



# Data Mining with R

วิชา การค้นพบองค์ความรู้และการทำเหมืองข้อมูลขั้นสูง

Veerasak Kritsanapraphan

Chulalongkorn University

Email : [veerasak.kr568@cbs.chula.ac.th](mailto:veerasak.kr568@cbs.chula.ac.th)

 @veerasakk

# Slide and Sample Data

[https://github.com/vkrit/chula\\_datamining](https://github.com/vkrit/chula_datamining)



# Agenda

- Overview and Data Visualization
- Data Preparation
- Predictive Data Mining
  - Decision Tree
  - K-Nearest Neighbor
  - Naive Bayes Classifier
  - Neural Network

# Slide and Source Codes

[https://github.com/vkrit/  
chula\\_datamining](https://github.com/vkrit/chula_datamining)



# Overview

- Predictive Data Mining
  - Two Phases of Processing
    - Training Phase : Learn a model from training data
    - Predicting Phase : Deploy the model to production and use that to predict the future outcome

# Data

- Iris Data Set from UCI Machine Learning Repository (<https://archive.ics.uci.edu/ml/datasets/Iris>)

Attribute Information:

1. Sepal Length in cm

2. Sepal width in cm

3. Petal length in cm

4. Petal length in cm

5. Classes:

- Iris Setosa

- Iris Versicolour

- Iris Virginica



# Getting Data

```
> iris <- read.csv("iris.data.csv", header=TRUE)
```

```
> head(iris)
  Sepal.Length Sepal.Width Petal.Length Petal.Width Species
1          5.1         3.5          1.4          0.2  setosa
2          4.9         3.0          1.4          0.2  setosa
3          4.7         3.2          1.3          0.2  setosa
4          4.6         3.1          1.5          0.2  setosa
5          5.0         3.6          1.4          0.2  setosa
6          5.4         3.9          1.7          0.4  setosa
> nrow(iris)
[1] 150
> table(iris$Species)

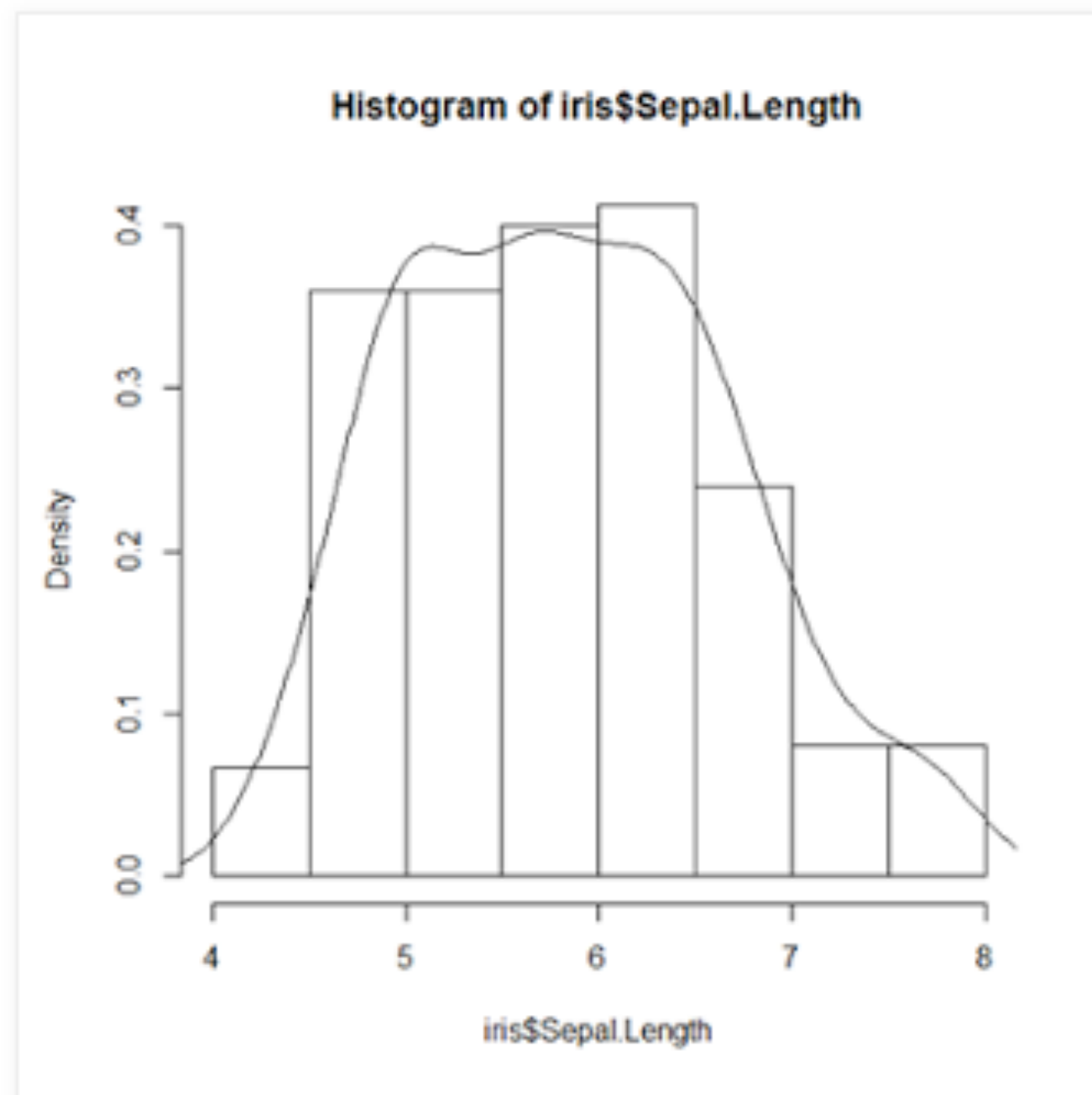
  setosa versicolor virginica 
     50         50         50 
>
```

# Data Visualization

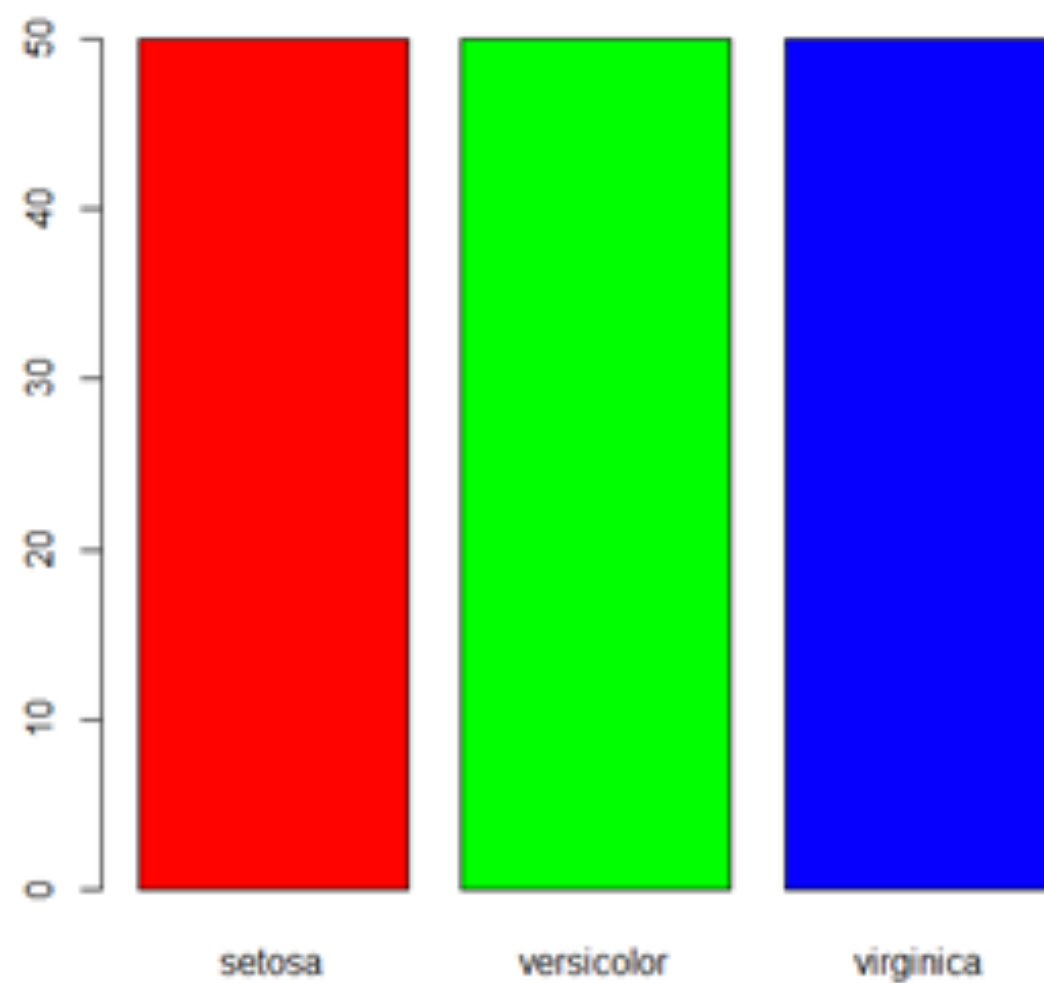
- Visualizing existing data is a very useful way to come up with ideas about what features should be included.
- "Dataframe" in R is a common way where data samples are organized in a tabular structure.



