

Strength

- * simple and effective
- * makes no assumptions about the data distribution
- * fast training phase.

weakness

- * does not produce a model limiting the ability to understand how features are related to the class.
- * Requires selection of an appropriate k.
- * nominal features and missing data requires additional processing.
- * slow classification phase.

Probabilistic learning: Understanding Naive Bayes -
Conditional probability and Bayes Theorem,
Naive Bayes algorithm for classification,
the Laplace estimator, Using numeric features with
Naive Bayes.

① Understanding Naive Bayes

- Weather estimates are based on probabilistic methods or those concerned with describing uncertainty. They use data on past events to extrapolate future events.
- Mathematician Thomas Bayes, who developed foundational principles to describe the probability of events, and how probabilities should be revised in the light of additional information. These principles formed the foundation for what are now known as Bayesian methods.
- Classifiers based on Bayesian methods utilize training data to calculate an observed probability of each outcome based on the evidence provided by feature values. When the classifier is later applied to unlabeled data, it uses the observed probabilities to predict the most likely class for the new feature.

→ Bayesian classifiers used for

- * Text classification, such as junk email (spam) filtering
- * Intrusion or anomaly detection in computer networks.
- * Diagnosing medical conditions given a set of observed symptoms.

- It suffices to say that a probability is a number between 0 and 1 (that is between 0 percent and 100 percent) which captures the chance that an event will occur in the light of the available evidence. The lower the probability, the less likely the event is to occur. The probability 0 that the event will definitely not occur, while probability 1 \Rightarrow the event will occur with 100 percent certainty.

Basic concept of Bayesian methods

- Bayesian probability theory is rooted in the idea that the estimated likelihood of an event or a potential outcome, should be based on the evidence at hand across multiple trials, or opportunities for the event to occur.

Event

Heads results

Rainy weather

Message is 8pm

Candidate becomes president

win the lottery

Trial

coin flip

A single day.

Incoming e-mail msg.

Presidential election

lottery ticket.

Bayesian methods provide insights into how the probability of these events can be estimated from the observed data.

Understanding probability

- The probability of an event is estimated from the observed data by dividing the no of trials in which the event occurred by the total no of trials.

$$\text{Eg : } \frac{3}{10} = 0.30 \text{ or } 30\%$$

- To denote these probabilities, use notation $P(A)$, signifies the probability of event A Eg : $P(\text{rain}) = 0.30$.

- The probability of all possible outcomes of a trial must always sum to 1, because a trial always results in some outcome happening.

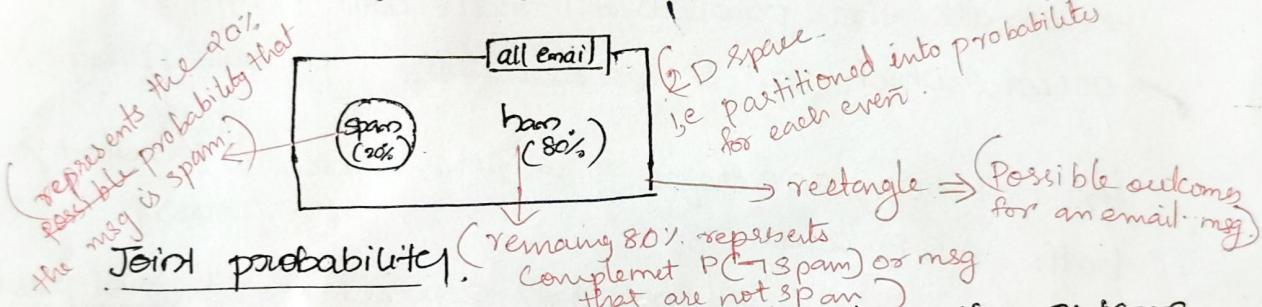
(2)

thus, if the trial has two outcomes that cannot occur simultaneously, such as rainy versus sunny.

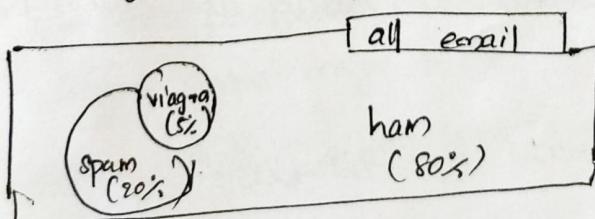
$$\text{eg: } p(\text{spam}) = 0.20 \quad p(\text{ham}) = 1 - 0.20 = 0.80$$

→ this concludes that spam & ham are mutually exclusive & exhaustive events, which implies that they cannot occur at the same time & are the only possible outcomes.

