



REVIEW OF PROBABILITY CONCEPTS

SYRACUSE UNIVERSITY
School of Information Studies

BAYES' THEOREM

What is Bayes' theorem mathematically?

How can computers use Bayes' theorem?

REVIEW OF PROBABILITY CONCEPTS

An **event** is an outcome of an experiment.

Experiment: Toss a fair coin

Possible outcomes: “head” or “tail”

The **probability** that either event (“head” or “tail”) occurs is 50%, i.e.:

$$P(\text{head}) = 0.5, P(\text{tail}) = 0.5$$

$$P(\text{head}) + P(\text{tail}) = 1$$

If the coin is not balanced, then $P(\text{head}) \neq P(\text{tail})$.

JOINT PROBABILITY

Joint probability $P(A, B)$ is the chance that two events A and B co-occurred.

Experiment: Toss two fair coins

Event A: coin1 = “head”

Event B: coin2 = “head”

$P(A, B) = P(\text{coin1} = \text{“head”}, \text{coin2} = \text{“head”})$

$P(A, B)$ and $P(B, A)$ are the same thing.

INDEPENDENT EVENTS

The occurrence of one event does not depend on the occurrence of the other event.

E.g., tossed two fair coins; both landed head up

Event A: coin1 = “head”

Event B: coin2 = “head”

$$P(A, B) = P(A) * P(B) = 0.5 * 0.5 = 0.25$$

DEPENDENT EVENTS

The occurrence of one event is dependent on the occurrence of another event.

E.g., choose a card from a standard deck of 52 cards. Without replacement, choose another card. What is the probability of choosing two queens?

Because the outcome of the second card is dependent on the outcome of the first card,
 $P(\text{card1} = \text{"queen"}, \text{card2} = \text{"queen"}) \neq P(\text{card1} = \text{"queen"}) * P(\text{card2} = \text{"queen"})$

CONDITIONAL PROBABILITY

Conditional probability $P(B | A)$

The probability that event B occurs if event A occurs

Event A: card1 = “queen”

Event B: card2 = “queen”

$P(B | A) = P(\text{card2} = \text{“queen”} \mid \text{card1} = \text{“queen”}) = 3/51$, because there are 51 remaining cards and only three queens among them.

RELATIONSHIP BETWEEN CONDITIONAL AND JOINT PROBABILITIES

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

$$\begin{aligned} &P(\text{card1} = \text{"queen"}, \text{card2} = \text{"queen"}) \\ &= P(\text{card1} = \text{"queen"}) * P(\text{card2} = \text{"queen"} \mid \text{card1} = \text{"queen"}) \\ &= 4/52 * 3/51 \\ &= 0.5\% \end{aligned}$$

PROBABILITY IS USEFUL IN DAILY LIFE

Which one is greater?

$P(\text{"will"} \mid \text{"I"})$

$P(\text{"has"} \mid \text{"I"})$

This is how your smartphone knows what words to recommend when you are texting!

All words that people typed become the training corpus, in which the conditional probabilities are calculated.



BAYES RULES IN MAMMOGRAM EXAMPLE

SYRACUSE UNIVERSITY
School of Information Studies

BAYES' THEOREM

Bayes' theorem lets us swap the order of the dependence between events.

Two events A and B

We know that $P(A,B) = P(B|A)P(A)$

Since $P(A,B) = P(B,A)$, we also know that $P(B,A) = P(B|A)P(A)$

Therefore:

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

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

BACK TO THE MAMMOGRAM EXAMPLE IN PROFESSOR STROGATZ'S ARTICLE

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

Event A: A patient has cancer

Event B: A patient's mammogram test result is positive

$P(B | A)$ = Among patients who have cancer, how many have positive test results?

We can calculate it from past data.

$P(A | B)$ = For a patient with a positive test result, what is the chance of cancer?

This is **prediction**!

Bayes' theorem provides an approach to predict the probability of future events based on prior experience.

THE MAMMOGRAM EXAMPLE

The probability that a woman has breast cancer is ... ?

If a woman has breast cancer, the probability that she will have a positive mammogram is ... ?

If a woman does not have breast cancer, the probability that she will still have a positive mammogram is ... ?

BACK TO THE MAMMOGRAM EXAMPLE

The probability that a woman has breast cancer is

$$P(\text{cancer}) = 0.008$$

If a woman has breast cancer, the probability that she will have a positive mammogram is

$$P(\text{positive} | \text{cancer}) = 0.9$$

If a woman does not have breast cancer, the probability that she will still have a positive mammogram is

$$P(\text{positive} | \text{no cancer}) = 0.07$$

TRANSLATE THEM INTO PROBABILITY NOTATIONS

Imagine a woman who has a positive mammogram. What is the probability that she actually has breast cancer?

Which probability is greater:

$P(\text{cancer} \mid \text{positive})$ or $P(\text{no cancer} \mid \text{positive})$?

To calculate $P(\text{cancer} \mid \text{positive})$, need to calculate:

$P(\text{positive})$

$P(\text{cancer})$

$P(\text{positive} \mid \text{cancer})$

To calculate $P(\text{no cancer} \mid \text{positive})$, need to calculate:

$P(\text{positive})$

$P(\text{no cancer})$

$P(\text{positive} \mid \text{no cancer})$

TRANSLATE THEM INTO PROBABILITY NOTATIONS

The probability that a woman has positive mammogram is:

$P(\text{positive}) = \text{Unknown}$

The probability that a woman has breast cancer is 0.008.

$P(\text{cancer}) = 0.008$, $P(\text{no cancer}) = 0.992$, (add up to 1)

If a woman has breast cancer, the probability that she will have a positive mammogram is 0.9.

$P(\text{positive} \mid \text{cancer}) = 0.9$

If a woman does not have breast cancer, the probability that she will still have a positive mammogram is 0.07.

$P(\text{positive} \mid \text{no cancer}) = 0.07$

TRANSLATE THEM INTO PROBABILITY NOTATIONS

All **prior probabilities** we have calculated:

$$P(\text{cancer}) = 0.008$$

$$P(\text{no cancer}) = 0.992$$

All **conditional probabilities** we have calculated:

$$P(\text{positive} \mid \text{cancer}) = 0.9$$

$$P(\text{positive} \mid \text{no cancer}) = 0.07$$

The **posterior probabilities** to be calculated:

$$P(\text{cancer} \mid \text{positive}) = ?$$

$$P(\text{no cancer} \mid \text{positive}) = ?$$

SO, OUR PREDICTION IS ...

$$\begin{aligned} &P(\text{cancer} \mid \text{positive}) \\ &= \frac{P(\text{positive} \mid \text{cancer}) \cdot P(\text{cancer})}{P(\text{positive})} = \frac{0.9 \cdot 0.008}{P(\text{positive})} = \frac{0.0072}{P(\text{positive})} \end{aligned}$$

$$\begin{aligned} &P(\text{no_cancer} \mid \text{positive}) \\ &= \frac{P(\text{positive} \mid \text{no_cancer}) \cdot P(\text{no_cancer})}{P(\text{positive})} = \frac{0.07 \cdot 0.992}{P(\text{positive})} = \frac{0.069}{P(\text{positive})} \end{aligned}$$