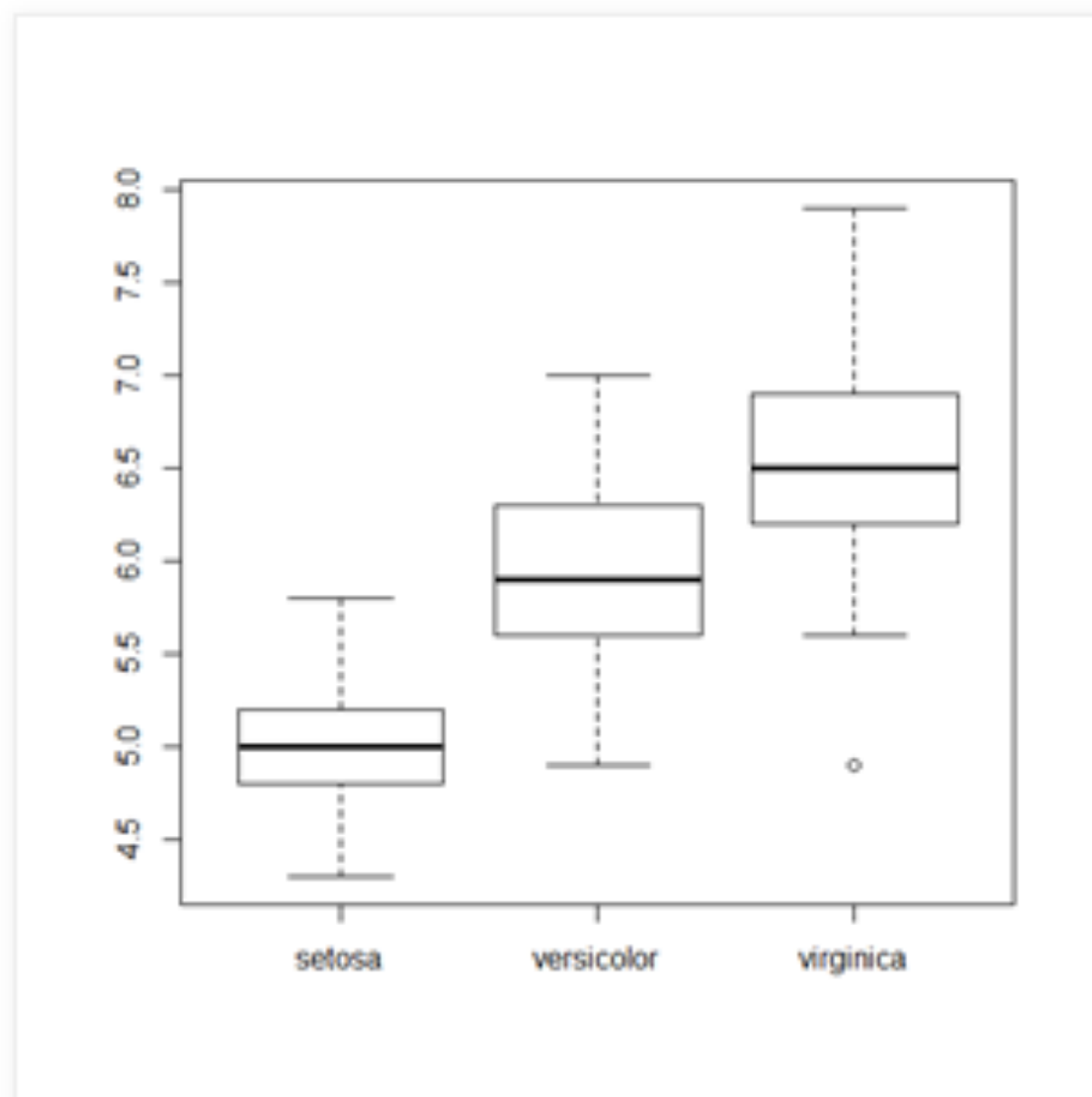
```
> # Plot the histogram  
> hist(iris$Sepal.Length, breaks=10, prob=T)  
> # Plot the density curve  
> lines(density(iris$Sepal.Length))  
>
```



```
> categories <- table(iris$Species)
> barplot(categories, col=c('red', 'green', 'blue'))
>
```



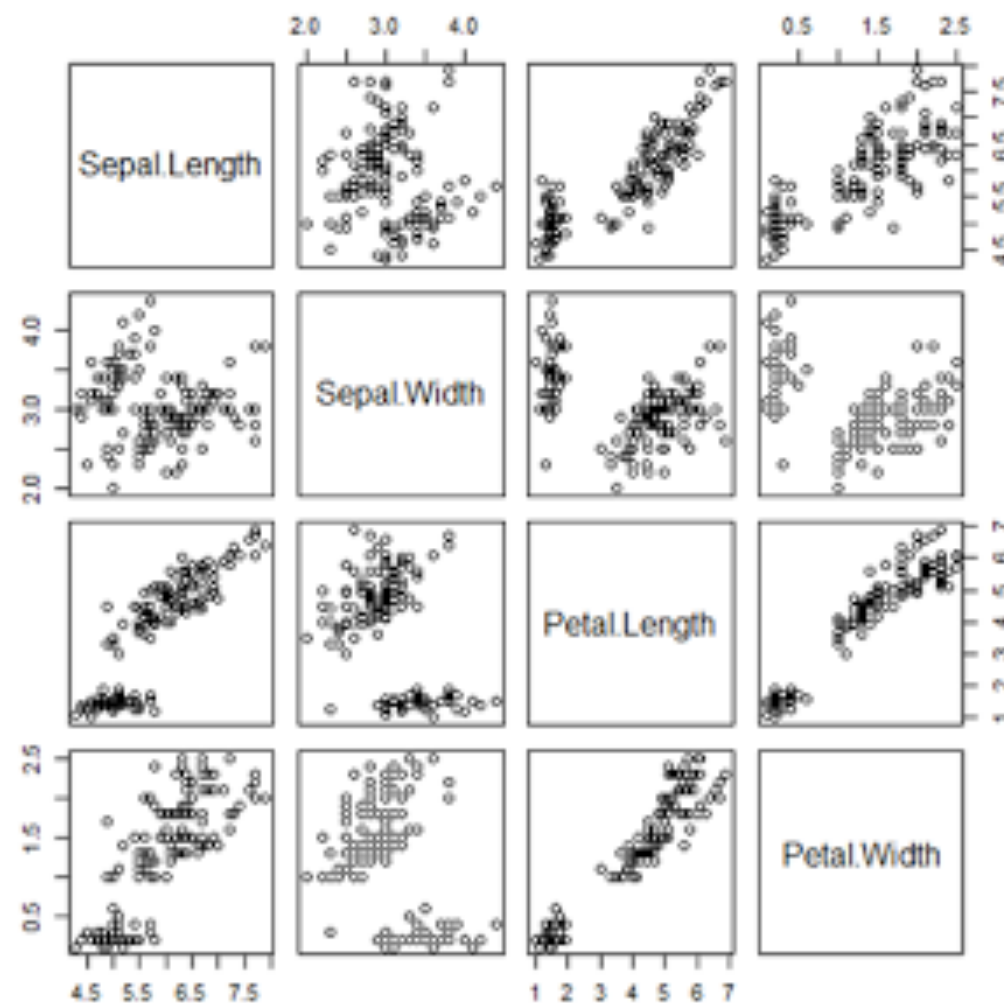
```
> boxplot(Sepal.Length~Species, data=iris)  
>
```



