Gaussian \rightarrow 1-5 classes

NB multinomial \rightarrow More than 5 classes

Naive Bayes on Text Data

7

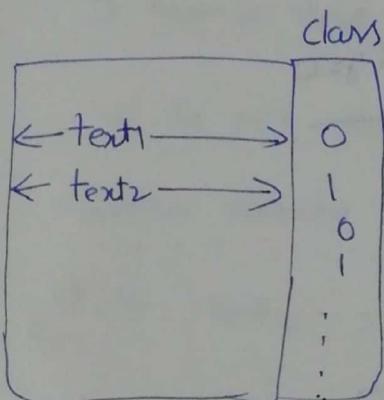
→ Navies Bayes is first method to create a Spam filter

⑧ Email → SPAM
 ↳ not spam.

Review \rightarrow +ve
 \nwarrow -ve. } Polarity

→ Textclassification → NB a very simple and heavily applied method.

→ NB can be applied on many languages.



Task :

$$P(Y_{qj}=1) + \epsilon_{2qj})$$

$$P(Y_{q_j} = 0 | \text{text}_{q_j})$$

text → Stopwords

stemming

n-grams

n-grams (Most ppl apply)

$\hookrightarrow w_1, w_2, w_3 \dots w_d \}$ Binary Bow

↳ whether a word is present or not

we are not counting 'number of occurrences'

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

$$P(y=1 | \text{text}) \propto P(y=1 | \omega_1, \omega_2, \omega_3 \dots \omega_d)$$

$$\alpha \underbrace{P(y=1) * P(w_1 | y=1) * P(w_2 | y=1) * \dots * P(w_d | y=1)}_{\text{Likelihood}}$$

class prior.

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

$$P(y=0 \mid \text{text}) \propto P(y=0) \cdot \prod_{i=1}^d P(w_i \mid y=0)$$

$$P(y=1) = \frac{\# \text{ of Train points with } y=1}{\text{Total } \# \text{ of Train points}}$$

$$P(y=0) = \frac{\# \text{ of Training points with } y=0}{\text{Total } \# \text{ of train points}}$$

$$P(w_i \mid y=1) = \frac{\# \text{ of data pts which contain } w_i \text{ & } y=1}{\# \text{ data points with } y=1}$$

$$P(w_i \mid y=0) = \frac{\# \text{ of data pts which contain } w_i \text{ & } y=0}{\# \text{ data points with } y=0}$$

→ NB is a very good baseline
benchmark

LAPLACE / ADDITIVE Smoothing

(19)

Training \div $P(y=1)$; $P(y=0)$ \leftarrow class priors.

$$\begin{array}{ll}
 P(w_1 | y=1) & P(w_1 | y=0) \\
 P(w_2 | y=1) & P(w_2 | y=0) \\
 \vdots & \\
 P(w_m | y=1) & P(w_m | y=0)
 \end{array}
 \quad \left. \quad \right\} \text{Likelihoods.}$$

Test \div

I have a text $\text{text}_q = (w_1, w_2, w_3, w_4)$

w_4 is not present in the set of words $\{w_1, w_2, w_3, \dots, w_m\}$ training corpus

\rightarrow Let us assume w_4 is a crazy word which is not often used

$$P(y=1 | \text{text}_q) = P(y=1 | w_1, w_2, w_3, w_4)$$

$$= P(y=1) * P(w_1 | y=1) * P(w_2 | y=1) * P(w_3 | y=1) * P(w_4 | y=1)$$

Now the challenge is how do we get

$$P(w_4 | y=1) \text{ or } P(w_4 | y=0)$$

Ignoring w_4 or dropping w_4 is equal to $P(w_4 | y=1) = 1$ &
 $P(w_4 | y=1) = 0$. This is incorrect.

$$\begin{aligned}
 P(w_4 | y=1) &= \frac{P(w_4, y=1)}{P(y=1)} \quad \# \text{pts where } w_4 \text{ occurs in } y=1 \\
 &= \frac{\# \text{pts where } w_4 \text{ occurs in } y=1}{\# \text{pts where } y=1} \\
 &= \frac{0}{\# \text{pts}} = 0
 \end{aligned}$$

If I put, it will be dangerous. The whole probability will be zero.

Work around for this problem is Laplace Smoothing/Additive Smoothing. (20)

Laplace Smoothing is different from Laplacian Smoothing

↓
This will be used in Image processing.

$$P(w^i | y=1) = \frac{0+\alpha}{n_i + \alpha K}$$

K : # of distinct values that w^i can take

w^i can be present & not present. So it can take 0 & 1
In text processing case.