NO CANCER!

$$\begin{aligned} &P(\text{cancer} \mid \text{positive}) \\ &= \frac{P(\text{positive} \mid \text{cancer}) \cdot P(\text{cancer})}{P(\text{positive})} = \frac{0.9 \cdot 0.008}{P(\text{positive})} = \frac{0.0072}{P(\text{positive})} \end{aligned}$$

$$\begin{aligned} &P(\text{no_cancer} \mid \text{positive}) \\ &= \frac{P(\text{positive} \mid \text{no_cancer}) \cdot P(\text{no_cancer})}{P(\text{positive})} = \frac{0.07 \cdot 0.992}{P(\text{positive})} = \frac{0.069}{P(\text{positive})} \end{aligned}$$

Although we don't know $P(\text{positive})$, it does not matter. We just need to know which posterior probability is greater.

THIS DIAGNOSIS IS DETERMINED BY THE MAMMOGRAM RESULT ONLY

$$\begin{aligned} &P(\text{cancer} | \text{positive}) \\ &= \frac{P(\text{positive} | \text{cancer}) \cdot P(\text{cancer})}{P(\text{positive})} = \frac{0.9 \cdot 0.008}{P(\text{positive})} = \frac{0.0072}{P(\text{positive})} \end{aligned}$$

$$\begin{aligned} &P(\text{no_cancer} | \text{positive}) \\ &= \frac{P(\text{positive} | \text{no_cancer}) \cdot P(\text{no_cancer})}{P(\text{positive})} = \frac{0.07 \cdot 0.992}{P(\text{positive})} = \frac{0.069}{P(\text{positive})} \end{aligned}$$

WHAT IF THE DIAGNOSIS IS DETERMINED BY MORE FACTORS THAN JUST THE MAMMOGRAM RESULT?

Attribute 1: Positive mammogram? Yes or no

Attribute 2: Family history? Yes or no

Attribute 3: Alcohol? Yes or no

How many posteriors to calculate?

Two posteriors for each possible combination of the attributes.

In this case, $2 * 2^3 = 16$

ALL POSTERIOR PROBABILITIES FOR THREE BINARY ATTRIBUTES

$P(\text{positive}, \text{yes}, \text{yes} \mid \text{cancer})$

$P(\text{positive}, \text{yes}, \text{no} \mid \text{cancer})$

$P(\text{positive}, \text{no}, \text{yes} \mid \text{cancer})$

$P(\text{positive}, \text{no}, \text{no} \mid \text{cancer})$

$P(\text{negative}, \text{yes}, \text{yes} \mid \text{cancer})$

$P(\text{negative}, \text{yes}, \text{no} \mid \text{cancer})$

$P(\text{negative}, \text{no}, \text{yes} \mid \text{cancer})$

$P(\text{negative}, \text{no}, \text{no} \mid \text{cancer})$

$P(\text{positive}, \text{yes}, \text{yes} \mid \text{cancer})$

$P(\text{positive}, \text{yes}, \text{no} \mid \text{no_cancer})$

$P(\text{positive}, \text{no}, \text{yes} \mid \text{no_cancer})$

$P(\text{positive}, \text{no}, \text{no} \mid \text{no_cancer})$

$P(\text{negative}, \text{yes}, \text{yes} \mid \text{no_cancer})$

$P(\text{negative}, \text{yes}, \text{no} \mid \text{no_cancer})$

$P(\text{negative}, \text{no}, \text{yes} \mid \text{no_cancer})$

$P(\text{negative}, \text{no}, \text{no} \mid \text{no_cancer})$



NAIVE BAYES

SYRACUSE UNIVERSITY
School of Information Studies

CHALLENGE OF BAYESIAN CLASSIFIER

$P(\text{positive, yes, yes} \mid \text{cancer})$

$P(\text{positive, yes, no} \mid \text{cancer})$

$P(\text{positive, no, yes} \mid \text{cancer})$

$P(\text{positive, no, no} \mid \text{cancer})$

$P(\text{negative, yes, yes} \mid \text{cancer})$

$P(\text{negative, yes, no} \mid \text{cancer})$

$P(\text{negative, no, yes} \mid \text{cancer})$

$P(\text{negative, no, no} \mid \text{cancer})$

$P(\text{positive, yes, yes} \mid \text{cancer})$

$P(\text{positive, yes, no} \mid \text{no_cancer})$

$P(\text{positive, no, yes} \mid \text{no_cancer})$

$P(\text{positive, no, no} \mid \text{no_cancer})$

$P(\text{negative, yes, yes} \mid \text{no_cancer})$

$P(\text{negative, yes, no} \mid \text{no_cancer})$

$P(\text{negative, no, yes} \mid \text{no_cancer})$

$P(\text{negative, no, no} \mid \text{no_cancer})$

CHALLENGE OF BAYESIAN CLASSIFIER

$P(\text{positive}, \text{yes}, \text{yes} \mid \text{cancer})$	$P(\text{positive}, \text{yes}, \text{yes} \mid \text{cancer})$
$P(\text{positive}, \text{yes}, \text{no} \mid \text{cancer})$	$P(\text{positive}, \text{yes}, \text{no} \mid \text{no_cancer})$
$P(\text{positive}, \text{no}, \text{yes} \mid \text{cancer})$	$P(\text{positive}, \text{no}, \text{yes} \mid \text{no_cancer})$
$P(\text{positive}, \text{no}, \text{no} \mid \text{cancer})$	$P(\text{positive}, \text{no}, \text{no} \mid \text{no_cancer})$
$P(\text{negative}, \text{yes}, \text{yes} \mid \text{cancer})$	$P(\text{negative}, \text{yes}, \text{yes} \mid \text{no_cancer})$
$P(\text{negative}, \text{yes}, \text{no} \mid \text{cancer})$	$P(\text{negative}, \text{yes}, \text{no} \mid \text{no_cancer})$
$P(\text{negative}, \text{no}, \text{yes} \mid \text{cancer})$	$P(\text{negative}, \text{no}, \text{yes} \mid \text{no_cancer})$
$P(\text{negative}, \text{no}, \text{no} \mid \text{cancer})$	$P(\text{negative}, \text{no}, \text{no} \mid \text{no_cancer})$