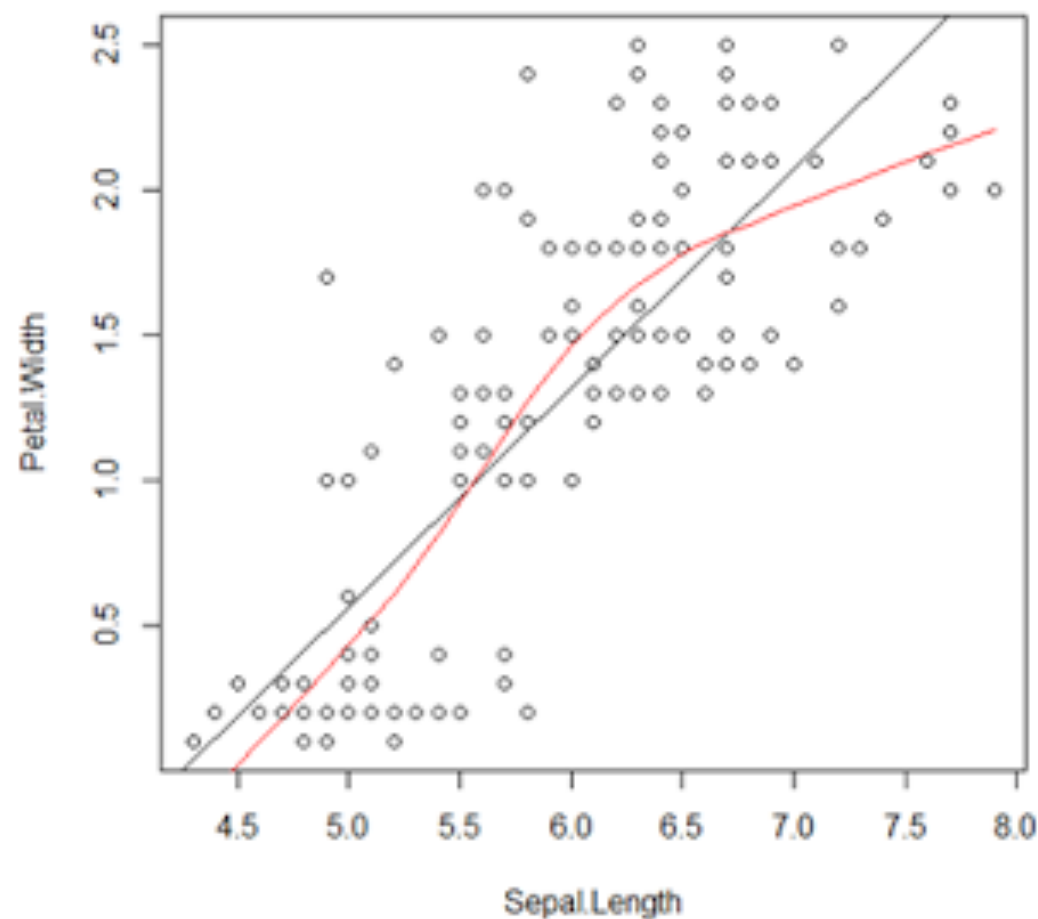
```

> # Scatter plot for all pairs
> pairs(iris[,c(1,2,3,4)])
> # Compute the correlation matrix
> cor(iris[,c(1,2,3,4)])
      Sepal.Length Sepal.Width Petal.Length Petal.Width
Sepal.Length      1.0000000 -0.1170695  0.8716902  0.8179410
Sepal.Width       -0.1170695  1.0000000 -0.4284401 -0.3661259
Petal.Length       0.8716902 -0.4284401  1.0000000  0.9628654
Petal.Width        0.8179410 -0.3661259  0.9628654  1.0000000
>

```



```
> # First plot the 2 variables
> plot(Petal.Width~Sepal.Length, data=iris)
> # Learn the regression model
> model <- lm(Petal.Width~Sepal.Length, data=iris)
> # Plot the regression line
> abline(model)
> # Now learn the local linear model
> model2 <- lowess(iris$Petal.Width~iris$Sepal.Length)
> lines(model2, col="red")
>
```



# Preparing Training Data

- At this step, the purpose is to transform the raw data in a form that can fit into the data mining model.
  - Data sampling
  - Data validation and handle missing data
  - Normalize numeric value into a uniform range
  - Compute aggregated value (a special case is to compute frequency counts)
  - Expand categorical field to binary fields
  - Discretize numeric value into categories
  - Create derived fields from existing fields
  - Reduce dimensionality
  - Power and Log transformation

# Data Sampling

```
> # select 10 records out from iris with replacement
> index <- sample(1:nrow(iris), 10, replace=T)
> index
[1] 133  36 107 140  66  67  36   3  97  37
> irissample <- iris[index,]
> irissample
```

	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
133	6.4	2.8	5.6	2.2	virginica
36	5.0	3.2	1.2	0.2	setosa
107	4.9	2.5	4.5	1.7	virginica
140	6.9	3.1	5.4	2.1	virginica
66	6.7	3.1	4.4	1.4	versicolor
67	5.6	3.0	4.5	1.5	versicolor
36.1	5.0	3.2	1.2	0.2	setosa
3	4.7	3.2	1.3	0.2	setosa
97	5.7	2.9	4.2	1.3	versicolor
37	5.5	3.5	1.3	0.2	setosa

```
>
```

# Impute missing data

- Discard the whole record
- Infer missing value based on the data of other record. Approach is to fill the missing data with the average or the median.

```
> # Create some missing data
> irissample[10, 1] <- NA
> irissample
```

	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
133	6.4	2.8	5.6	2.2	virginica
36	5.0	3.2	1.2	0.2	setosa
107	4.9	2.5	4.5	1.7	virginica
140	6.9	3.1	5.4	2.1	virginica
66	6.7	3.1	4.4	1.4	versicolor
67	5.6	3.0	4.5	1.5	versicolor
36.1	5.0	3.2	1.2	0.2	setosa
3	4.7	3.2	1.3	0.2	setosa
97	5.7	2.9	4.2	1.3	versicolor
37	NA	3.5	1.3	0.2	setosa



```

> library(e1071)
Loading required package: class
Warning message:
package 'e1071' was built under R version 2.14.2
> fixIris1 <- impute(irissample[,1:4], what='mean')
> fixIris1
      Sepal.Length Sepal.Width Petal.Length Petal.Width
133      6.400000      2.8      5.6      2.2
36       5.000000      3.2      1.2      0.2
107      4.900000      2.5      4.5      1.7
140      6.900000      3.1      5.4      2.1
66       6.700000      3.1      4.4      1.4
67       5.600000      3.0      4.5      1.5
36.1     5.000000      3.2      1.2      0.2
3        4.700000      3.2      1.3      0.2
97       5.700000      2.9      4.2      1.3
37       5.655556      3.5      1.3      0.2
> fixIris2 <- impute(irissample[,1:4], what='median')
> fixIris2
      Sepal.Length Sepal.Width Petal.Length Petal.Width
133      6.4      2.8      5.6      2.2
36       5.0      3.2      1.2      0.2
107      4.9      2.5      4.5      1.7
140      6.9      3.1      5.4      2.1
66       6.7      3.1      4.4      1.4
67       5.6      3.0      4.5      1.5
36.1     5.0      3.2      1.2      0.2
3        4.7      3.2      1.3      0.2
97       5.7      2.9      4.2      1.3
37       5.6      3.5      1.3      0.2
>

```

# Normalize numeric value

```
> # scale the columns
> # x-mean(x)/standard deviation
> scaleiris <- scale(iris[, 1:4])
> head(scaleiris)
      Sepal.Length Sepal.Width Petal.Length Petal.Width
[1,]   -0.8976739    1.01560199   -1.335752   -1.311052
[2,]   -1.1392005   -0.13153881   -1.335752   -1.311052
[3,]   -1.3807271    0.32731751   -1.392399   -1.311052
[4,]   -1.5014904    0.09788935   -1.279104   -1.311052
[5,]   -1.0184372    1.24503015   -1.335752   -1.311052
[6,]   -0.5353840    1.93331463   -1.165809   -1.048667
>
```

# Reduce dimensionality

There are two ways to reduce the number of input attributes.

