

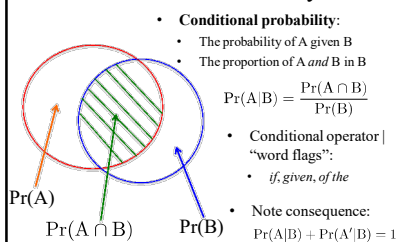
Naïve Bayes Classification

Things We'd Like to Do

- Spam Classification
 - Given an email, predict whether it is spam or not
- Medical Diagnosis
 - Given a list of symptoms, predict whether a patient has disease X or not
- Weather
 - Based on temperature, humidity, etc.... predict if it will rain tomorrow

A Little Bit of Probability Conditional Probability and Bayes' Theorem

Conditional Probability



Example

In a large soil database 72% of the of the samples contain mica and 43% mica and schist. Assuming the database reflective of a relevant population, what is the probability that a randomly selected soil sample (from the same population) that contains mica also contains schist?

$$\Pr(\text{schist}|\text{mica}) = \frac{\Pr(\text{schist} \cap \text{mica})}{\Pr(\text{mica})}$$

Example

$$\Pr(\text{schist}|\text{mica}) = \frac{\Pr(\text{schist} \cap \text{mica})}{\Pr(\text{mica})}$$

$$\Pr(\text{schist} \cap \text{mica}) = 43\%$$

$$\Pr(\text{mica}) = 72\%$$

$$= 0.43/0.72 = 0.597 \approx 60\% \text{ chance}$$

Multiplication Rule

- Another important consequence of conditional probability is the **multiplication rule**:

$$\Pr(A \cap B) = \Pr(A|B)\Pr(B)$$

Statistical Independence

- If A is independent of B then the probability of A is not affected by knowledge of B.
- If A and B are **statistically independent** *if*

$$\Pr(A|B) = \Pr(A)$$

$$\Pr(B|A) = \Pr(B)$$

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

- If A and B do not satisfy the above they are **statistically dependent**

Example

Using the information from the large soil database:

- 72% of the of the samples contain mica
- 43% contain mica and schist.
- 100% contain mica or schist

Compute:

$$\Pr(\text{mica})$$

$$\Pr(\text{schist}|\text{mica})$$

$$\Pr(\text{schist} \cap \text{mica})$$

$$\Pr(\text{schist})$$

$$\Pr(\text{mica}|\text{schist})$$

Example

```
# Data from the question:
M <- 0.72 # Schist
MandS <- 0.43 # Mica and Schist

# Pr(M)
M
[1] 0.72
> SgivenM
[1] 0.5972222
# Pr(S|M) = Pr(SandM)/Pr(M)
SgivenM <- MandS/M
SgivenM
[1] 0.43
# Pr(S and M) = Pr(M and S)
MandS
[1] 0.43
> S
[1] 0.71
> MgivenS
[1] 0.6056338

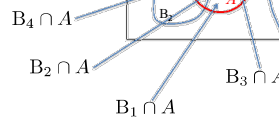
# Pr(S) = Pr(S or M) = Pr(S and M) + Pr(M)
S <- 1 + 0.43 - 0.72
S
[1] 0.6056338

# Pr(M|S) = Pr(M and S)/Pr(S)
MgivenS <- 0.43/S
MgivenS
[1] 0.6956338
```

The Law of Total Probability

- Suppose a sample space can be partitioned into a set of *disjoint* events B_i such that

$$B_1 \cup B_2 \cup \dots \cup B_m = \Omega$$



The Law of Total Probability

- Suppose a sample space can be partitioned into a set of *disjoint* events B_i such that

$$B_1 \cup B_2 \cup \dots \cup B_m = \Omega$$

- The probability of an arbitrary event A in Ω can be written as:

$$\Pr(A) = \Pr(B_1 \cap A) + \Pr(B_2 \cap A) + \Pr(B_3 \cap A) + \dots + \Pr(B_m \cap A)$$

$$= \Pr(A|B_1)\Pr(B_1) + \Pr(A|B_2)\Pr(B_2) + \Pr(A|B_3)\Pr(B_3) + \dots + \Pr(A|B_m)\Pr(B_m)$$

$$= \sum_{i=1}^m \Pr(A|B_i)\Pr(B_i)$$

$$\Pr(A) = \sum_{i=1}^m \Pr(A|B_i)\Pr(B_i) \quad \text{Law of total probability}$$