→ Bz an event cannot simultaneously happen and not happen, an event is always mutually exclusive and exhaustive with its complement. The complement of event A is typically denoted by A^c , $A!$. $P(\neg A) \stackrel{\text{short hand notation}}{=} P(A')$

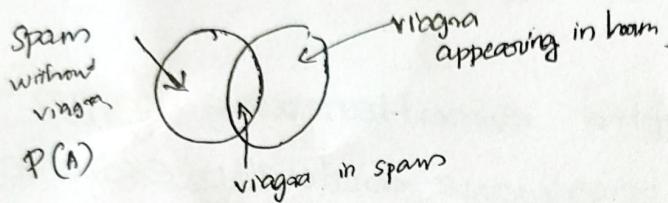


→ for instance, a second event based on the outcome appears that an e-mail message contains the word viagra. In most cases, this word is likely to appear only in a spam message; its presence in an incoming e-mail is therefore a very strong piece of evidence that the message is spam.



(4) \rightarrow diagram that the viagra circle does not fill the spam circle, nor is it completely contained in the spam circle. This implies that not all spam messages contain the word viagra and not every e-mail with word viagra is spam.

\rightarrow Venn diagram \rightarrow a visualisation.



\rightarrow 20% of all msgs were spam & 5% of all msgs contained the word viagra.

\rightarrow estimate the probability that both $p(\text{spam})$ and $p(\text{viagra})$ occur, which can be ~~written as~~ \downarrow $p(\text{spam} \cap \text{viagra})$

the notation $A \cap B$ refers to event in which "intersection of event both A & B occurs".

\rightarrow calculating $p(\text{spam} \cap \text{viagra})$ depends on the joint probability of two events. If the two events are totally unrelated they are called independent events, implies that knowing the outcome of one event does not provide any information about the outcome of the other. E.g. outcome of a heads result on a coin flip is independent from whether the weather is rainy or sunny on any given day.

\rightarrow dependent events on the basis of predictive modeling.

not contained in the event A or B

→ independent events A & B the probability of both happening can be expressed as $P(A \cap B) = P(A) * P(B)$

conditional probability with Bayes' theorem

→ the relationships between dependent events can be described using Bayes' theorem.

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

$P(A|B)$ → probability of event A, given that event B occurred
this is known as conditional probability

By definition, $P(A \cap B) = P(A|B) * P(B)$

Rearranging this formula that $P(A \cap B) = P(B \cap A)$
results in $P(A \cap B) = P(B|A) * P(A)$

formulation of Bayes' theorem.

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

likelihood prior probability

$$P(\text{spam}|\text{viagra}) = \frac{P(\text{viagra}|\text{spam})P(\text{spam})}{P(\text{viagra})}$$

posterior probability marginal likelihood

prior probability → probability that any prior msg was spam. This estimate is known as the prior probability

likelihood → the probability that ~~the~~ viagra was used in previous spam msgs. $P(VI)$

marginal likelihood → probability that viagra appeared in any msg at all

(6)

→ To calculate the components of Bayes' theorem,
to construct a frequency table.

		Viagra		Total
Frequency	Yes	No		
spam	4	16	20	
ham	1	79	80	
Total	5	100	100	

Rows of likelihood table indicate the conditional probabilities for yes/no that em was either spam or ham

		Viagra		Total
Likelihood	yes	no		
spam	4/20	16/20	20	
ham	1/80	79/80	80	
Total	5/100	95/100	100	

→ The frequency table can be used to construct a likelihood table.

$$P(\text{Viagra} = \text{Yes} | \text{spam}) = \frac{4}{20} = 0.20 \Rightarrow \text{the probability of } 20\% \text{ that a msg contains the term Viagra, given that msg is spam.}$$

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

$$P(\text{spam} \cap \text{Viagra}) \text{ as } P(\text{Viagra} | \text{spam}) * P(\text{spam}) = \frac{4}{20} * \frac{20}{100} = \underline{\underline{0.04}}$$

which notes that 4 out of 100 msgs were Spam with the term Viagra illustrates the importance of Bayes' theorem while calculating conditional probability

To compute the posterior probability,

$$P(\text{spam} | \text{Viagra})$$

$$P(\text{Viagra} | \text{spam}) * P(\text{spam}) / P(\text{Viagra}) \text{ or} \\ (\underline{\underline{4/20}}) * (\underline{\underline{20/100}}) / (\underline{\underline{5/100}}) = 0.80$$

Therefore the probability is 80% that a message is spam, given that it contains the word Viagra.