1. Removing irrelevant input variables.
2. Removing redundant input variables.

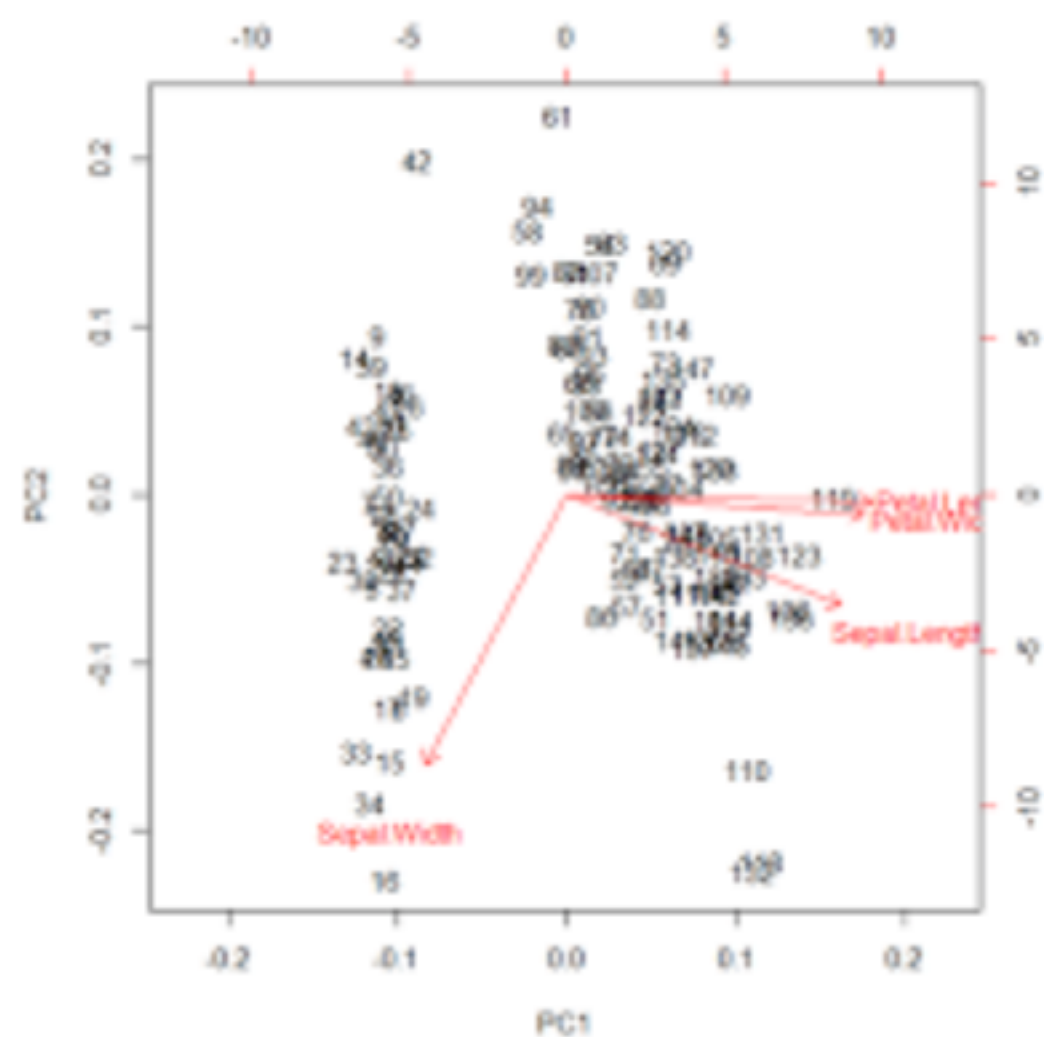
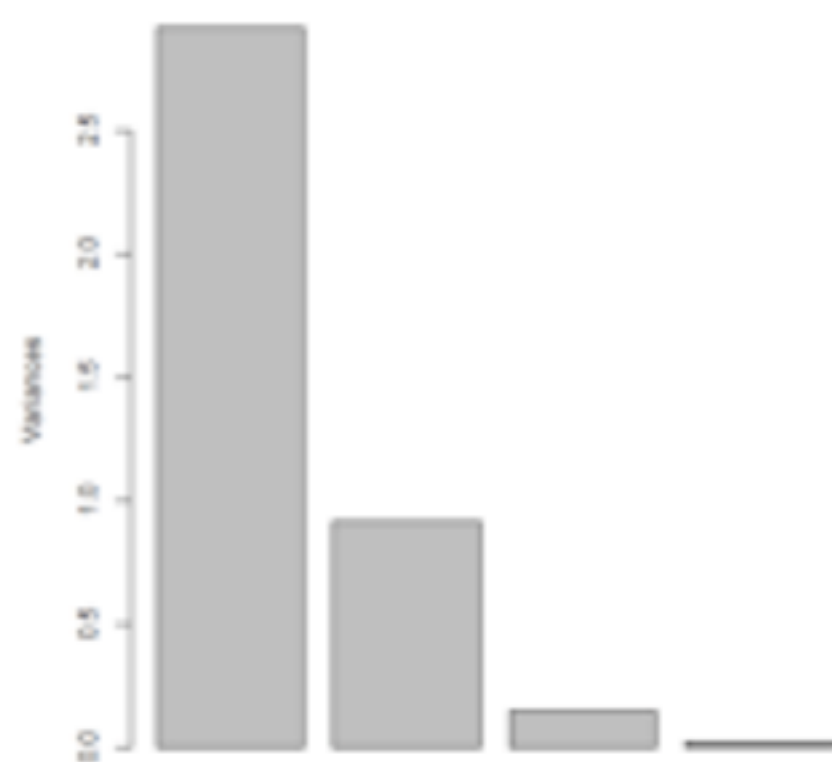
```

> # Use iris data set
> cor(iris[, -5])
      Sepal.Length Sepal.Width Petal.Length Petal.Width
Sepal.Length  1.00000000000 -0.1175697841  0.8717537759  0.8179411263
Sepal.Width   -0.1175697841  1.00000000000 -0.4284401043 -0.3661259325
Petal.Length   0.8717537759 -0.4284401043  1.00000000000  0.9628654314
Petal.Width    0.8179411263 -0.3661259325  0.9628654314  1.00000000000
> # Some attributes shows high correlation, compute PCA
> pca <- prcomp(iris[, -5], scale=T)
> summary(pca)
Importance of components:
              PC1          PC2          PC3          PC4
Standard deviation  1.708361  0.9560494  0.3830886  0.1439265
Proportion of Variance 0.729620 0.2285100 0.0366900 0.0051800
Cumulative Proportion 0.729620 0.9581300 0.9948200 1.0000000
> # Notice PC1 and PC2 covers most variation
> plot(pca)
> pca$rotation
              PC1          PC2          PC3          PC4
Sepal.Length  0.5210659147 -0.37741761556  0.7195663527  0.2612862800
Sepal.Width   -0.2693474425 -0.92329565954 -0.2443817795 -0.1235096196
Petal.Length   0.5804130958 -0.02449160909 -0.1421263693 -0.8014492463
Petal.Width    0.5648565358 -0.06694198697 -0.6342727371  0.5235971346
> # Project first 2 records in PCA direction
> predict(pca)[1:2,]
              PC1          PC2          PC3          PC4
[1,] -2.257141176 -0.4784238321  0.1272796237  0.02408750846
[2,] -2.074013015  0.6718826870  0.2338255167  0.10266284468
> # plot all points in top 2 PCA direction
> biplot(pca)