Example: A medical test

- Professor P LOVES hamburgers. But he's also a hypochondriac. He thinks he is infected with "Mad Cow Disease" (MCD), so he gets himself tested (T).
- The true positive rate of the test is: $\Pr(T^+ | MCD^+) = 0.7$
- The false positive rate of the test is: $\Pr(T^+ | MCD^-) = 0.1$
- The background prevalence of MCD in the yummy cow population is: $\Pr(MCD^+) = 0.02$

What is the probability that Prof. P tests positive for MCD, $\Pr(T^+)$?

Example: A medical test

```
# Data from the question:
Tp.given.MCDp <- 0.7
Tp.given.MCDm <- 0.1
MCDp <- 0.02

# Pr(T+) = Pr(T+ | MCD+) Pr(MCD+) + Pr(T+ | MCD-) Pr(MCD-)
Tp.given.MCDp * MCDp + Tp.given.MCDm * (1-MCDp)

> Tp.given.MCDp * MCDp + Tp.given.MCDm * (1-MCDp)
[1] 0.112
```

There's more than one way to condition: Bayes' Theorem

- Intersection commutes: $A \cap B = B \cap A$
- So: $\Pr(A \cap B) = \Pr(B \cap A)$
- But from the multiplication rule we know:

$$\Pr(A \cap B) = \Pr(A|B)\Pr(B)$$

$$\Pr(B \cap A) = \Pr(B|A)\Pr(A)$$
- So: $\Pr(A|B)\Pr(B) = \Pr(B|A)\Pr(A)$

$$\Pr(A|B) = \frac{\Pr(B|A)\Pr(A)}{\Pr(B)} \quad \text{Bayes' Theorem}$$

Bayes' Theorem

- A slightly more general form for Bayes' Theorem:

$$\Pr(B_i|A) = \frac{\Pr(A|B_i)}{\Pr(A)}$$

- Suppose a sample space can be partitioned into a set of *disjoint* events B_i such that

$$B_1 \cup B_2 \cup \dots \cup B_m = \Omega$$

$$\Pr(B_i|A) = \frac{\Pr(A|B_i)}{\sum_{j=1}^m \Pr(A|B_j)\Pr(B_j)}$$

Bayes' Theorem

Probability

- Bayes' theorem* – allows to calculate conditional probabilities for one variable when conditional probabilities for another variable are known

$$P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)}$$

- Also known as Bayes' rule
- Multiplication rule for the joint distribution is used: $P(X,Y) = P(Y|X)P(X)$
- By symmetry, we also have: $P(Y,X) = P(X|Y)P(Y)$

- The terms are referred to as:

- $P(X)$, the *prior probability*, the initial degree of belief for X
- $P(X|Y)$, the *posterior probability*, the degree of belief after incorporating the knowledge of Y
- $P(Y|X)$, the *likelihood* of Y given X
- $P(Y)$, the *evidence*

- Bayes' theorem: $\text{posterior probability} = \frac{\text{likelihood} \times \text{prior probability}}{\text{evidence}}$

Example: A medical test again...

- Suppose Professor P is positive for MCD. What is the probability that he truly has MCD, $\Pr(MCD^+ | T^+)$?

```
# Data from the question:
Tp.given.MCDp <- 0.7
Tp.given.MCDm <- 0.1
MCDp <- 0.02

# Pr(T+) = Pr(T+ | MCD+) Pr(MCD+) + Pr(T+ | MCD-) Pr(MCD-)
Tp <- Tp.given.MCDp * MCDp + Tp.given.MCDm * (1-MCDp)
Tp

# Pr(MCD+ | T+) = Pr(T+ | MCD+) Pr(MCD+) / Pr(T+)
(Tp.given.MCDp * MCDp) / Tp

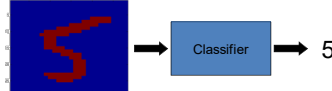
> Tp.given.MCDp * MCDp / Tp
[1] 0.112
[2] 0.325
```

Bayesian Classification

- Problem statement:
 - Given features X_1, X_2, \dots, X_n
 - Predict a label Y

Another Application

Digit Recognition



- $X_1, \dots, X_n \in \{0, 1\}$ (Black vs. White pixels)
- $Y \in \{5, 6\}$ (predict whether a digit is a 5 or a 6)

The Bayes Classifier

- A good strategy is to predict:

$$\arg \max_Y P(Y | X_1, \dots, X_n)$$

- (for example: what is the probability that the image represents a 5 given its pixels?)