2^{n+1} possible combinations of values for n
binary attributes

CHALLENGE OF BAYESIAN CLASSIFIER

$P(\text{positive}, \text{yes}, \text{yes} \mid \text{cancer})$

$P(\text{positive}, \text{yes}, \text{no} \mid \text{cancer})$

$P(\text{positive}, \text{no}, \text{yes} \mid \text{cancer})$

$P(\text{positive}, \text{no}, \text{no} \mid \text{cancer})$

$P(\text{negative}, \text{yes}, \text{yes} \mid \text{cancer})$

$P(\text{negative}, \text{yes}, \text{no} \mid \text{cancer})$

$P(\text{negative}, \text{no}, \text{yes} \mid \text{cancer})$

$P(\text{negative}, \text{no}, \text{no} \mid \text{cancer})$

$P(\text{positive}, \text{yes}, \text{yes} \mid \text{cancer})$

$P(\text{positive}, \text{yes}, \text{no} \mid \text{no_cancer})$

$P(\text{positive}, \text{no}, \text{yes} \mid \text{no_cancer})$

$P(\text{positive}, \text{no}, \text{no} \mid \text{no_cancer})$

$P(\text{negative}, \text{yes}, \text{yes} \mid \text{no_cancer})$

$P(\text{negative}, \text{yes}, \text{no} \mid \text{no_cancer})$

$P(\text{negative}, \text{no}, \text{yes} \mid \text{no_cancer})$

$P(\text{negative}, \text{no}, \text{no} \mid \text{no_cancer})$

2ⁿ possible combinations of values for n binary attributes

What if n = 10,000? Quite common for text classification when each word is an attribute.

CHALLENGE OF BAYESIAN CLASSIFIER

$P(\text{positive}, \text{yes}, \text{yes} \mid \text{cancer})$

$P(\text{positive}, \text{yes}, \text{no} \mid \text{cancer})$

$P(\text{positive}, \text{no}, \text{yes} \mid \text{cancer})$

$P(\text{positive}, \text{no}, \text{no} \mid \text{cancer})$

$P(\text{negative}, \text{yes}, \text{yes} \mid \text{cancer})$

$P(\text{negative}, \text{yes}, \text{no} \mid \text{cancer})$

$P(\text{negative}, \text{no}, \text{yes} \mid \text{cancer})$

$P(\text{negative}, \text{no}, \text{no} \mid \text{cancer})$

$P(\text{positive}, \text{yes}, \text{yes} \mid \text{cancer})$

$P(\text{positive}, \text{yes}, \text{no} \mid \text{no_cancer})$

$P(\text{positive}, \text{no}, \text{yes} \mid \text{no_cancer})$

$P(\text{positive}, \text{no}, \text{no} \mid \text{no_cancer})$

$P(\text{negative}, \text{yes}, \text{yes} \mid \text{no_cancer})$

$P(\text{negative}, \text{yes}, \text{no} \mid \text{no_cancer})$

$P(\text{negative}, \text{no}, \text{yes} \mid \text{no_cancer})$

$P(\text{negative}, \text{no}, \text{no} \mid \text{no_cancer})$

Too many probabilities to calculate!

CHALLENGE OF BAYESIAN CLASSIFIER

$P(\text{positive, yes, yes} \mid \text{cancer})$

$P(\text{positive, yes, no} \mid \text{cancer})$

$P(\text{positive, no, yes} \mid \text{cancer})$

$P(\text{positive, no, no} \mid \text{cancer})$

$P(\text{negative, yes, yes} \mid \text{cancer})$

$P(\text{negative, yes, no} \mid \text{cancer})$

$P(\text{negative, no, yes} \mid \text{cancer})$

$P(\text{negative, no, no} \mid \text{cancer})$

$P(\text{positive, yes, yes} \mid \text{no_cancer})$

$P(\text{positive, yes, no} \mid \text{no_cancer})$

$P(\text{positive, no, yes} \mid \text{no_cancer})$

$P(\text{positive, no, no} \mid \text{no_cancer})$

$P(\text{negative, yes, yes} \mid \text{no_cancer})$

$P(\text{negative, yes, no} \mid \text{no_cancer})$

$P(\text{negative, no, yes} \mid \text{no_cancer})$

$P(\text{negative, no, no} \mid \text{no_cancer})$

There may not be enough data for accurately estimating each probability.

NAIVE BAYES CLASSIFIER

Assume independence among attributes A_i when class is given:

$$\begin{aligned} P(A_1, A_2, \dots, A_n \mid C) &= P(A_1 \mid C_j) * P(A_2 \mid C_j) * \dots * P(A_n \mid C_j) \\ &= \prod P(A_i \mid C_j) \end{aligned}$$

This assumption means that given the target class C , the probability of observing the conjunction $A_1, A_2, A_3, \dots, A_n$ is just the product of the probabilities for the individual attributes.

For example, for people with cancer, the assumption is that whether there is family history of cancer has nothing to do with race, which may or may not be scientifically true.

WHY IS THE INDEPENDENCE ASSUMPTION NEEDED?

Assume independence among attributes A_i when class is given:

$$\begin{aligned} P(A_1, A_2, \dots, A_n \mid C) &= P(A_1 \mid C_j) P(A_2 \mid C_j) \dots P(A_n \mid C_j) \\ &= \prod P(A_i \mid C_j) \end{aligned}$$

To greatly reduce the number of probabilities to estimate

WHY IS THE INDEPENDENCE ASSUMPTION NEEDED?

$$P(\text{positive}, \text{yes}_{\text{family}}, \text{yes}_{\text{race}} \mid \text{cancer}) \\ = P(\text{positive} \mid \text{cancer}) \cdot P(\text{yes}_{\text{family}} \mid \text{cancer}) \cdot P(\text{yes}_{\text{race}} \mid \text{cancer})$$

$P(\text{positive}, \text{yes}, \text{yes} \mid \text{cancer})$
 $P(\text{positive}, \text{yes}, \text{no} \mid \text{cancer})$
 $P(\text{positive}, \text{no}, \text{yes} \mid \text{cancer})$
 $P(\text{positive}, \text{no}, \text{no} \mid \text{cancer})$
 $P(\text{negative}, \text{yes}, \text{yes} \mid \text{cancer})$
 $P(\text{negative}, \text{yes}, \text{no} \mid \text{cancer})$
 $P(\text{negative}, \text{no}, \text{yes} \mid \text{cancer})$
 $P(\text{negative}, \text{no}, \text{no} \mid \text{cancer})$

Reduce
calculation
 $2^3 \rightarrow 2*3$

$P(\text{positive} \mid \text{cancer})$
 $P(\text{yes}_{\text{family}} \mid \text{cancer})$
 $P(\text{yes}_{\text{race}} \mid \text{cancer})$
 $P(\text{negative} \mid \text{cancer})$
 $P(\text{no}_{\text{family}} \mid \text{cancer})$
 $P(\text{no}_{\text{race}} \mid \text{cancer})$

SIMILARLY, ...