```

Pin it

pca





# Add derived attributes

```
> iris2 <- transform(iris, ratio=round(Sepal.Length/Sepal.Width, 2))  
> head(iris2)
```

	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species	ratio
1	5.1	3.5	1.4	0.2	setosa	1.46
2	4.9	3.0	1.4	0.2	setosa	1.63
3	4.7	3.2	1.3	0.2	setosa	1.47
4	4.6	3.1	1.5	0.2	setosa	1.48
5	5.0	3.6	1.4	0.2	setosa	1.39
6	5.4	3.9	1.7	0.4	setosa	1.38

# Discretize numeric value into categories

```
> # Equal width cuts
> segments <- 10
> maxL <- max(iris$Petal.Length)
> minL <- min(iris$Petal.Length)
> theBreaks <- seq(minL, maxL,
                   by=(maxL-minL)/segments)
> cutPetalLength <- cut(iris$Petal.Length,
                       breaks=theBreaks,
                       include.lowest=T)
> newdata <- data.frame(orig.Petal.Len=iris$Petal.Length,
                       cut.Petal.Len=cutPetalLength)
> head(newdata)
  orig.Petal.Len cut.Petal.Len
1           1.4      [1,1.59]
2           1.4      [1,1.59]
3           1.3      [1,1.59]
4           1.5      [1,1.59]
5           1.4      [1,1.59]
6           1.7    (1.59,2.18]
>
> # Constant frequency / height
> myBreaks <- quantile(iris$Petal.Length,
                      probs=seq(0,1,1/segments))
> cutPetalLength2 <- cut(iris$Petal.Length,
                       breaks=myBreaks,
                       include.lowest=T)
> newdata2 <- data.frame(orig.Petal.Len=iris$Petal.Length,
                       cut.Petal.Len=cutPetalLength2)
> head(newdata2)
  orig.Petal.Len cut.Petal.Len
1           1.4      [1,1.4]
2           1.4      [1,1.4]
3           1.3      [1,1.4]
4           1.5    (1.4,1.5]
5           1.4      [1,1.4]
6           1.7    (1.7,3.9]
>
```

# Binarize categorical attributes

```
> cat <- levels(iris$Species)
> cat
[1] "setosa"      "versicolor" "virginica"
> binarize <- function(x) {return(iris$Species == x)}
> newcols <- sapply(cat, binarize)
> colnames(newcols) <- cat
> data <- cbind(iris[,c('Species')], newcols)
> data[45:55,]
```

		setosa	versicolor	virginica
[1,]	1	1	0	0
[2,]	1	1	0	0
[3,]	1	1	0	0
[4,]	1	1	0	0
[5,]	1	1	0	0
[6,]	1	1	0	0
[7,]	2	0	1	0
[8,]	2	0	1	0
[9,]	2	0	1	0
[10,]	2	0	1	0
[11,]	2	0	1	0



# **Data Mining**

Techniques

# Iris Data Preparation

```
> set.seed(1234)
```

```
> ind <- sample(2, nrow(iris), replace=TRUE, prob=c(0.7, 0.3))
```

```
> trainData <- iris[ind==1,]
```

```
> testData <- iris[ind==2,]
```

# Decision Tree

## Conditional Inference Trees

### Description

Recursive partitioning for continuous, censored, ordered, nominal and multivariate response variables in a conditional inference framework.

### Usage

```
ctree(formula, data, subset = NULL, weights = NULL,  
      controls = ctree_control(), xtrafo = ptrrafo, ytrafo = ptrrafo,  
      scores = NULL)
```

### Arguments

- formula** a symbolic description of the model to be fit. Note that symbols like `:` and `-` will not work and the tree will make use of all variables listed on the rhs of `formula`.
- data** a data frame containing the variables in the model.
- subset** an optional vector specifying a subset of observations to be used in the fitting process.
- weights** an optional vector of weights to be used in the fitting process. Only non-negative integer valued weights are allowed.
- controls** an object of class [TreeControl](#), which can be obtained using [ctree\\_control](#).
- xtrafo** a function to be applied to all input variables. By default, the [ptrrafo](#) function is applied.
- ytrafo** a function to be applied to all response variables. By default, the [ptrrafo](#) function is applied.
- scores** an optional named list of scores to be attached to ordered factors.

# Decision Tree - Create Model

```
> library(party)
> myFormula <- Species ~ Sepal.Length + Sepal.Width + Petal.Length + Petal.Width
> iris_ctree <- ctree(myFormula, data=trainData)
> # check the prediction
> table(predict(iris_ctree), trainData$Species)
```

	setosa	versicolor	virginica
setosa	40	0	0
versicolor	0	37	3
virginica	0	1	31

```
> print(iris_ctree)
```

```
Conditional inference tree with 4 terminal nodes
```

```
Response: Species
```

```
Inputs: Sepal.Length, Sepal.Width, Petal.Length, Petal.Width
```

```
Number of observations: 112
```

```
1) Petal.Length <= 1.9; criterion = 1, statistic = 104.643
```

```
2)* weights = 40
```

```
1) Petal.Length > 1.9
```

```
3) Petal.Width <= 1.7; criterion = 1, statistic = 48.939
```

```
4) Petal.Length <= 4.4; criterion = 0.974, statistic = 7.397
```

```
5)* weights = 21
```

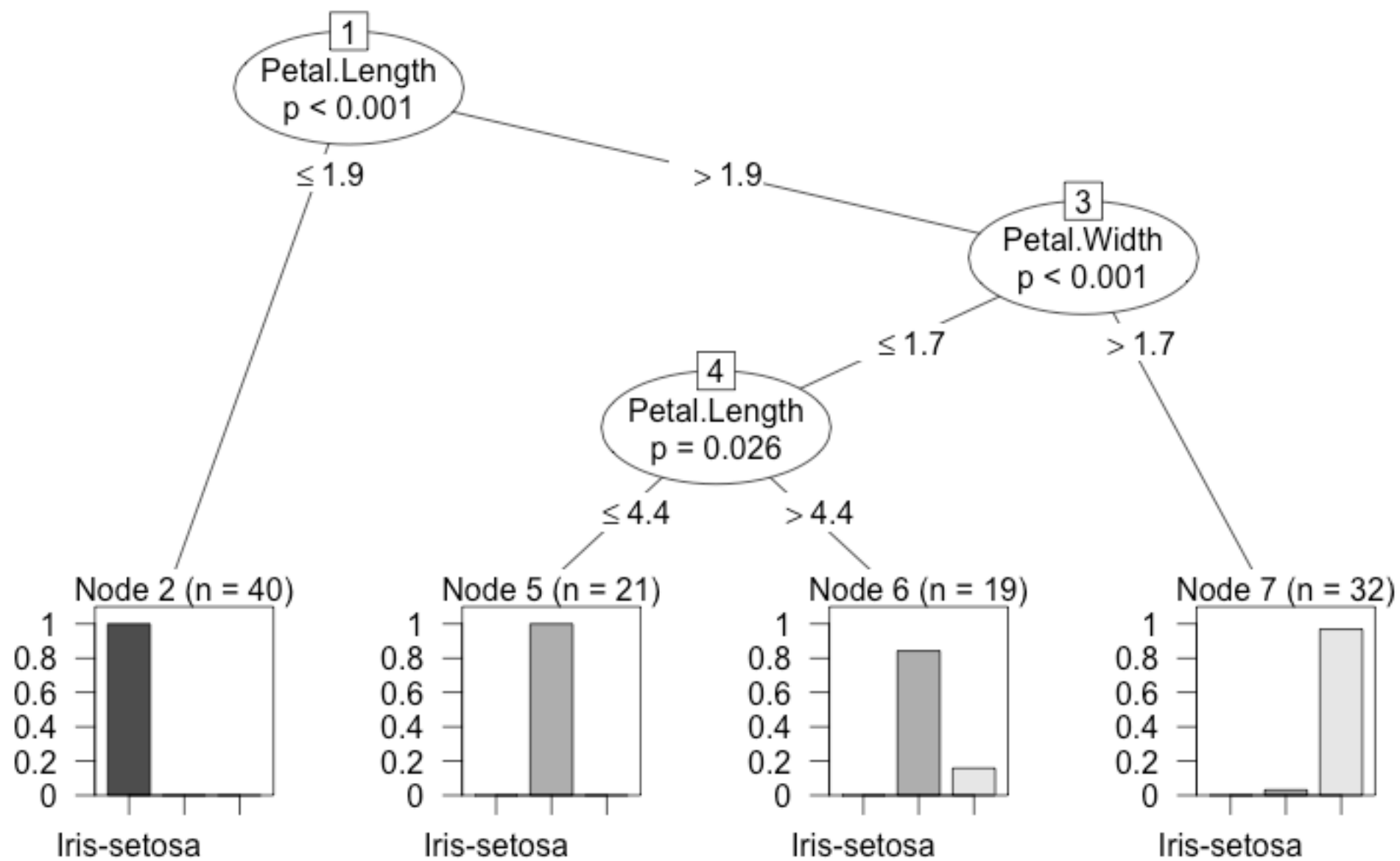
```
4) Petal.Length > 4.4
```

```
6)* weights = 19
```