- So ... How do we compute that?

The Bayes Classifier

- Use Bayes Rule!

$$P(Y | X_1, \dots, X_n) = \frac{\overset{\text{Likelihood}}{P(X_1, \dots, X_n | Y)} \overset{\text{Prior}}{P(Y)}}{\underset{\text{Normalization Constant}}{P(X_1, \dots, X_n)}}$$

- Why did this help? Well, we think that we might be able to specify how features are "generated" by the class label

The Bayes Classifier

- Let's expand this for our digit recognition task:

$$P(Y = 5 | X_1, \dots, X_n) = \frac{P(X_1, \dots, X_n | Y = 5) P(Y = 5)}{P(X_1, \dots, X_n | Y = 5) P(Y = 5) + P(X_1, \dots, X_n | Y = 6) P(Y = 6)}$$

$$P(Y = 6 | X_1, \dots, X_n) = \frac{P(X_1, \dots, X_n | Y = 6) P(Y = 6)}{P(X_1, \dots, X_n | Y = 5) P(Y = 5) + P(X_1, \dots, X_n | Y = 6) P(Y = 6)}$$

- To classify, we'll simply compute these two probabilities and predict based on which one is greater

Model Parameters

- For the Bayes classifier, we need to "learn" two functions, the likelihood and the prior

- How many parameters are required to specify the prior for our digit recognition example?

Model Parameters

- How many parameters are required to specify the likelihood?
 - (Supposing that each image is 30x30 pixels)

?

Model Parameters

- The problem with explicitly modeling $P(X_1, \dots, X_n | Y)$ is that there are usually way too many parameters:
 - We'll run out of space
 - We'll run out of time
 - And we'll need tons of training data (which is usually not available)

The Naïve Bayes Model

- The *Naïve Bayes Assumption*: Assume that all features are independent **given the class label Y**
- Equationally speaking:

$$P(X_1, \dots, X_n | Y) = \prod_{i=1}^n P(X_i | Y)$$

- (We will discuss the validity of this assumption later)

Why is this useful?

- # of parameters for modeling $P(X_1, \dots, X_n | Y)$:
 - $2(2^n - 1)$
- # of parameters for modeling $P(X_1 | Y), \dots, P(X_n | Y)$
 - $2n$

Naïve Bayes Training

- Now that we've decided to use a Naïve Bayes classifier, we need to train it with some data:



MNIST Training Data

Naïve Bayes Training

- Training in Naïve Bayes is **easy**:
 - Estimate $P(Y=v)$ as the fraction of records with $Y=v$

$$P(Y=v) = \frac{\text{Count}(Y=v)}{\# \text{ records}}$$
 - Estimate $P(X_i=u | Y=v)$ as the fraction of records with $Y=v$ for which $X_i=u$

$$P(X_i=u | Y=v) = \frac{\text{Count}(X_i=u \wedge Y=v)}{\text{Count}(Y=v)}$$
- (This corresponds to Maximum Likelihood estimation of model parameters)

Naïve Bayes Training

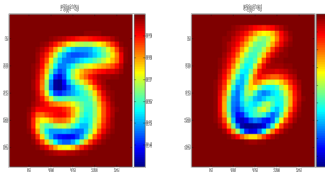
- In practice, some of these counts can be zero
- Fix this by adding "virtual" counts:

$$P(X_i = u | Y = v) = \frac{\text{Count}(X_i = u \wedge Y = v) + 1}{\text{Count}(Y = v) + 2}$$

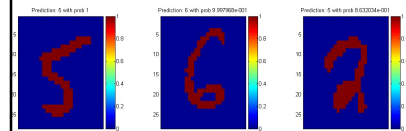
- (This is like putting a prior on parameters and doing MAP estimation instead of MLE)
- This is called *Smoothing*

Naïve Bayes Training

- For binary digits, training amounts to averaging all of the training fives together and all of the training sixes together.