$P(\text{positive}, \text{yes}, \text{yes} \mid \text{no_cancer})$
 $P(\text{positive}, \text{yes}, \text{no} \mid \text{no_cancer})$
 $P(\text{positive}, \text{no}, \text{yes} \mid \text{no_cancer})$
 $P(\text{positive}, \text{no}, \text{no} \mid \text{no_cancer})$
 $P(\text{negative}, \text{yes}, \text{yes} \mid \text{no_cancer})$
 $P(\text{negative}, \text{yes}, \text{no} \mid \text{no_cancer})$
 $P(\text{negative}, \text{no}, \text{yes} \mid \text{no_cancer})$
 $P(\text{negative}, \text{no}, \text{no} \mid \text{no_cancer})$

Reduce
calculation
 $2^3 \rightarrow 2*3$

$P(\text{positive} \mid \text{no_cancer})$
 $P(\text{yes}_{\text{family}} \mid \text{no_cancer})$
 $P(\text{yes}_{\text{race}} \mid \text{no_cancer})$
 $P(\text{negative} \mid \text{no_cancer})$
 $P(\text{no}_{\text{family}} \mid \text{no_cancer})$
 $P(\text{no}_{\text{race}} \mid \text{no_cancer})$

SIMILARLY, ...

$P(\text{positive}, \text{yes}, \text{yes} \mid \text{no_cancer})$

$P(\text{positive}, \text{yes}, \text{no} \mid \text{no_cancer})$

$P(\text{positive}, \text{no}, \text{yes} \mid \text{no_cancer})$

$P(\text{positive}, \text{no}, \text{no} \mid \text{no_cancer})$

$P(\text{negative}, \text{yes}, \text{yes} \mid \text{no_cancer})$

$P(\text{negative}, \text{yes}, \text{no} \mid \text{no_cancer})$

$P(\text{negative}, \text{no}, \text{yes} \mid \text{no_cancer})$

$P(\text{negative}, \text{no}, \text{no} \mid \text{no_cancer})$

Reduce
calculation
 $2^3 \rightarrow 2*3$

$P(\text{positive} \mid \text{no_cancer})$

$P(\text{yes}_{\text{family}} \mid \text{no_cancer})$

$P(\text{yes}_{\text{race}} \mid \text{no_cancer})$

$P(\text{negative} \mid \text{no_cancer})$

$P(\text{no}_{\text{family}} \mid \text{no_cancer})$

$P(\text{no}_{\text{race}} \mid \text{no_cancer})$

Total calculation reduction $2*2^n \rightarrow 2*2*n$

If $n = 10$, then 2,048 probabilities reduced to 40!

WHY DOES “NAIVE” BAYES WORK?

This assumption is often violated in real-world problems.

E.g., family history of cancer may be correlated with race.

However, naive Bayes algorithm works quite well in many problems, even when the assumption is violated. There are some mathematical explanation for this phenomenon.

Domingos, P., & Pazzani, M. (1997). On the optimality of the simple Bayesian classifier under zero-one loss. *Machine Learning*, 29, 103–30.

What if the correlations between attributes have to be addressed?

Use Bayesian belief network (BBN)



SMOOTHING

SYRACUSE UNIVERSITY
School of Information Studies

PROBLEM OF ZERO PROBABILITIES

$$\begin{aligned} & P(\text{cancer} = \text{no} \mid \text{pos_mammo} = \text{yes}, \text{fam_hist} = \text{yes}, \text{caucasian} = \text{yes}) \\ &= \frac{P(\text{pos_mammo} = \text{yes} \mid \text{cancer} = \text{no}) \cdot P(\text{fam_hist} = \text{yes} \mid \text{cancer} = \text{no}) \cdot P(\text{caucasian} = \text{yes} \mid \text{cancer} = \text{no}) \cdot P(\text{cancer} = \text{no})}{P(\text{pos_mammo} = \text{yes}, \text{fam_hist} = \text{yes}, \text{caucasian} = \text{yes})} \\ &= \frac{(2/4) \cdot (0/4) \cdot (2/4) \cdot (4/7)}{P(\text{pos_mammo} = \text{yes}, \text{fam_hist} = \text{yes}, \text{caucasian} = \text{yes})} \\ &= 0 \end{aligned}$$

Zero probability of no cancer?

PROBLEM OF ZERO PROBABILITIES

If one of the conditional probabilities is zero, then the entire product becomes zero.

But the zero probability may just be caused by lack of data.

$$P(\text{fam_hist}=\text{yes} \mid \text{cancer}=\text{yes}) = 1$$

$$P(\text{fam_hist}=\text{no} \mid \text{cancer}=\text{yes}) = 0$$

$$P(\text{fam_hist}=\text{yes} \mid \text{cancer}=\text{no}) = 0$$

$$P(\text{fam_hist}=\text{no} \mid \text{cancer}=\text{no}) = 1$$

SOLUTION TO ZERO PROBABILITY

Probability estimation should replace these zero probabilities with a very small probability, which means such events still occur in real world but are so rare that the training data did not include any of them.

Since all probabilities should add up to one, the other nonzero probabilities need to “shrink” a little bit, in order to “lend” the small amount to the zero probabilities.

This technique is called “smoothing.”

SMOOTHING FOR ZERO PROBABILITIES

Using a **smoothing** algorithm, such as Laplacian smoothing, also called “add-one” smoothing

$$\text{Original: } P(A_i | C) = \frac{N_{ic}}{N_c}$$

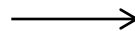
$$\text{Laplace: } P(A_i | C) = \frac{N_{ic} + 1}{N_c + c}$$

A binary classifier has two classes ($c = 2$).

ADD-ONE SMOOTHING

Example	Pos_mammo	Fam_hist	Alcohol	Cancer
1	Yes	Yes	Yes	Yes
2	Yes	Yes	No	Yes
3	No	Yes	Yes	Yes
4	Yes	No	No	No
5	Yes	No	Yes	No
6	No	No	Yes	No
7	No	No	No	No

	Cancer = Yes	Cancer = No
Fam_hist = Yes	3	0
Fam_hist = No	0	4



	Cancer = Yes	Cancer = No
Fam_hist = Yes	3 + 1	0 + 1
Fam_hist = No	0 + 1	4 + 1

Add-one smoothing means to add an example to each category.

ADD-ONE SMOOTHING

	Cancer = Yes	Cancer = No		Cancer = Yes	Cancer = No
Fam_hist = Yes	3	0	→	3 + 1	0 + 1
Fam_hist = No	0	4		0 + 1	4 + 1

Add-one smoothing means to add an example to each category.