In non text processing

$$P(\text{feature} = a | y=1) = \frac{0+\alpha}{n_i + \alpha}$$

Here K = # distinct values w^i can take.

α = Any numeric value

Typically we will take $\alpha = 1$

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

Case 1: $\alpha = 1 \approx 1/102 \neq 0$

$\approx P(w^i | y=1) \neq 0$

Why we are doing this entire process why can we replace with a simple Epsilon which is $\epsilon = 0.001$

Case 2: $\alpha = 10,000 = \frac{10,000}{100 + 20,000} = \frac{10,000}{20,100} \approx \frac{1}{2}$

when your alpha is large, we are saying Two probabilities are

Equal $w^i = 0$

$w^i = 0$

$$P(w^i | y=1) = P(w^i | y=0) = \frac{1}{2}$$

The probability of w_i occur for positive class and the probability of w_i occurs for negative class and these both probabilities are equal to half

→ In Laplace Smoothing

$$P(w_i | y=1) = \frac{\# \text{ of datapoints with } w_i \text{ & } y=1}{\# \text{ of datapts with } y=1 + \alpha K}$$

Present in my
Training data

This is called additive smoothing as we are adding αK

$$\rightarrow P(w_i | y=1) = \frac{2/50}{50 + \alpha K} = \frac{2}{50} \text{ when } \alpha = 0$$

$$\alpha = 1 \quad \frac{2+2}{50+4} = \frac{3}{54}$$

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

$$\alpha = 10 \quad \frac{2+10}{50+20} = \frac{12}{70}$$

$$\alpha = 100 \quad \frac{2+100}{50+200} = \frac{102}{250}$$

As alpha is increasing ~~ratio~~ $\alpha/p_1 < \alpha/p_2 < \alpha/p_3 < \alpha/p_4$

$\alpha \uparrow$ we are moving likelihood probabilities to uniform distribution

If ~~numerators~~ 3) denominator is small we have less confidence in ratio

Often times people apply $\alpha=1$ which is add-one-smoothing

We are smoothing it towards uniform smoothing $\left(\frac{1002}{2050}\right) \approx 1/2$

α is directly proportional to bias & variance.

Log Probabilities & Numerical Stability

$$P(y=1 | w_1, w_2, \dots, w_d) = \underbrace{P(y=1)}_{0 \leq 1} * \underbrace{P(w_1 | y=1)}_{0 \leq 1} * \underbrace{P(w_2 | y=1)}_{0 \leq 1} * \dots * \underbrace{P(w_d | y=1)}_{0 \leq 1}$$

Suppose d is large, let $d=100$

we are multiplying 100 numbers which lie b/w 0 & 1

$$(0.2 * 0.1 * 0.2 * 0.1) = 0.0004 \text{ (It is a very small decimal)}$$

If I multiply 100 numbers all of which lies b/w 0 & 1 we can run into numerical stability

In python if we take double precision floating point #. it has 16 significant values.

when we multiplying lot of small number we will end up at Numeric Underflow

Then python starts doing rounding and result in Errors

$$\log(0.0004) \approx -2.23979$$

Log probabilities :

we initially have.

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

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

log don't change the order $x \uparrow ; \log(x) \uparrow$

↳ monotonic function.

~~$\log(P(y=1 | w_1, w_2, \dots, w_d))$~~

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

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

log converts mult into addition $\log(a \cdot b) = \log(a) + \log(b)$

log converts Exponentiation into multiple $\log(a^b) = b \cdot \log(a)$

Bias and Variance Trade off

(25)

Naive Bayes : high Bias \rightarrow underfitting
high variance \rightarrow overfitting.

- in Laplace smoothing defines underfitting & overfitting
- changes either high bias & high variance.

Case 1 $\alpha = 0$

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

Let's assume a word occurs two times $= \frac{2}{1000}$

$$n = 2000 \text{ words} \rightarrow 1000 \text{ true} \\ \rightarrow 1000 \text{ false}$$

Even for words that are rare we are giving no probability $P(w_i | y=1)$
result in overfitting. (high variance)

\rightarrow Small Training data changes result in ^{large} _{Model change} } \downarrow ^{large} _{Var of variance}

Since w_i occurs only in 2 out of 2000 cases
1000 true 1000 false.