Naïve Bayes Classification



Another Example of the Naïve Bayes Classifier

The weather data, with counts and probabilities													
outlook	temperature					humidity		windy		play			
	yes	no		yes	no		yes	no	yes	no	yes	no	
sunny	2	3	hot	2	2	high	3	4	false	6	2	9	5
overcast	4	0	mild	4	2	normal	6	1	true	3	3		
rainy	3	2	cool	3	1								
sunny	2/9	3/5	hot	2/9	2/5	high	3/9	4/5	false	6/9	2/5	9/14	5/14
overcast	4/9	0/5	mild	4/9	2/5	normal	6/9	1/5	true	3/9	3/5		
rainy	3/9	2/5	cool	3/9	1/5								

A new day				
outlook	temperature	humidity	windy	play
sunny	cool	high	true	?

- Likelihood of yes

$$= \frac{2}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{9}{14} = 0.0053$$
- Likelihood of no

$$= \frac{3}{5} \times \frac{1}{5} \times \frac{4}{5} \times \frac{3}{5} \times \frac{5}{14} = 0.0206$$
- Therefore, the prediction is No

The Naïve Bayes Classifier for Data Sets with Numerical Attribute Values

- One common practice to handle numerical attribute values is to assume normal distributions for numerical attributes.

The numeric weather data with summary statistics													
outlook	temperature				humidity		windy		play				
	yes	no	yes	no	yes	no	yes	no	yes	no	yes	no	
sunny	2	3	83	85	86	85	false	6	2	9	5		
overcast	4	0	70	80	96	90	true	3	3				
rainy	3	2	68	65	80	70							
			64	72	65	95							
			69	71	70	91							
			75		80								
			75		70								
			72		90								
			81		75								
sunny	2/9	3/5	mean	73	74.6	mean	79.1	86.2	false	6/9	2/5	9/14	5/14
overcast	4/9	0/5	std	6.2	7.9	std	10.2	9.7	true	3/9	3/5		
			dev			dev							
rainy	3/9	2/5											

- Let x_1, x_2, \dots, x_n be the values of a numerical attribute in the training data set.

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\sigma = \frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)^2$$

$$f(w) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(w-\mu)^2}{\sigma^2}}$$

- For examples,

$$f(\text{temperature} = 66 | \text{Yes}) = \frac{1}{\sqrt{2\pi}(6.2)} e^{-\frac{(66-74.6)^2}{2 \times 6.2^2}} = 0.0340$$

- Likelihood of Yes = $\frac{2}{9} \times 0.0340 \times 0.0221 \times \frac{3}{9} \times \frac{9}{14} = 0.000036$

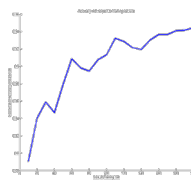
- Likelihood of No = $\frac{3}{5} \times 0.0291 \times 0.038 \times \frac{3}{5} \times \frac{5}{14} = 0.000136$

Outputting Probabilities

- What's nice about Naïve Bayes (and generative models in general) is that it returns probabilities
 - These probabilities can tell us how confident the algorithm is
 - So... don't throw away those probabilities!

Performance on a Test Set

- Naïve Bayes is often a good choice if you don't have much training data!



Naïve Bayes Assumption

- Recall the Naïve Bayes assumption:
 - that all features are independent **given the class label Y**
- Does this hold for the digit recognition problem?

Exclusive-OR Example

- For an example where conditional independence fails:
 - $Y = \text{XOR}(X_1, X_2)$

X_1	X_2	$P(Y=0 X_1, X_2)$	$P(Y=1 X_1, X_2)$
0	0	1	0
0	1	0	1
1	0	0	1
1	1	1	0

- Actually, the Naïve Bayes assumption is almost never true

- Still... Naïve Bayes often performs surprisingly well even when its assumptions do not hold

Numerical Stability

- It is often the case that machine learning algorithms need to work with very small numbers
 - Imagine computing the probability of 2000 independent coin flips
 - MATLAB thinks that $(.5)^{2000} = 0$

Underflow Prevention

- Multiplying lots of probabilities
→ floating-point underflow.
- Recall: $\log(xy) = \log(x) + \log(y)$,
→ better to sum logs of probabilities rather than multiplying probabilities.