Original probabilities:

$$P(\text{fam_hist}=\text{yes} \mid \text{cancer}=\text{yes}) = 3/3$$

$$P(\text{fam_hist}=\text{no} \mid \text{cancer}=\text{yes}) = 0/3 = 0$$

$$P(\text{fam_hist}=\text{yes} \mid \text{cancer}=\text{no}) = 0/4 = 0$$

$$P(\text{fam_hist}=\text{no} \mid \text{cancer}=\text{no}) = 4/4 = 1$$

Smoothed probabilities:

$$P(\text{fam_hist}=\text{yes} \mid \text{cancer}=\text{yes}) = 4/5$$

$$P(\text{fam_hist}=\text{no} \mid \text{cancer}=\text{yes}) = 1/5$$

$$P(\text{fam_hist}=\text{yes} \mid \text{cancer}=\text{no}) = 1/6$$

$$P(\text{fam_hist}=\text{no} \mid \text{cancer}=\text{no}) = 5/6$$

The nonzero probabilities become smaller,
“lending” values to zero probabilities.

SMOOTHING FOR ZERO PROBABILITIES

Using a **smoothing** algorithm, such as Laplacian smoothing, also called “add-one” smoothing

$$\text{Original : } P(A_i | C) = \frac{N_{ic}}{N_c}$$
$$\text{Laplace : } P(A_i | C) = \frac{N_{ic} + 1}{N_c + c}$$

A binary classifier has two classes ($c = 2$).

$$\frac{0}{3} = 0$$
$$\frac{0+1}{3+2} = \frac{1}{5} = 0.2$$

Well, this is not a very small probability!

SMOOTHING FOR ZERO PROBABILITIES

Using a **smoothing** algorithm, such as Laplacian smoothing, also called “add-one” smoothing

$$\text{Original: } P(A_i | C) = \frac{N_{ic}}{N_c}$$

$$\text{Laplace: } P(A_i | C) = \frac{N_{ic} + 1}{N_c + c}$$

The probabilities in real-world problems are usually very small; thus, add-one smoothing would change the probability just a little bit.

$$\frac{1}{7000} = 0.000143$$

$$\frac{1+1}{7000+2} = \frac{2}{7002} = 0.000285$$

LOG PROBABILITIES

The probabilities are so small that we usually use $\log(P)$ instead, so that we don't have to store that many preceding zeros, as in 0.000222.

$$\log\left(\frac{1}{7000}\right) = 8.854$$

$$\log\left(\frac{1+1}{7000+2}\right) = \log\left(\frac{2}{7002}\right) = 8.161$$



PROBABILITY OF CONTINUOUS VARIABLE

SYRACUSE UNIVERSITY
School of Information Studies

HOW TO ESTIMATE PROBABILITIES OF CONTINUOUS ATTRIBUTES

Two approaches for continuous attributes:

Discretization

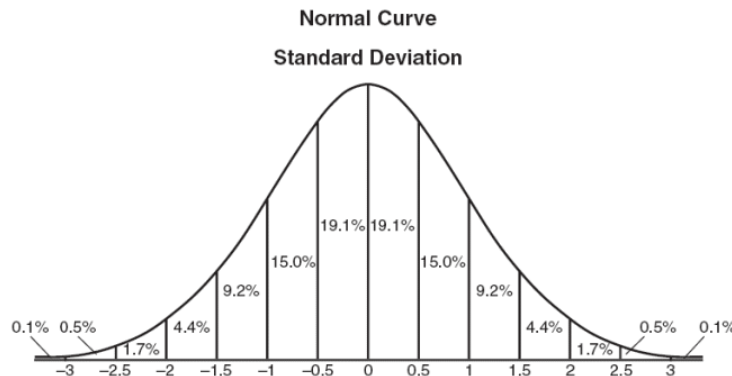
(0, 60k), (60k, 100k), (100k, ...)

Probability density estimation

Assume attribute follows a normal distribution.

Use data to estimate parameters of distribution
(e.g., mean and standard deviation).

Once probability distribution is known, can use the probability density function to estimate the conditional probability $P(A_i | c)$.



HOW TO ESTIMATE PROBABILITIES OF CONTINUOUS ATTRIBUTES

<i>Tid</i>	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

Normal distribution:

$$P(A_i | c_j) = \frac{1}{\sqrt{2\pi} \sigma_{ij}} e^{-\frac{(A_i - \mu_{ij})^2}{2\sigma_{ij}^2}}$$

One for each (A_i, c_j) pair

For (Income, Class = No):

If Class = No

Sample mean = 110

Sample variance = 2,975

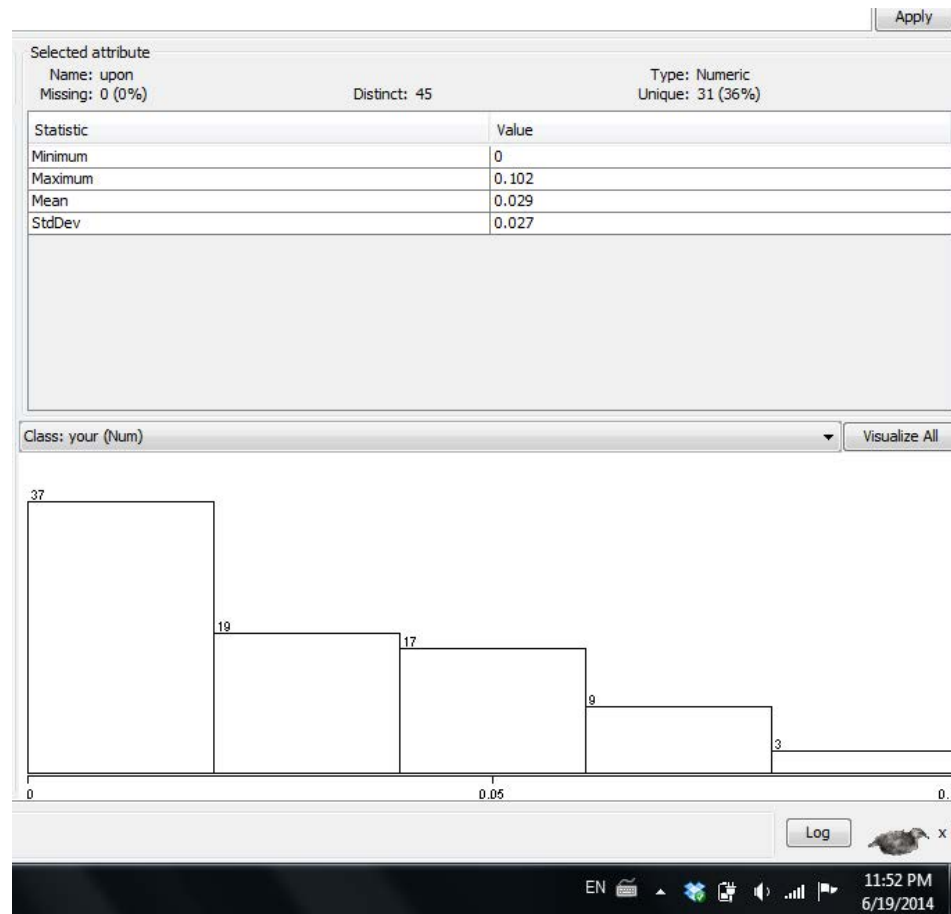
$$P(\text{Income} = 120 | \text{No}) = \frac{1}{\sqrt{2\pi} (54.54)} e^{-\frac{(120 - 110)^2}{2(2975)}} = 0.0072$$

HOW DO I KNOW IF A VARIABLE FOLLOWS NORMAL DISTRIBUTION?