Classification with Naive Bayes

spam	+
ham	-

The Naive Bayes learner is trained by constructing a likelihood table for the appearance of these four words (labeled w_1, w_2, w_3 , and w_4) for 100 e-mails.

likelihood	Viagra (w_1)		Money (w_2)		Groceries (w_3)		Unsubscribe (w_4)		Total
	Yes	No	Yes	No	Yes	No	Yes	No	
spam	4/20	16/20	10/20	10/20	0/20	20/20	12/20	8/20	20
ham	1/80	79/80	14/80	66/80	8/80	71/80	23/80	57/80	80
Total	5/100	95/100	24/100	76/100	8/100	91/100	35/100	65/100	100

Using Bayes's theorem,

Problem : suppose that a msg contains the terms Viagra and Unsubscribe, but does not contain either Money or Groceries.

Given that Viagra = Yes, Money = No, Grocery = No, Unsubscribe = Yes we can define the problem using the formula.

$$P(\text{spam} | w_1 \cap w_2 \cap w_3 \cap w_4) = \frac{P(w_1 \cap w_2 \cap w_3 \cap w_4 | \text{spam}) P(\text{spam})}{P(w_1 \cap w_2 \cap w_3 \cap w_4)}$$

Assuming conditional independence allows us to simplify the formula using the probability rule for independent events, which states that

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

conditional probability of spam expressed as.

$$P(\text{spam} | w_1 \cap w_2 \cap w_3 \cap w_4) \propto P(w_1 | \text{spam}) P(\neg w_2 | \text{spam}) \\ P(\neg w_3 | \text{spam}) P(w_4 | \text{spam}) * P(\text{spam})$$

Probability that the message is ham can be expressed

$$P(\text{ham} | W_1 \cap \neg W_2 \cap \neg W_3 \cap \neg W_4) \propto P(W_1 | \text{ham}) P(\neg W_2 | \text{ham}) \\ P(\neg W_3 | \text{ham}) P(W_4 | \text{ham}) P(\text{ham})$$

equals symbol has been replaced by the \propto symbol to indicate the fact that the denominator has been omitted.

Using values in the likelihood table.

Overall likelihood of spam is:

$$(4/20) * (10/20) * (20/20) * (12/20) * (60/100) = 0.012$$

likelihood of ham is:

$$(1/80) * (66/80) * (71/80) * (23/80) * (80/100) = 0.002$$

$$0.012 / 0.002 = 6 \checkmark$$

this msg is 6 times more likely to be Spam than ham.

To convert these numbers into probabilities.

$$0.012 / (0.012 + 0.002) = 0.857$$

[the probability of spam is equal to the likelihood that the message is spam divided by the likelihood that msg is either spam or ham]

Similarly, the probability of ham is equal to the likelihood that the msg is ham divided by the likelihood that msg is either spam or ham.

$$0.002 / (0.012 + 0.002) = 0.143$$

message is spam with 85.7 percent probability and ham with 14.3 percent probability. Because these are mutually exclusive and exhaustive events, the probabilities sum to 1.

Finally Naive Bayes classification algorithm we used in the preceding example can be summarized by.

probability of level L for class C , given the evidence provided by features F_1 through F_n , is equal to the product of the probabilities of each piece of evidence conditioned on the class level, the prior probability of the class level, and a scaling factor $1/Z$, which converts the likelihood values into probabilities

$$P(C_L | F_1, \dots, F_n) = \frac{1}{Z} P(C_L) \prod_{i=1}^n P(F_i | C_L)$$

The Laplace estimator

Using
nume

Using the Naive Bayes algorithm as before, we can compute the likelihood of spam as:

$$(4/20) * (10/20) * (0/20) * (12/20) * (20/100) = 0$$

The likelihood of ham is

$$(1/80) * (14/80) * (8/80) * (23/80) * (80/100) = 0.00005$$

Therefore, the probability of spam is

$$0 / (0 + 0.00005) = 0$$

the probability of ham is

$$0.00005 / (0 + 0.00005) = 1$$

These result suggest that the msg is spam with 0 percent probability and ham with 100% probability.