Underflow Prevention

- Class with highest final un-normalized log probability score is still the most probable.

$$c_{NB} = \operatorname{argmax}_{c_j \in C} \log P(c_j) + \sum_{i \text{ positions}} \log P(x_i | c_j)$$

Numerical Stability

- Instead of comparing $P(Y=5|X_1, \dots, X_n)$ with $P(Y=6|X_1, \dots, X_n)$,
– Compare their logarithms

$$\begin{aligned} \log(P(Y|X_1, \dots, X_n)) &= \log\left(\frac{P(X_1, \dots, X_n|Y) \cdot P(Y)}{P(X_1, \dots, X_n)}\right) \\ &= \text{constant} + \log\left(\prod_{i=1}^n P(X_i|Y)\right) + \log P(Y) \\ &= \text{constant} + \sum_{i=1}^n \log P(X_i|Y) + \log P(Y) \end{aligned}$$

Recovering the Probabilities

- What if we want the probabilities though??
- Suppose that for some constant K , we have:
– And $\log P(Y=5|X_1, \dots, X_n) + K$
 $\log P(Y=6|X_1, \dots, X_n) + K$
- How would we recover the original probabilities?

Recovering the Probabilities

- Given: $\alpha_i = \log a_i + K$
- Then for any constant C :
$$\begin{aligned} \frac{a_i}{\sum_j a_j} &= \frac{e^{\alpha_i}}{\sum_j e^{\alpha_j}} \\ &= \frac{e^{\alpha_i + C}}{\sum_j e^{\alpha_j + C}} \\ &= \frac{e^{\alpha_i}}{\sum_j e^{\alpha_j}} \end{aligned}$$
- One suggestion: set C such that the greatest α_i is shifted to zero:
$$C = -\max_i \{\alpha_i\}$$

See <https://stackoverflow.com/questions/100822/example-of-how-the-log-sum-exp-trick-works-in-naive-bayes> for more info

Recap

- We defined a *Bayes classifier* but saw that it's intractable to compute $P(X_1, \dots, X_n|Y)$
- We then used the *Naive Bayes assumption* – that everything is independent given the class label Y
- A natural question: is there some happy compromise where we only assume that *some* features are conditionally independent?
– Stay Tuned...

Conclusions

- Naive Bayes is:
– Really easy to implement and often works well
– Often a good first thing to try
– Commonly used as a “punching bag” for smarter algorithms

Evaluating classification algorithms

- You have designed a new classifier.
- You give it to me, and I try it on my image dataset

Evaluating classification algorithms

- I tell you that it achieved 95% accuracy on my data.
- Is your technique a success?

Types of errors

- But suppose that
 - The 95% is the correctly classified pixels
 - Only 5% of the pixels are actually edges
 - It misses all the edge pixels
- How do we count the effect of different types of error?

Types of errors

		Prediction	
		Edge	Not edge
Ground Truth	Edge	True Positive	False Negative
	Not Edge	False Positive	True Negative

Two parts to each: whether you got it correct or not, and what you guessed. For example for a particular pixel, our guess might be labelled...

True Positive
 Did we get it correct? True, we did get it correct.
 What did we say? We said 'positive', i.e. edge.
 or maybe it was labelled as one of the others, maybe...

False Negative
 Did we get it correct? False, we did not get it correct.
 What did we say? We said 'negative', i.e. not edge.

Sensitivity and Specificity

Count up the total number of each label (TP, FP, TN, FN) over a large dataset. In ROC analysis, we use two statistics:

Sensitivity = $\frac{TP}{TP+FN}$
 Can be thought of as the likelihood of spotting a positive case when presented with one.
 Or... the proportion of edges we find.

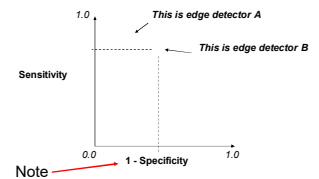
Specificity = $\frac{TN}{TN+FP}$
 Can be thought of as the likelihood of spotting a negative case when presented with one.
 Or... the proportion of non-edges that we find

$$\text{Sensitivity} = \frac{TP}{TP+FN} = ? \quad \text{Specificity} = \frac{TN}{TN+FP} = ?$$