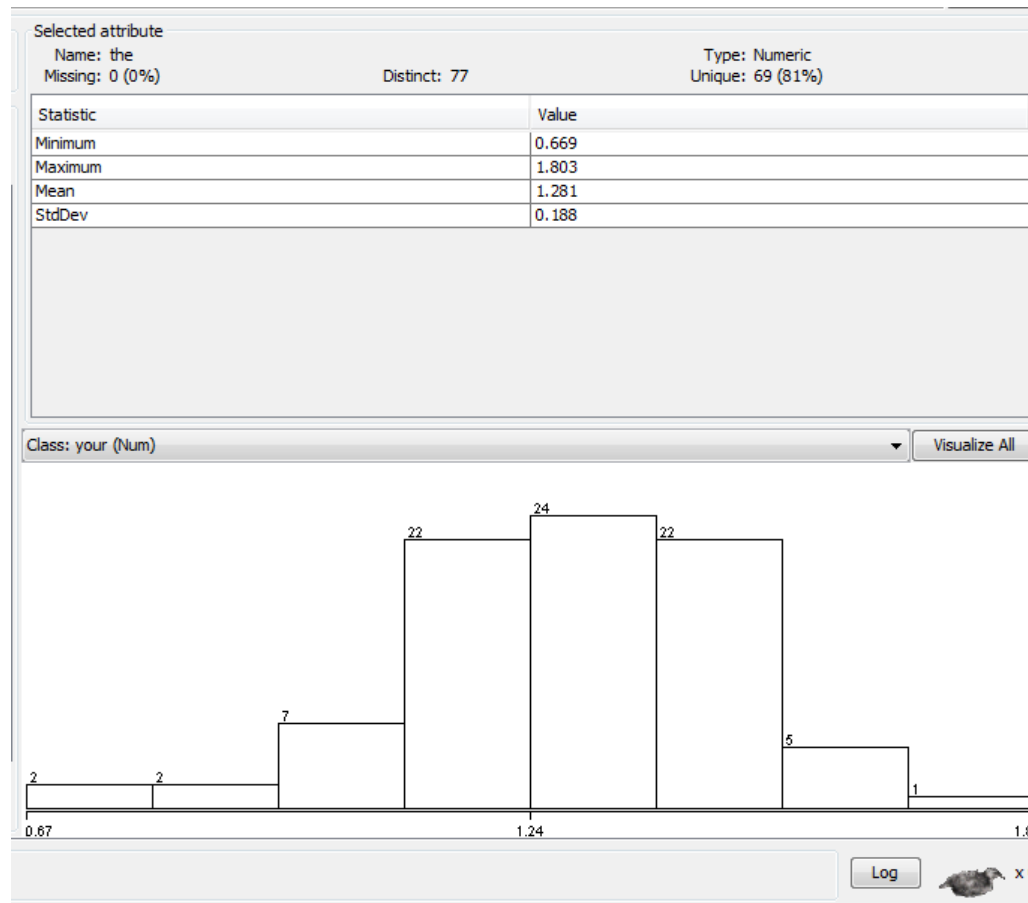
Approach 1: Use statistics test.

Approach 2 (recommended to our class):
Use visualization. (Does it look like a bell
curve?)

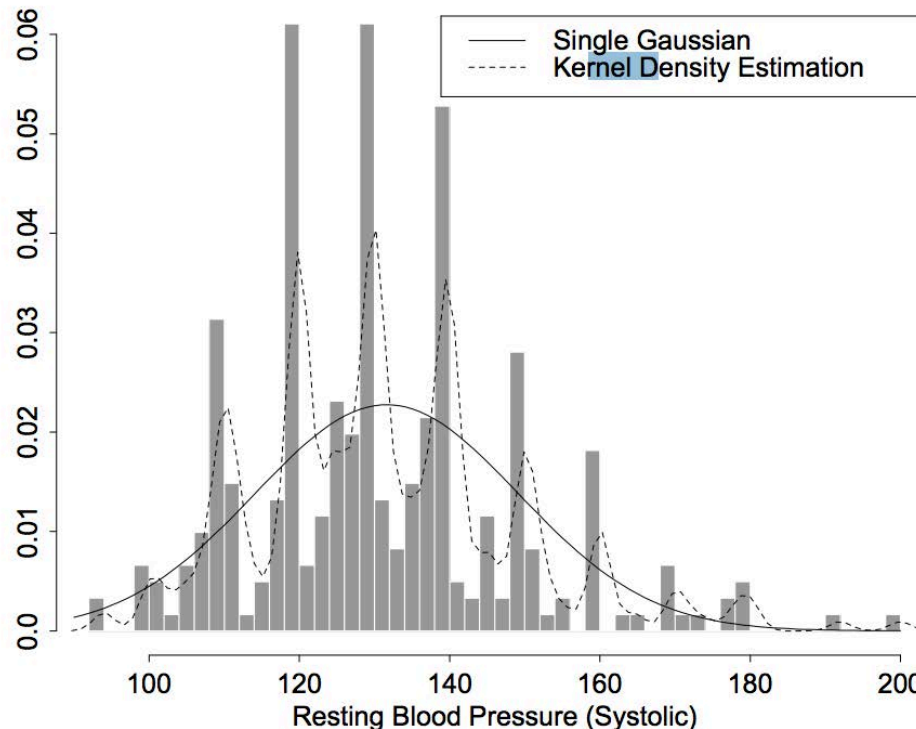
A VARIABLE THAT SEEMS NOT TO FOLLOW NORMAL DISTRIBUTION



A VARIABLE THAT SEEMS TO FOLLOW NORMAL DISTRIBUTION



NORMAL VS. KERNEL DENSITY ESTIMATION



Normal distribution and Gaussian distribution are two names for the same thing.

Figure 3: Systematic measurement errors in the Cleveland heart disease database.