Small change in my D_{train}

\hookrightarrow remove the 2 texts (x_i^1) that contain w_i

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

As probability changed the model changes.

when $\alpha = 0$ small change in D_{train} results in large change in the
model \rightarrow high variance \rightarrow overfitting

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

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

0.81

when α become very large
for all w_i

$$P(w_i | y=1) \approx P(w_i | y=0) \approx \frac{1}{2}$$

Model cannot distinguish whether w_i is 0 or 1

α = v-large

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

$\approx \frac{1}{2}$

→ How to find the right alpha

→ known: right $\alpha \rightarrow$ using simple CV or 10 fold CV

→ Same ~~is~~ with α also.

α in NB \rightarrow hyper Parameters
 κ in KNB \rightarrow

→ Right ^{hyper} parameters in most model are decided using Cross Validation

Feature importance and interpretability

(27)

How to get Feature importance:

Likelihood

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

$w_i, P(w_i | y=0)$

① Sort w_i s based on $P(w_i | y=1)$ in desc.

Top words occur very often

w_i s which has high value of probabilities & important words & features in determining that pair of the class.

w_i	$P(w_i y=1)$	$P(w_i y=0)$
w_1		
w_2		
w_3		
w_4		
\vdots		
w_n		

→ class: Find words/feature (w_i) with highest Value of $P(w_i | y=1)$

→ vector = " " " of $P(w_i | y=0)$

In NB feature importance determined/obtained directly from the model.

which is not possible in KNN

Interpretability

Suppose a query text x_q is given $\rightarrow y_q = 1$

$\{w_1, w_2, w_3, \dots, w_d\}$

From concluding $y_q = 1$ because x_q contains words w_3, w_6, w_{10} which have a high value of $P(w_3 | y=1), P(w_{10} | y=1), P(w_6 | y=1)$

Similarly

I am not saying $y_q \neq 0$; w_9, w_{15}, w_{20}

$P(w_8 | y=0)$ is very small

$P(w_{16} | y=0)$ is very small.

IMBALANCED DATA

29

What happens to NB with imbalanced data.

$$n \rightarrow n_1 + n_2$$
$$n \rightarrow n_1 - n_2$$

$$n_1 > n_2$$
$$n_2 > n_1$$

Basic formula for NB

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

$$P(y=0 | w_1, w_2, w_3, \dots, w_d) = P(y=0) \prod_{i=1}^d P(w_i | y=0)$$

Compare

Let's assume 90% of pts are +ve pts

10% of pts are -ve pts

$$P(y=1) = n_1/n = 0.9$$

$$P(y=0) = n_2/n = 0.1$$

When

$$\prod_{i=1}^d P(w_i | y=1) \text{ exactly equal to } \prod_{i=1}^d P(w_i | y=0)$$

The first one has ~~no~~ advantage to be a positive class as it is 90%.

when we have an imbalanced dataset, because of class purity \rightarrow

majority/dominating class has an advantage.

SOL ① up sampling

② down sampling

$$n_1 \approx n_2 \quad P(y=1) = P(y=0) = 1/2$$

② Simply drop probability terms. $P(y=1) \approx P(y=0)$

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

③ modified NB to account for class imbalances
 \hookrightarrow not often used

Imbalance data

+ve $n_1 = 900$

-ve $n_2 = 100$

$$\rightarrow \text{majority class}$$

$$P(w_i | y=1) = \frac{900}{900} \rightarrow 0 \text{ to } 900$$

$$P(w_i | y=0) = \frac{100}{100} \rightarrow 0 \text{ to } 100$$

The minority class tend to have a small numerators when compared to majority class.

Laplace Smoothing

use the same alpha for positive & negative class

& impact more for minority class & less for majority class

Ex: minority \rightarrow majority

$$P(w_i | y=0) = \frac{2}{100}, \frac{18}{900} \leftarrow P(w_i | y=1)$$

$$\rightarrow \frac{2}{20}, \frac{1}{10} \rightarrow$$

Let's assume $\alpha = 100$

$$\frac{2+10}{100+20}, \frac{18+10}{900+20} = \frac{12}{120}$$

$$\frac{12}{120} = 28/920 = 0.18104 \rightarrow 3.04\%$$

\rightarrow got change to 10%.

Now a Hack

We can create a ~~one~~ α for the class & α for the class.

\rightarrow RNN is impacted by imbalance data

Outliers

How are outliers are handled in NB?