This prediction make no sense.

Probabilities in the Naive Bayes formula are multiplied in a chain, this 0 percent value causes the posterior probability of spam to be zero, giving the word groceries the ability to effectively nullify and overrule all of the other evidence.

- A solution to this problem involves using something called the Laplace estimator, which is named after the french mathematician Pierre-Simon Laplace. The Laplace estimator essentially adds a small number to each of the counts in the frequency table, which ensures that each feature has a nonzero probability of occurring with each class. Typically, the Laplace estimator is set to 1, which ensures that each class-feature combination is found in the data at least once.

Using a Laplace value of 1, we add one to each numerator in the likelihood function.

The Likelihood of spam is

$$(5/24) * (11/24) * (1/24) * (13/24) * (20/100) = 0.0004$$

The Likelihood of ham is

$$(2/84) * (15/84) * (9/84) * (24/84) * (80/100) = 0.0001$$

This means that the probability of spam is 80%,
and the probability of ham is 20%.

==

Using numeric features with Naive Bayes

- Naive Bayes uses frequency tables to learn the data, each feature must be categorical in order to create the combinations of class and feature values comprising of the matrix.

Naive Bayes algorithm does not work with ~~numeric~~ numeric data.

One easy and effective solution is to discretize numeric features, which simply means that the numbers are put into categories known as bins.

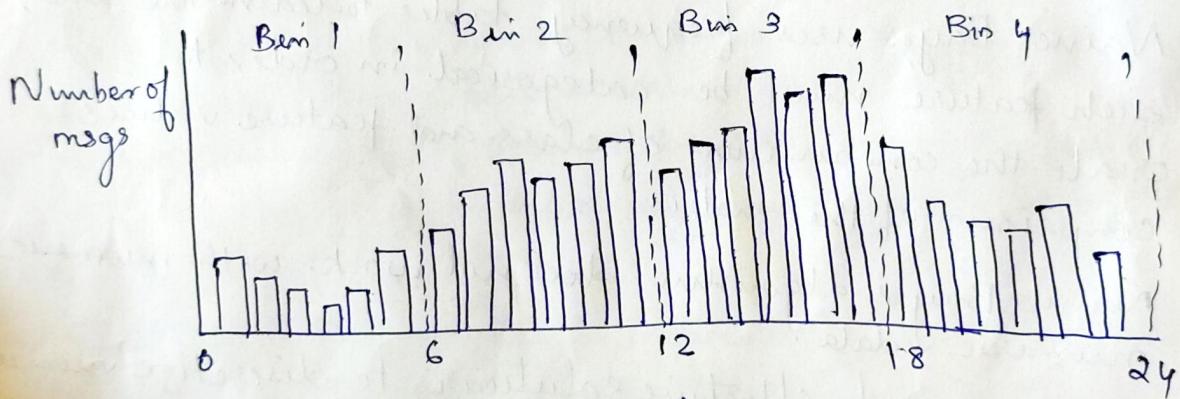
For this reason, discretization is also sometimes called binning. This method is ideal when there are large amounts of training data, a common condition while working with Naive Bayes.

- There are several different ways to discretize a numeric feature. Perhaps the most common is to split the data for natural categories or cut points in the distribution of data.

for eg : Suppose that you added a feature to the spam dataset that recorded the time of night or day the e-mail was sent, from 0 to 24 hrs past midnight.

Depicted using a histogram, the time data might look in the following diagram.

- ⇒ In the early hours of the morning, the msg frequency is low.
- ⇒ The activity picks up during business hours and tapers off in the evening.
- ⇒ This seems to create four natural bins of activity, as partitioned by the dashed lines indicating places where the numeric data are divided into levels of a new nominal feature,



Discretizing a numeric feature always results in a reduction of information as the feature's original granularity is reduced to a smaller number of categories.

A few bins can result in important trends being obscured.

Too many bins can result in small counts in the Naive Bayes frequency table, which can increase the algorithm's sensitivity to noisy data.