SUMMARY OF NAIVE BAYES

SYRACUSE UNIVERSITY
School of Information Studies

BUILD A NAIVE BAYES CLASSIFIER

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

$P(\text{refund=no}) = 7/10$

$P(\text{refund=yes}) = 3/10$

Prior probability

naive Bayes Classifier:

$P(\text{Refund=Yes|No}) = 3/7$

$P(\text{Refund=No|No}) = 4/7$

$P(\text{Refund=Yes|Yes}) = 0$

$P(\text{Refund=No|Yes}) = 1$

$P(\text{Marital Status=Single|No}) = 2/7$

$P(\text{Marital Status=Divorced|No}) = 1/7$

$P(\text{Marital Status=Married|No}) = 4/7$

$P(\text{Marital Status=Single|Yes}) = 2/3$

$P(\text{Marital Status=Divorced|Yes}) = 1/3$

$P(\text{Marital Status=Married|Yes}) = 0$

For taxable income:

If class=No: sample mean=110

sample variance=2975

If class=Yes: sample mean=90

sample variance=25

Conditional probability

USE NAIVE BAYES CLASSIFIER FOR PREDICTION

Given a test record, calculate posterior probabilities, and choose decision with maximum posterior probabilities.

$$X = (\text{Refund} = \text{No}, \text{Married}, \text{Income} = 120\text{K})$$

$$\begin{aligned} P(X | \text{Class} = \text{No}) &= P(\text{Refund} = \text{No} | \text{Class} = \text{No}) \\ &\quad \times P(\text{Married} | \text{Class} = \text{No}) \\ &\quad \times P(\text{Income} = 120\text{K} | \text{Class} = \text{No}) \\ &= 4/7 \times 4/7 \times 0.0072 = 0.0024 \end{aligned}$$

$$P(\text{Class} = \text{No} | X) = P(X | \text{Class} = \text{No}) \times P(\text{Class} = \text{No}) / P(X)$$

USE NAIVE BAYES CLASSIFIER FOR PREDICTION

Given a test record:

$X = (\text{Refund} = \text{No}, \text{Married}, \text{Income} = 120\text{K})$

$$\begin{aligned} P(X|\text{Class}=\text{No}) &= P(\text{Refund}=\text{No}|\text{Class}=\text{No}) \\ &\quad \times P(\text{Married}|\text{Class}=\text{No}) \\ &\quad \times P(\text{Income}=120\text{K}|\text{Class}=\text{No}) \\ &= 4/7 \times 4/7 \times 0.0072 = 0.0024 \end{aligned}$$

$$P(\text{Class}=\text{No}|X) = P(X|\text{Class}=\text{No}) \times P(\text{Class}=\text{No}) / P(X)$$

$$\begin{aligned} P(X|\text{Class}=\text{Yes}) &= P(\text{Refund}=\text{No}|\text{Class}=\text{Yes}) \\ &\quad \times P(\text{Married}|\text{Class}=\text{Yes}) \\ &\quad \times P(\text{Income}=120\text{K}|\text{Class}=\text{Yes}) \\ &= 1 \times 0 \times 1.2 \times 10^{-9} = 0 \end{aligned}$$

$$P(\text{Class}=\text{Yes}|X) = P(X|\text{Class}=\text{Yes}) \times P(\text{Class}=\text{Yes}) / P(X)$$

Don't forget smoothing!

FEATURES OF BAYESIAN LEARNING METHODS

Each observed training example can incrementally decrease or increase the estimated probability that a hypothesis is correct, and therefore the algorithm is robust to inconsistent examples.

Prior knowledge can be combined with observed data to determine the final probability of a hypothesis.

E.g., an unbalanced coin has 60% chance heads, 40% tails.

Bayesian methods provide probabilistic predictions.

E.g., “This pneumonia patient has a 93% chance of complete recovery.”

CHALLENGE OF BAYESIAN METHODS

Practical difficulty

- Requires initial knowledge of many probabilities

- Estimates the probabilities when they are unknown

- May need to assume normal distribution for continuous variables

Significant cost to compute all probabilities

- Specialized assumptions to reduce the computational cost

 - E.g., naive Bayes is fast.

- Independence assumption may not hold for some attributes

 - Use other techniques, such as Bayesian belief networks (BBN).

Domingos, P., & Pazzani, M. (1997). On the optimality of the simple Bayesian classifier under zero-one loss. *Machine Learning*, 29, 103–30.

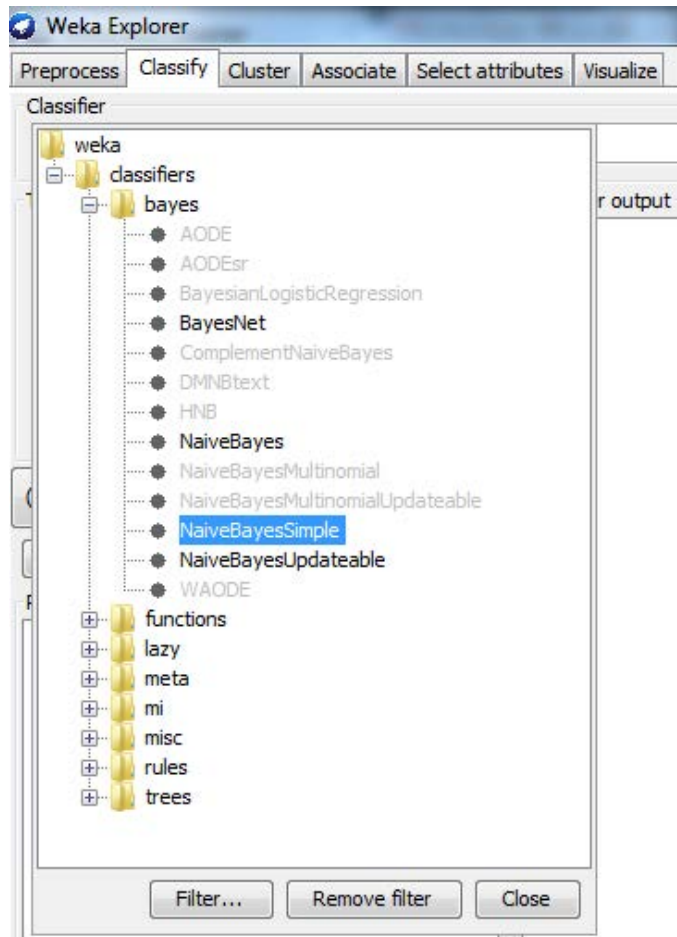
Mitchell, T. (1990). *Machine learning*. New York: McGraw-Hill.



VARIATIONS IN NAIVE BAYES IMPLEMENTATIONS

SYRACUSE UNIVERSITY
School of Information Studies

WEKA'S IMPLEMENTATION OF NAIVE BAYES ALGORITHM



Many different versions of naive Bayes' algorithms.