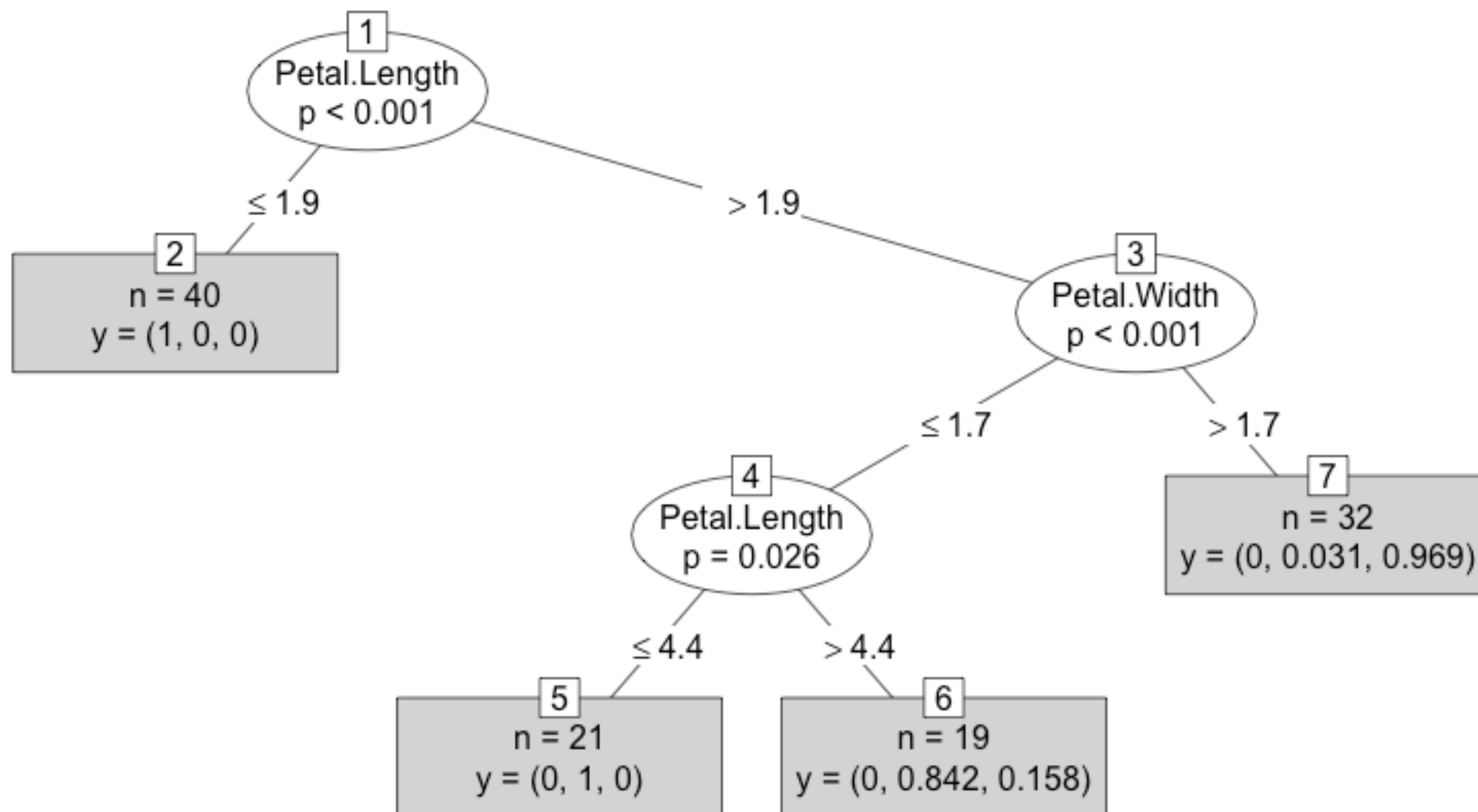
```
3) Petal.Width > 1.7
```

```
7)* weights = 32
```

```
> plot(iris_ctree)
```



```
> plot(iris_ctree, type="simple")
```



# Decision Tree - Prediction

```
> # predict on test data  
> testPred <- predict(iris_ctree, newdata = testData)  
> table(testPred, testData$Species)
```

testPred	setosa	versicolor	virginica
setosa	10	0	0
versicolor	0	12	2
virginica	0	0	14



# K-NN

## k-Nearest Neighbour Classification

### Description

k-nearest neighbour classification for test set from training set. For each row of the test set, the  $k$  nearest (in Euclidean distance) training set vectors are found, and the classification is decided by majority vote, with ties broken at random. If there are ties for the  $k$ th nearest vector, all candidates are included in the vote.

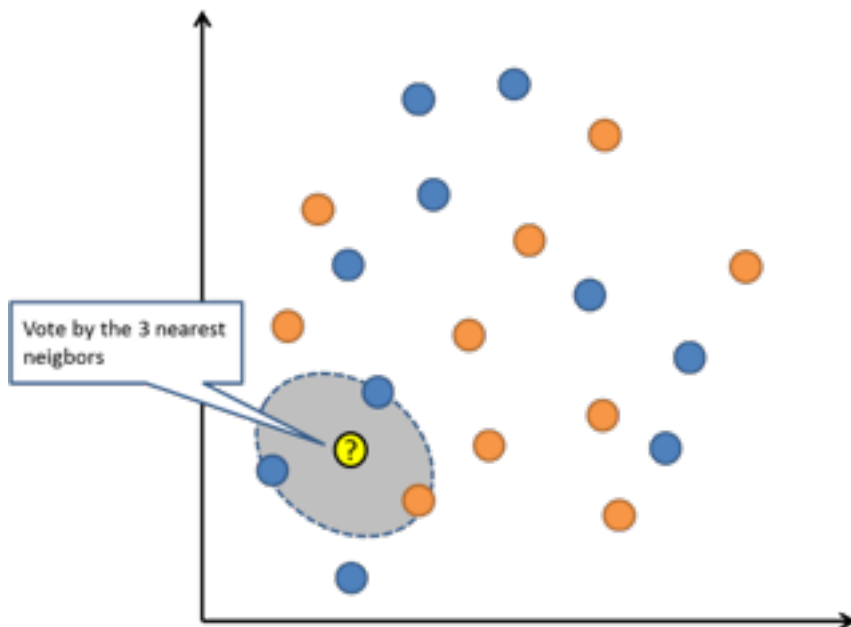
### Usage

```
knn(train, test, cl, k = 1, l = 0, prob = FALSE, use.all = TRUE)
```

### Arguments

<code>train</code>	matrix or data frame of training set cases.
<code>test</code>	matrix or data frame of test set cases. A vector will be interpreted as a row vector for a single case.
<code>cl</code>	factor of true classifications of training set
<code>k</code>	number of neighbours considered.
<code>l</code>	minimum vote for definite decision, otherwise <code>doubt</code> . (More precisely, less than $k-1$ dissenting votes are allowed, even if $k$ is increased by ties.)
<code>prob</code>	If this is true, the proportion of the votes for the winning class are returned as attribute <code>prob</code> .
<code>use.all</code>	controls handling of ties. If true, all distances equal to the $k$ th largest are included. If false, a random selection of distances equal to the $k$ th is chosen to use exactly $k$ neighbours.

