

We need to find probability of a class label given X

X is n -d

variable. Using Bayes' theorem, the conditional probability can be derived

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

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

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

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:

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

whichever class label's probability is the highest, we'll pick that one

denominator is constant, we can ignore it

$$p(C_k | \mathbf{x}) = \frac{p(C_k) p(\mathbf{x} | C_k)}{p(\mathbf{x})} \rightarrow p(C_k \cap \mathbf{x}) = \underline{\underline{p(C_k, \mathbf{x})}}$$

the above is called a joint probability

$$\begin{aligned} p(\underbrace{x_1}_{A} | \underbrace{x_2, \dots, x_n, C_k}_{B}) &= p(x_1 | x_2, \dots, x_n, C_k) \\ &\quad * p(x_2, \dots, x_n, C_k) \\ \downarrow p(A, B) &= p(A | B) p(B) \\ &= \end{aligned}$$

Bayesthm

which can be rewritten as follows, using the [chain rule](#) for repeated applications of the definition of [conditional probability](#):

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

Thus, the joint model can be expressed as

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

$$p(x_2, x_3, \dots, x_n, c_k)$$

----->

which is
proportional to
what we just
computed

Contents - Google Docs

Naive Bayes classifier - Wikipedia

ShatterLine Blog » Not-so-Naive

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

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:

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

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

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

Thus, the joint model can be expressed as

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

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

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

we assume conditional independence

for ex if we assume conditional independence among A and B

means that

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

$$\dots \{ p(A|B,C) = p(A|C) \}$$

B and C

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

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

$$p(x_i | x_{i+1}, \dots, x_n, C_k) = p(x_i | C_k) \Rightarrow x_i \text{ is indep of } x_{i+1}, x_{i+2}, \dots, x_n \text{ given } C_k$$

Thus, the joint model can be expressed as

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

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

Thus, the joint model can be expressed as

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

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

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

where the evidence $Z = p(\mathbf{x}) = \sum_k p(C_k) p(\mathbf{x} | C_k)$ 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:

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



Naive Bayes

	Predictors	Response
	Outlook f_1	Class $y_i = Y \text{ or } N$
Day1	Sunny	No
Day2	Sunny	No
Day3	Overcast	Yes
Day4	Rain	Yes
Day5	Rain	Yes
Day6	Rain	No
Day7	Overcast	Yes
Day8	Sunny	No
Day9	Sunny	Yes
Day10	Rain	Yes
Day11	Sunny	Yes
Day12	Overcast	Yes
Day13	Overcast	Yes
Day14	Rain	No

✓ binary classifcn.
✓ f_1, f_2, f_3, f_4 : {categorical feat}

The Learning Phase

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

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

$P(\text{Temperature}=t | \text{Class}_{\text{play}}=b)$, where $t \in [\text{Hot, Mild, Cool}]$ and $b \in [\text{yes, no}]$.

$P(\text{Humidity}=h | \text{Class}_{\text{play}}=b)$, where $h \in [\text{High, Normal}]$ and $b \in [\text{yes, no}]$.

$P(\text{Wind}=w | \text{Class}_{\text{play}}=b)$, where $w \in [\text{Weak, Strong}]$ and $b \in [\text{yes, no}]$.

we compute stuff in learning phase

remember we need to compute this, (NB just stands for naive bayes)

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

NB

$$p(f_1 | C) \cdot p(f_2 | C) \cdot p(f_3 | C) \cdot p(f_4 | C) \cdot p(C)$$

Response
Class
Play=Yes
Play=No
No
No
Yes
Yes
Yes
No

$P(c)$ is easy to compute, just count how many yes or no are there

Temperature \in [Hot, Mild, Cool]
 Humidity \in [High, Normal]
 Windy \in [Weak, Strong]

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

Play \in [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

Handwritten notes on the right:

$$p(c | f_1, f_2, f_3, f_4)$$

$$p(f_1 | c)$$

$$* p(f_2 | c)$$

$$* p(f_3 | c)$$

$$* p(f_4 | c)$$

$$* p(c)$$

for each feature count this shit below

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

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

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

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

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

Handwritten note: $p(\text{outlook}=\text{Sunny} | c=\text{yes}) = 2/9$

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

P(Temperature=t Class _{play=Yes No})	Frequency		Probability in Class	
Temperature =	Play=Yes	Play=No	Play=Yes	Play=No
Hot				
Mild				
Cool				

similarly count for all features

		total= 9		total=5	
$P(\text{Temperature}=t \text{Class}_{\text{play=Yes No}})$		Frequency		Probability in Class	
	Temperature =	Play=Yes	Play=No	Play=Yes	Play=No
	Hot	2	2	2/9	2/5
	Mild	4	2	4/9	2/5
	Cool	3	1	3/9	1/5
	total= 9	total=5			