NaiveBayesSimple

NaiveBayes

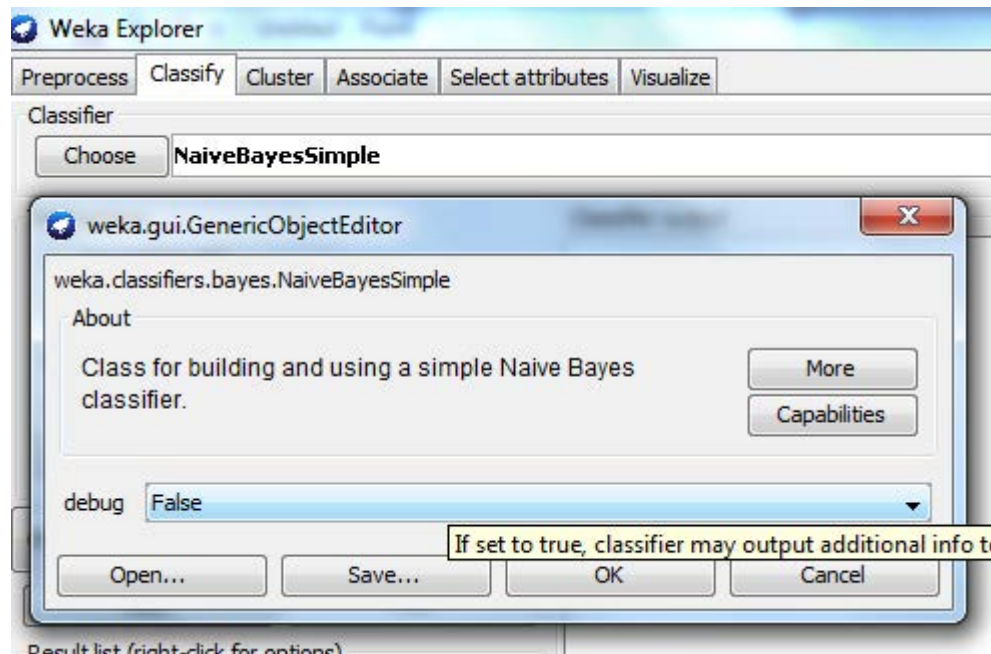
NaiveBayesUpdateable

NaiveBayesMultinomial

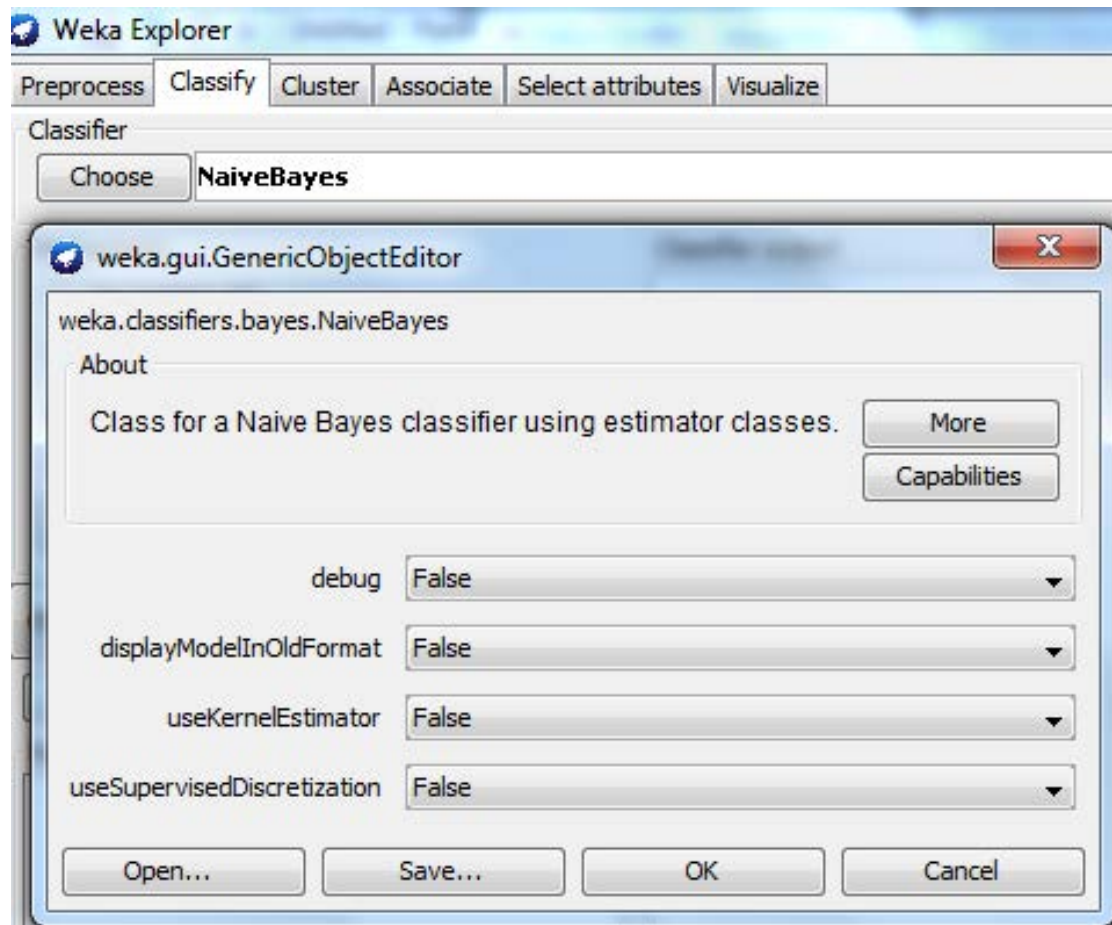
What are their differences?

NAIVE BAYES SIMPLE

NaiveBayesSimple assumes all numeric variables follow normal distribution, and thus does not allow discretization. There are no performance-related parameters to tune in this algorithm.



NAIVE BAYES



NAIVE BAYES

NaiveBayes allows you to choose from three different methods to handle numeric variables.

Method 1: Use default parameter setting; assume they follow normal distribution (same as `NaiveBayesSimple`).

Method 2: Turn on the “`useKernelEstimator`” option to use kernel estimation for numeric variables. The kernel method does not assume normal distribution.

Method 3: Turn on the “`useSupervisedDiscretization`” option to discretize numeric variables.

Which method to choose?

If your numeric variables follow normal distribution, choose method 1 or `NaiveBayesSimple`.

Otherwise, try the other two methods.

NAIVE BAYES UPDATEABLE

Sometimes training data may trickle in, instead of coming in a batch. It is computationally expensive to retrain the entire model whenever a new training example is added to the training set.

NaiveBayesUpdateable can take new training examples and “update” the existing model without complete retraining.

NAIVE BAYES MULTINOMIAL

This algorithm is particularly designed for text categorization. For implementation details, see the supplementary reading [mitchell-nb-text-classifier.pdf](#), an excerpt from Tom Mitchell's textbook on machine learning.

THE E1071 PACKAGE IN R

Library(e1071)

<https://cran.r-project.org/web/packages/e1071/e1071.pdf>

BAYES' THEOREM IN THE NEWS

MIT Technology Review: “How Statisticians Found Air France Flight 447 Two Years After It Crashed Into Atlantic”

<http://www.technologyreview.com/view/527506/how-statisticians-found-air-france-flight-447-two-years-after-it-crashed-into-atlantic/>

NPR News: “Can a 250-Year-Old Mathematical Theorem Find a Missing Plane?”

<http://www.npr.org/blogs/thetwo-way/2014/03/25/294390476/can-a-250-year-old-mathematical-theorem-find-a-missing-plane>