# K-NN



```
> library(class)
> train_input <- as.matrix(iris.train[, -5])
> train_output <- as.vector(iris.train[, 5])
> test_input <- as.matrix(iris.test[, -5])
> prediction <- knn(train_input, test_input,
                    train_output, k=5)
> table(prediction, iris.test$Species)
```

prediction	setosa	versicolor	virginica
setosa	10	0	0
versicolor	0	10	1
virginica	0	0	9

```
>
```

# Naive Bayes

## Naive Bayes Classifier

### Description

Computes the conditional a-posterior probabilities of a categorical class variable given independent predictor variables using the Bayes rule.

### Usage

```
## S3 method for class 'formula'
naiveBayes(formula, data, laplace = 0, ..., subset, na.action = na.pass)
## Default S3 method:
naiveBayes(x, y, laplace = 0, ...)

## S3 method for class 'naiveBayes'
predict(object, newdata,
  type = c("class", "raw"), threshold = 0.001, eps = 0, ...)
```

# Naive Bayes

## Arguments

<code>x</code>	A numeric matrix, or a data frame of categorical and/or numeric variables.
<code>y</code>	Class vector.
<code>formula</code>	A formula of the form <code>class ~ x1 + x2 + ....</code> . Interactions are not allowed.
<code>data</code>	Either a data frame of predictors (categorical and/or numeric) or a contingency table.
<code>laplace</code>	positive double controlling Laplace smoothing. The default (0) disables Laplace smoothing.
<code>...</code>	Currently not used.
<code>subset</code>	For data given in a data frame, an index vector specifying the cases to be used in the training sample. (NOTE: If given, this argument must be named.)
<code>na.action</code>	A function to specify the action to be taken if NAs are found. The default action is not to count them for the computation of the probability factors. An alternative is <code>na.omit</code> , which leads to rejection of cases with missing values on any required variable. (NOTE: If given, this argument must be named.)
<code>object</code>	An object of class <code>"naiveBayes"</code> .
<code>newdata</code>	A dataframe with new predictors (with possibly fewer columns than the training data). Note that the column names of <code>newdata</code> are matched against the training data ones.
<code>type</code>	If <code>"raw"</code> , the conditional a-posterior probabilities for each class are returned, and the class with maximal probability else.
<code>threshold</code>	Value replacing cells with probabilities within <code>eps</code> range.
<code>eps</code>	double for specifying an epsilon-range to apply laplace smoothing (to replace zero or close-zero probabilities by <code>threshold</code> .)

# Naive Bayes

```
> library(e1071)
> # Can handle both categorical and numeric input,
> # but output must be categorical
> model <- naiveBayes(Species~., data=iristrain)
> prediction <- predict(model, iristest[,5])
> table(prediction, iristest[,5])
```

prediction	setosa	versicolor	virginica
setosa	10	0	0
versicolor	0	10	2
virginica	0	0	8



# Neural Network

## Training of neural networks

### Description

`neuralnet` is used to train neural networks using backpropagation, resilient backpropagation (RPROP) with (Riedmiller, 1994) or without weight backtracking (Riedmiller and Braun, 1993) or the modified globally convergent version (GRPROP) by Anastasiadis et al. (2005). The function allows flexible settings through custom-choice of error and activation function. Furthermore the calculation of generalized weights (Intrator O. and Intrator N., 1993) is implemented.

### Usage

```
neuralnet(formula, data, hidden = 1, threshold = 0.01,  
          stepmax = 1e+05, rep = 1, startweights = NULL,  
          learningrate.limit = NULL,  
          learningrate.factor = list(minus = 0.5, plus = 1.2),  
          learningrate=NULL, lifesign = "none",  
          lifesign.step = 1000, algorithm = "rprop+",  
          err.fct = "sse", act.fct = "logistic",  
          linear.output = TRUE, exclude = NULL,  
          constant.weights = NULL, likelihood = FALSE)
```

# Neural Network

## Arguments

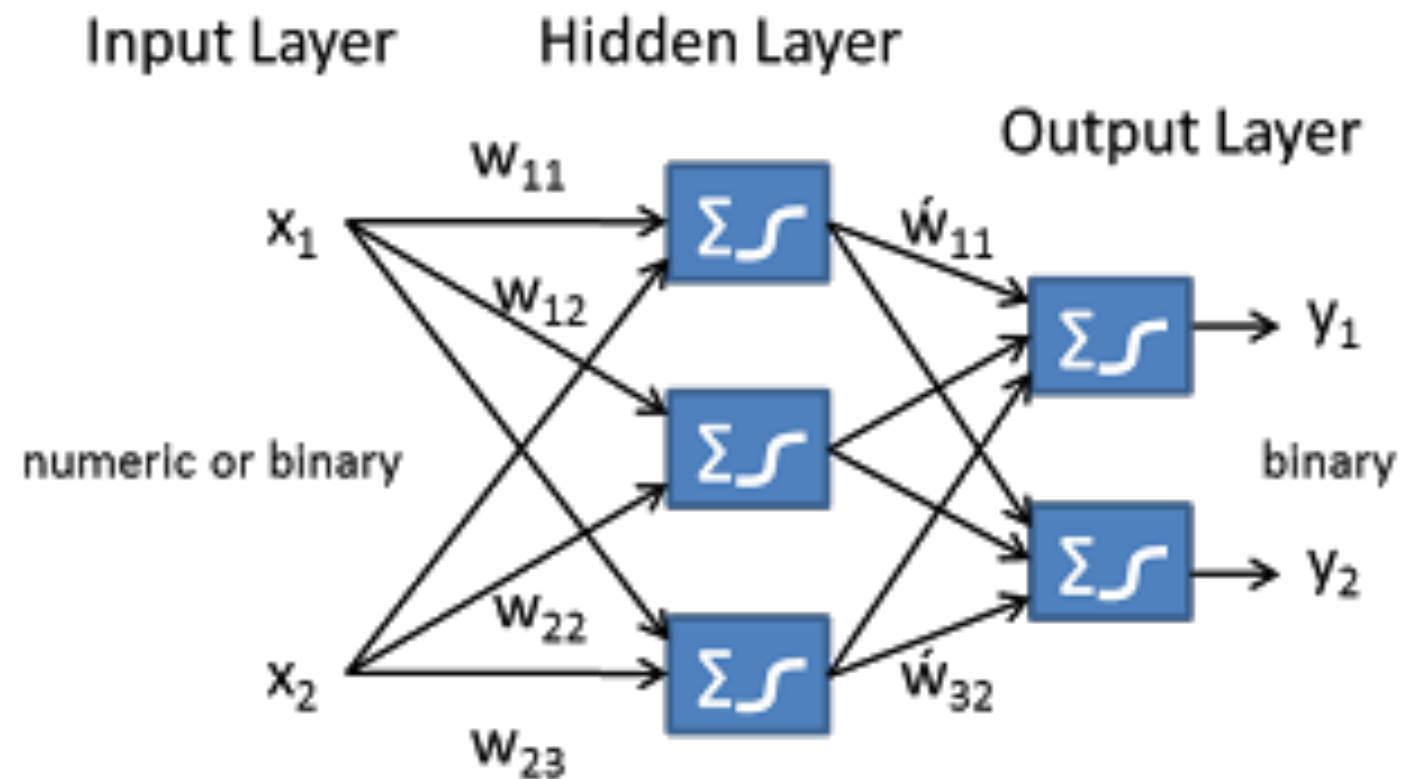
<code>formula</code>	a symbolic description of the model to be fitted.
<code>data</code>	a data frame containing the variables specified in <code>formula</code> .
<code>hidden</code>	a vector of integers specifying the number of hidden neurons (vertices) in each layer.
<code>threshold</code>	a numeric value specifying the threshold for the partial derivatives of the error function as stopping criteria.
<code>stepmax</code>	the maximum steps for the training of the neural network. Reaching this maximum leads to a stop of the neural network's training process.
<code>rep</code>	the number of repetitions for the neural network's training.
<code>startweights</code>	a vector containing starting values for the weights. The weights will not be randomly initialized.
<code>learningrate.limit</code>	a vector or a list containing the lowest and highest limit for the learning rate. Used only for RPROP and GRPROP.
<code>learningrate.factor</code>	a vector or a list containing the multiplication factors for the upper and lower learning rate. Used only for RPROP and GRPROP.
<code>learningrate</code>	a numeric value specifying the learning rate used by traditional backpropagation. Used only for traditional backpropagation.
<code>lifesign</code>	a string specifying how much the function will print during the calculation of the neural network. 'none', 'minimal' or 'full'.