$P(\text{Humidity}=h \text{Class}_{\text{play=Yes No}})$		Frequency		Probability in Class	
	Humidity =	Play=Yes	Play=No	Play=Yes	Play=No
	High	3	4	3/9	4/5
	Normal	6	1	6/9	1/5
	total= 9	total=5			

$P(\text{Wind}=w \text{Class}_{\text{play=Yes No}})$		Frequency		Probability in Class	
	Wind =	Play=Yes	Play=No	Play=Yes	Play=No
	strong	3	3	3/9	3/5
	weak	6	2	6/9	2/5
	total= 9	total=5			

$P(\text{Wind}=w \text{Class}_{\text{play=Yes No}})$		Frequency		Probability in Class	
	Wind =	Play=Yes	Play=No	Play=Yes	Play=No
	strong	3	3	3/9	3/5
	weak	6	2	6/9	2/5
	total= 9	total=5			

We also calculate $P(\text{Class}_{\text{play=Yes}})$ and $P(\text{Class}_{\text{play=No}})$.

$P(\text{Class}_{\text{play=Yes}}) \& P(\text{Class}_{\text{play=No}})$		Frequency		Probability in Class	
	Play	Play=Yes	Play=No	Play=Yes	Play=No
		9	5	9/14	5/14
	total= 9	total=5			

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

loc

Handwritten notes:
 $p(C|f_1, f_2, f_3, f_4)$
 $= p(C) +$
 $p(f_1|C) +$
 $p(f_2|C) +$
 $p(f_3|C) +$
 $p(f_4|C)$

if we assume number of categories each feature takes is small, $o(1)$

then time complexity for training is $o(ndc)$

Contents - Google Docs x ShatterLine Blog » Not-so-Naive Naive Bayes classifier - Wikip... Chekuri Srikan...

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

	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
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 notes:
 Training ph: → likelihood prob → $p(\text{class})$ → $O(ndc)$

Diagram: A vertical arrow on the left is labeled 'n'. An arrow points from 'Temperature' to 'Humidity'.

The Learning Phase

In the training phase, we build a model by looking at the data and computing the probabilities for each class.

space complexity :

if there are C classes, order of C space

$d \times c$ space because , d features, c classes

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.

Handwritten: $x' = (\text{Sunny, Cool, High, Strong})$

$$P(\text{Class}_{\text{play}}=\text{Yes} | x') = \frac{P(\text{Sunny} | \text{Class}_{\text{play}}=\text{Yes}) \times P(\text{Cool} | \text{Class}_{\text{play}}=\text{Yes}) \times P(\text{High} | \text{Class}_{\text{play}}=\text{Yes}) \times P(\text{Strong} | \text{Class}_{\text{play}}=\text{Yes})}{P(\text{Class}_{\text{play}}=\text{Yes})}$$

Handwritten: $y_{\text{aj}} = ?$

$$= \frac{2/9 \times 3/9 \times 3/9 \times 3/9 \times 9/14}{9/14} = 0.0053$$

$$P(\text{Class}_{\text{play}}=\text{No} | x') = \frac{P(\text{Sunny} | \text{Class}_{\text{play}}=\text{No}) \times P(\text{Cool} | \text{Class}_{\text{play}}=\text{No}) \times P(\text{High} | \text{Class}_{\text{play}}=\text{No}) \times P(\text{Strong} | \text{Class}_{\text{play}}=\text{No})}{P(\text{Class}_{\text{play}}=\text{No})}$$

$$= \frac{3/5 \times 1/5 \times 4/5 \times 3/5 \times 5/14}{5/14} = 0.0205$$

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

Handwritten: $P(C=\text{Yes} | x')$

we'll store all data required for this in a dictionary

test complexity :

for c classes, we need to lookup each feature, time = $d * c$

if d and c is small, training time is v v small

space is also small as compared to k-nn