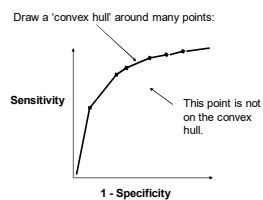
		Prediction	
		1	0
Ground Truth	1	60	30
	0	80	20

60+30 = 90 cases in the dataset were class 1 (edge)
 80+20 = 100 cases in the dataset were class 0 (non-edge)
 90+100 = 190 examples (pixels) in the data overall

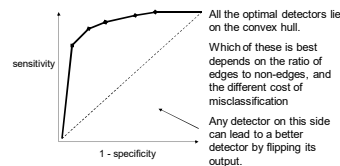
The ROC space



The ROC Curve



ROC Analysis



Take-home point: You should always quote sensitivity and specificity for your algorithm, if possible plotting an ROC graph. Remember also though, any statistic you quote should be an average over a suitable range of tests for your algorithm.

Holdout estimation

- What to do if the amount of data is limited?
- The *holdout* method reserves a certain amount for testing and uses the remainder for training
 → Usually: one third for testing, the rest for training

Holdout estimation

- Problem: the samples might not be representative
 - Example: class might be missing in the test data
- Advanced version uses *stratification*
 - Ensures that each class is represented with approximately equal proportions in both subsets

Repeated holdout method

- Repeat process with different subsamples
 - more reliable
 - In each iteration, a certain proportion is randomly selected for training (possibly with stratification)
 - The error rates on the different iterations are averaged to yield an overall error rate

Repeated holdout method

- Still not optimum: the different test sets overlap
 - Can we prevent overlapping?
 - Of course!

Cross-validation

- *Cross-validation* avoids overlapping test sets
 - First step: split data into k subsets of equal size
 - Second step: use each subset in turn for testing, the remainder for training
- Called *k-fold cross-validation*

Cross-validation

- Often the subsets are stratified before the cross-validation is performed
- The error estimates are averaged to yield an overall error estimate

More on cross-validation

- Standard method for evaluation: stratified ten-fold cross-validation
- Why ten?
 - Empirical evidence supports this as a good choice to get an accurate estimate
 - There is also some theoretical evidence for this
- Stratification reduces the estimate's variance
- Even better: repeated stratified cross-validation
 - E.g. ten-fold cross-validation is repeated ten times and results are averaged (reduces the variance)

Leave-One-Out cross-validation

- Leave-One-Out: a particular form of cross-validation:
 - Set number of folds to number of training instances
 - I.e., for n training instances, build classifier n times
- Makes best use of the data
- Involves no random subsampling
- Very computationally expensive
 - (exception: NN)

Leave-One-Out-CV and stratification

- Disadvantage of Leave-One-Out-CV: stratification is not possible
 - It *guarantees* a non-stratified sample because there is only one instance in the test set!

Hands-on Example

```
# Import Bayes.csv from class webpage

# Select training data
traindata <- Bayes[1:14,]

# Select test data
testdata <- Bayes[15,]
```


Construct Naïve Bayes Classifier the hard way

```
# Calculate the Prior for Play
Pplay <- table(traindata$Play)
Pplay <- Pplay/sum(Pplay)

# Calculate P(Sunny | Play)
sunny <- table(traindata[,c("Play", "Sunny")])
sunny <- sunny/rowSums(sunny)
```

```
# Calculate P(Hot | Play)
hot <- table(traindata[,c("Play", "Hot")])
hot <- hot/rowSums(hot)

# and Calculate P(Windy | Play)
windy <- table(traindata[,c("Play", "Windy")])
windy <- windy/rowSums(windy)
```

```
# Evaluate testdata
Pyes <- sunny["Yes", "Yes"] * hot["Yes", "No"] * windy["Yes", "Yes"]

Pno <- sunny["No", "Yes"] * hot["No", "No"] * windy["No", "Yes"]

# Do we play or not?
Max(Pyes, Pno)
```

```
# Do it again, but use the naiveBayes package

# install the package if you don't already have it
install.packages("e1071")

# load package
library(e1071)

# train model
m <- naiveBayes(traindata[,1:3], traindata[,4])

# evaluate testdata
predict(m, testdata[,1:3])
```

```
# use the naiveBayes classifier on the iris data
m <- naiveBayes(iris[,1:4], iris[,5])
table(predict(m, iris[,1:4]), iris[,5])
```

Questions?