# Neural Network

<code>lifesign.step</code>	an integer specifying the stepsize to print the minimal threshold in full lifesign mode.
<code>algorithm</code>	a string containing the algorithm type to calculate the neural network. The following types are possible: 'backprop', 'rprop+', 'rprop-', 'sag', or 'slr'. 'backprop' refers to backpropagation, 'rprop+' and 'rprop-' refer to the resilient backpropagation with and without weight backtracking, while 'sag' and 'slr' induce the usage of the modified globally convergent algorithm (grprop). See Details for more information.
<code>err.fct</code>	a differentiable function that is used for the calculation of the error. Alternatively, the strings 'sse' and 'ce' which stand for the sum of squared errors and the cross-entropy can be used.
<code>act.fct</code>	a differentiable function that is used for smoothing the result of the cross product of the covariate or neurons and the weights. Additionally the strings, 'logistic' and 'tanh' are possible for the logistic function and tangent hyperbolicus.
<code>linear.output</code>	logical. If act.fct should not be applied to the output neurons set linear output to TRUE, otherwise to FALSE.
<code>exclude</code>	a vector or a matrix specifying the weights, that are excluded from the calculation. If given as a vector, the exact positions of the weights must be known. A matrix with n-rows and 3 columns will exclude n weights, where the first column stands for the layer, the second column for the input neuron and the third column for the output neuron of the weight.
<code>constant.weights</code>	a vector specifying the values of the weights that are excluded from the training process and treated as fix.
<code>likelihood</code>	logical. If the error function is equal to the negative log-likelihood function, the information criteria AIC and BIC will be calculated. Furthermore the usage of confidence.interval is meaningfull.



# Neural Network

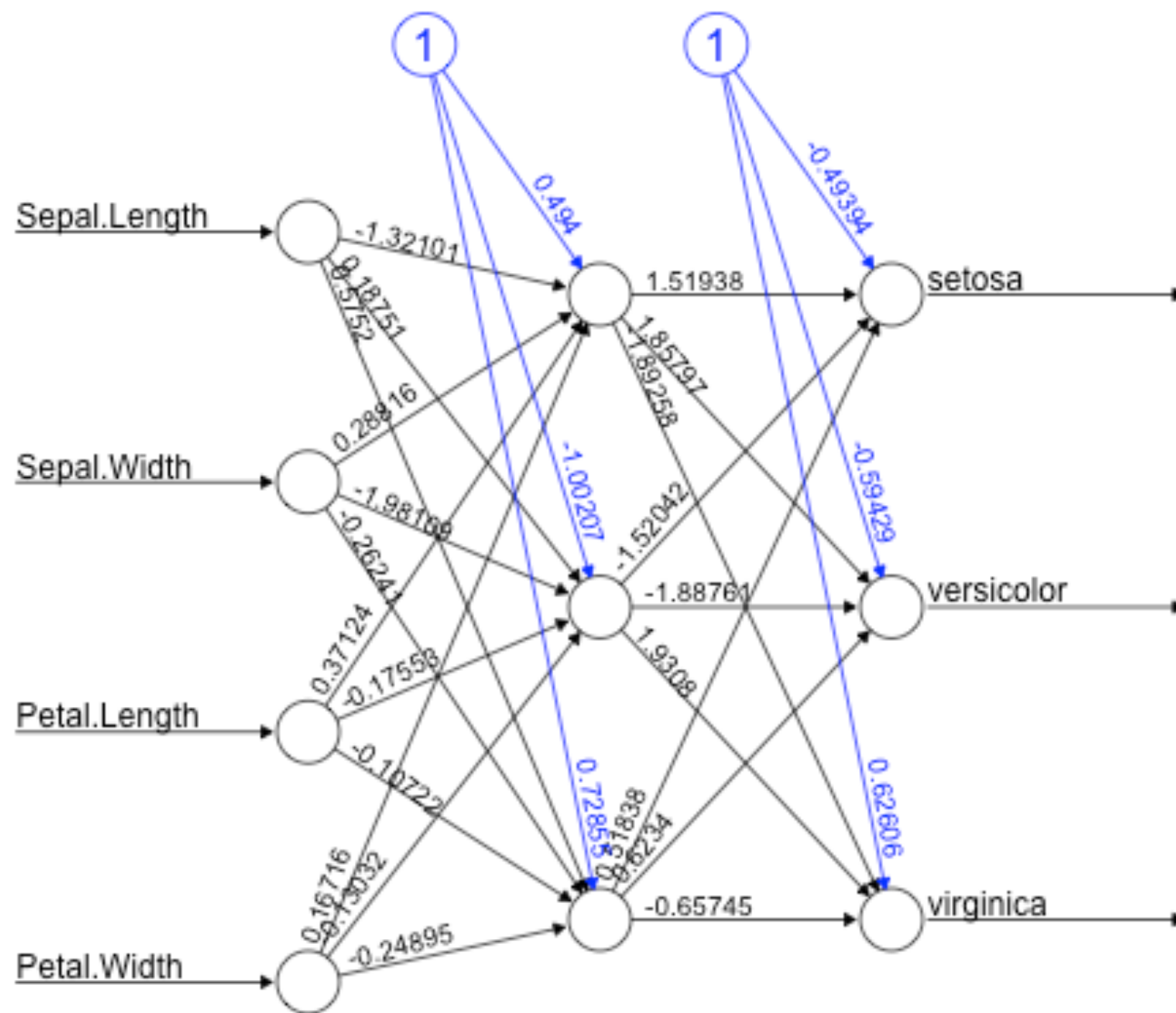


```

> library(neuralnet)
> nnet_irisrain <- irisrain
> #Binarize the categorical output
> nnet_irisrain <- cbind(nnet_irisrain,
                        irisrain$Species == 'setosa')
> nnet_irisrain <- cbind(nnet_irisrain,
                        irisrain$Species == 'versicolor')
> nnet_irisrain <- cbind(nnet_irisrain,
                        irisrain$Species == 'virginica')
> names(nnet_irisrain)[6] <- 'setosa'
> names(nnet_irisrain)[7] <- 'versicolor'
> names(nnet_irisrain)[8] <- 'virginica'
> nn <- neuralnet(setosa+versicolor+virginica ~
                  Sepal.Length+Sepal.Width
                  +Petal.Length
                  +Petal.Width,
                  data=nnet_irisrain,
                  hidden=c(3))
> plot(nn)
> mypredict <- compute(nn, irisrain[-5])$net.result
> # Put multiple binary output to categorical output
> maxidx <- function(arr) {
  return(which(arr == max(arr)))
}
> idx <- apply(mypredict, c(1), maxidx)
> prediction <- c('setosa', 'versicolor', 'virginica')[idx]
> table(prediction, irisrain$Species)

```

prediction	setosa	versicolor	virginica
setosa	10	0	0
versicolor	0	10	3
virginica	0	0	7



Error: 0.001277 Steps: 281