Text classificat \Rightarrow ~~w_q~~ \leftarrow outlier at test time.
 $w_q = w_1, w_2, w_3, w_4$

$w \notin \{w_1, w_2, w_3, \dots, w_m\} \leftarrow$ set of words that I have
not present seen in Dtrain.

Laplace smoothing

$$P(w_i | y=1) = \frac{0}{n_1} + \alpha$$

When we get a outlier at test time it is handled by α

Outlier in training data:

$\{w_1, w_2, w_3, \dots, w_m\} \leftarrow$ set of words in Dtrain.

Let us consider w_8 occur very few times in the Dtrain.

→ To handle outlier in Training

① if a word (w_j) occurs fewer than 10 times, then just ignore that word.

$w_j \notin \{w_1, w_2, w_3, \dots, w_m\}$
occurs very few times in Dtrain.

② Laplace smoothing (α)

Missing Values

Case 1

If we have a text data : Text : $\{w_1, w_2, \dots, w_d\}$
 ↑ no care of missing data

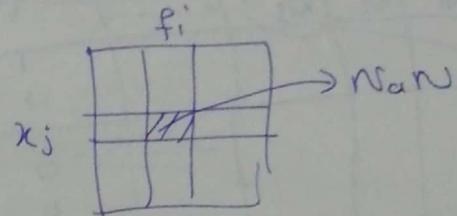
Case 2

Categorical features

$$f_i \in \{a_1, a_2, a_3\}$$

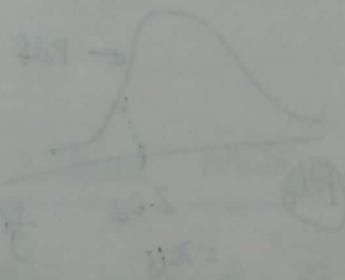
we can consider NaN is a category

$$f_i \in \{a_1, a_2, a_3, \text{NaN}\}$$



Case 3 Numerical features

↳ Standard feature Imputation (mean, median, knn)



NB with Numerical features

f_1	f_2	\dots	f_j	f_d	class label (y_i)
x_{i1}	x_{i2}	\dots	x_{ij}	\dots	1
2.62					1
					1
					0
					0
					0
					0
					0

n_1
+ve

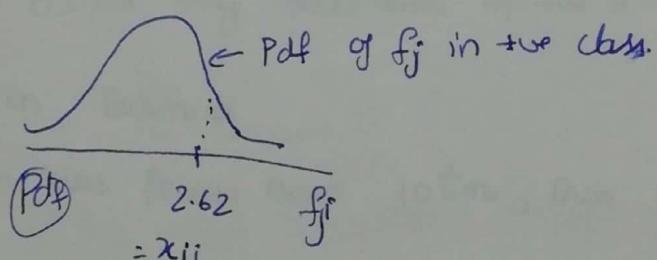
n_2
-ve

$$P(y_i=1 | x_{i1}, x_{i2}, \dots, x_{id}) \propto P(y_i=1) \prod_{j=1}^d P(x_{ij} | y_i=1)$$

class prior $\frac{n_1}{n_1+n_2}$

$$\frac{P(x_{ij} | y=1)}{P(x_{ij} | y=0)}$$

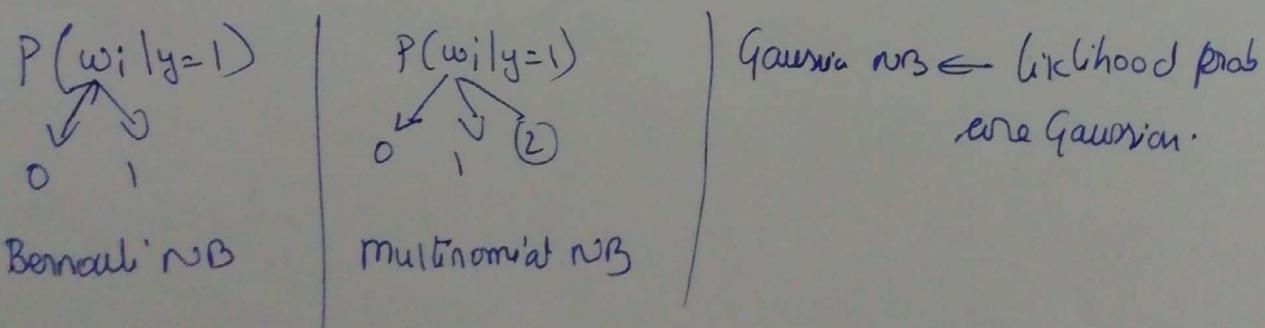
we can obtain this probability from the
Pdf of f_j in 0



Assumption

→ Ppl assume f_j in $\mathbb{R}_{\text{non-neg}}$ will come from Gaussian distribution $N(\mu_j, \sigma_j^2)$

Such a model is called Gaussian Naive Bayes



NB: Conditional Independencies

It is assuming that each of your features are independent of others given the class label.

MULTI CLASS CLASSIFICATION

NB fundamentally can do multiclass classification.

$$\left. \begin{array}{l} P(y_i=1 | w_1, w_2, w_3 \dots w_d) \\ P(y_i=0 | w_1, w_2, \dots, w_d) \\ P(y_i=2 | w_1, w_2, \dots, w_d) \\ \vdots \\ P(y=c-1 | w_1, w_2, \dots, w_d) \end{array} \right\} c\text{-class classification.}$$

Similarity & Distance Matrix

Can naive bayes classification instead of giving points explicit, we are given a distance/similarity metrics.

→ KNN can do.

→ NB cannot use distance/similarity matrix.

Large Dimensionality

→ Can NB handle large dimensionality

↳ we have to use log probabilities.

Best & Worst Cases

① Conditional Independence of features

↳ True: NB perform very well.

↳ False: NB deteriorates/ degrades

When some features are dependent NB works reasonable well.

True → good
great
phenomenal } These words are independent

② Text classification → Email spam

↳ review polarity } high dimension

③ NB is Extensively used when Categorical features/binary features

Typically when we have real valued features → NB is not used much.

④ NB is super ~~interpretable~~ interpretable

↳ we can get feature importance

↳ runtime complexity is low

↳ Train Time complexity is low

↳ run time Space is low (Big) (Likelihood)

⑤ Easily overfit if you do not Laplace smoothing and we have to

find ^{with} _{cross} validation

CFDE (Example)

Points

→ alpha is typically taken to be valued b/w 0 & 1. In some cases where we want to regularize heavily as we may be overfitting a lot, we could have alpha > 1 also.

How to find feature importance?

```
from sklearn.naive_bayes import BernoulliNB
```

```
bnb = BernoulliNB()
```

```
bnb.fit(X_train, y_train)
```

```
neg-class-prob-sorted = bnb.features_log_prob_[0, :].argsort()
```

```
pos-class-prob-sorted = bnb.features_log_prob_[1, :].argsort()
```

```
point(np.take(count_vect.get_feature_names(), neg-class-prob-sorted[:10]))
```

Gaussian NB

```
from sklearn import datasets
iris = datasets.load_iris()
from sklearn.naive_bayes import GaussianNB
gnb = GaussianNB()
y_pred = gnb.fit(iris.data, iris.target).  
predict(iris.data)
```

```
print("number of mislabeled points out of a  
total %d points : %d" % (iris.data.shape[0],  
(iris.target != y_pred).sum()))
```

Attributes

class_prior : Probability of Each class.

class_count : Number of training samples in Each class

theta : Mean of Each feature per class

sigma : Variance of Each feature per class.

Methods

fit (X, y[, sample_weight])

get_params

partial_fit

predict

predict_log_proba

predict_proba

score

set_params

multinomial NB
Bernoulli NB
out of core NB

BernoulliNB

`BernoulliNB(alpha=1.0, binarize=0.0, fit_prior=True, class_prior=None)`

→ Naive Bayes classifier for multivariate Bernoulli models

→ Like multinomialNB, this classifier is suitable for discrete data. The difference is that while multinomialNB works with occurrence counts,

→ BernoulliNB is designed for binary / boolean features

Parameters

alpha : float, optional (default=1.0)

Additive (Laplace/Lidstone) smoothing parameter (0 for no smoothing)

binarize : float, None, optional, default=0.0

Threshold for binarizing (mapping to booleans) of sample features. If none input is presumed to already consist of binary vectors

fit_prior : boolean, optional (default=True)

whether to learn class prior probabilities or not. If false, a uniform prior will be used

class_prior : array-like, size=[n_classes], optional (default=None)

Prior probabilities of the classes. If specified the priors are not adjusted according to the data.

Attributes

class_log_prior_ : array, shape = [n_classes]

Log probabilities of each class.

feature_log_prob:

Empirical log probability of feature given a class

class_count:

feature_count:

Number of samples encountered for each (class, feature) during fitting.

This value is weighted by the sample weight when provided.