



Data Mining with R

part 2

วิชา การค้นพบองค์ความรู้และการทำเหมืองข้อมูลขั้นสูง

Veerasak Kritsanapraphan

Chulalongkorn University

Email : veerasak.kr568@cbs.chula.ac.th

 @veerasakk

Agenda

- Data Mining Techniques using R
 - Neural Network
 - Clustering
 - K-Means Clustering
 - Hierarchical Clustering
 - Association Rules
 - Multi-mobile Learning
- Parallel Computing using R
 - Hadoop

Slide and Source Codes

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



Prepare Data

- `iris <- read.csv("iris.data.csv", header=TRUE)`
- `# Prepare iris`
- `set.seed(567)`
- `ind <- sample(2, nrow(iris), replace=TRUE, prob=c(0.7, 0.3))`
- `traindata <- iris[ind==1,]`
- `testdata <- iris[ind==2,]`

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

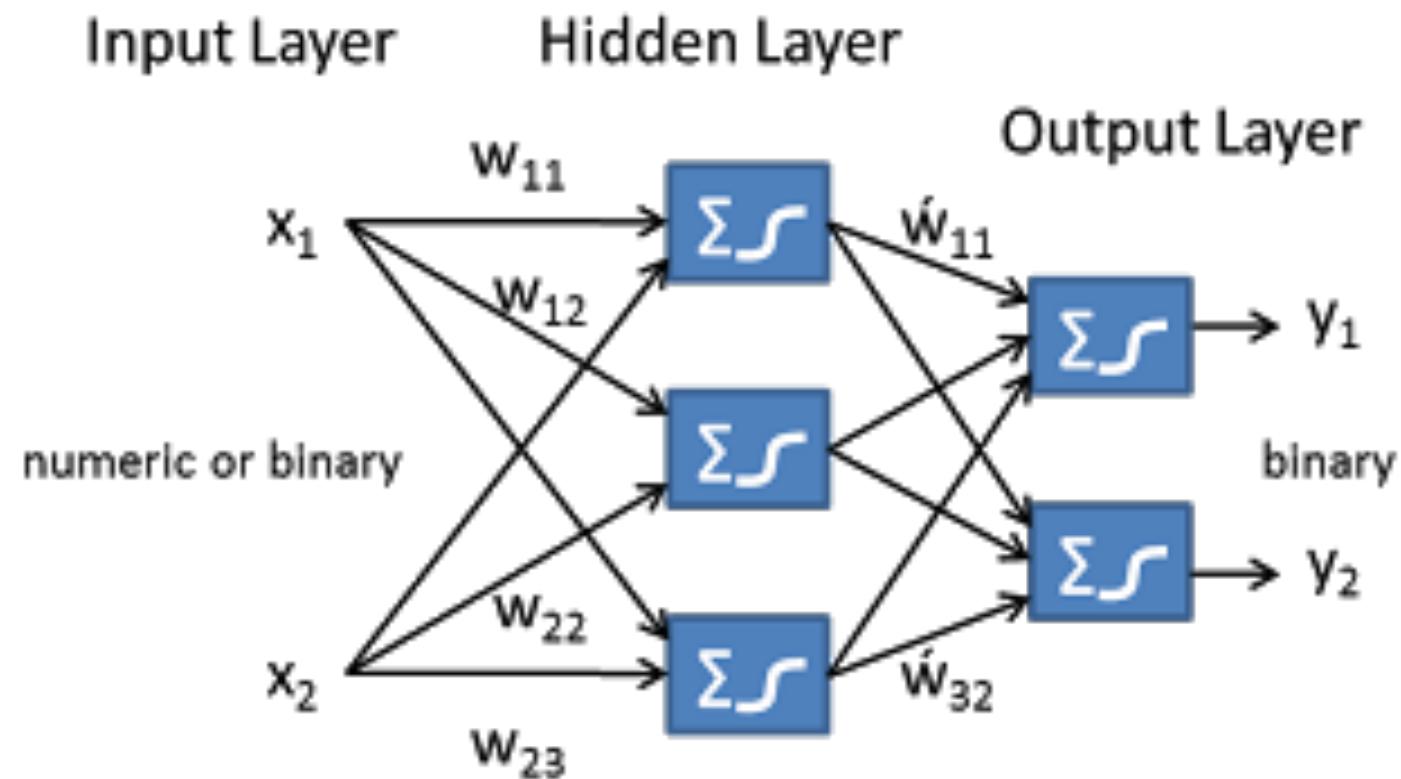
Arguments

<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



Code

```
> Library neuralnet
```

```
> nnet_istrain <- traindata
```

```
# Binarize Categorical Output
```

```
> nnet_istrain <- cbind(nnet_istrain, traindata$Species == "Iris-setosa")
```

```
> nnet_istrain <- cbind(nnet_istrain, traindata$Species == "Iris-versicolor")
```

```
> nnet_istrain <- cbind(nnet_istrain, traindata$Species == "Iris-verginiaca")
```

```
> names(nnet_istrain) [6] <- "setosa"
```

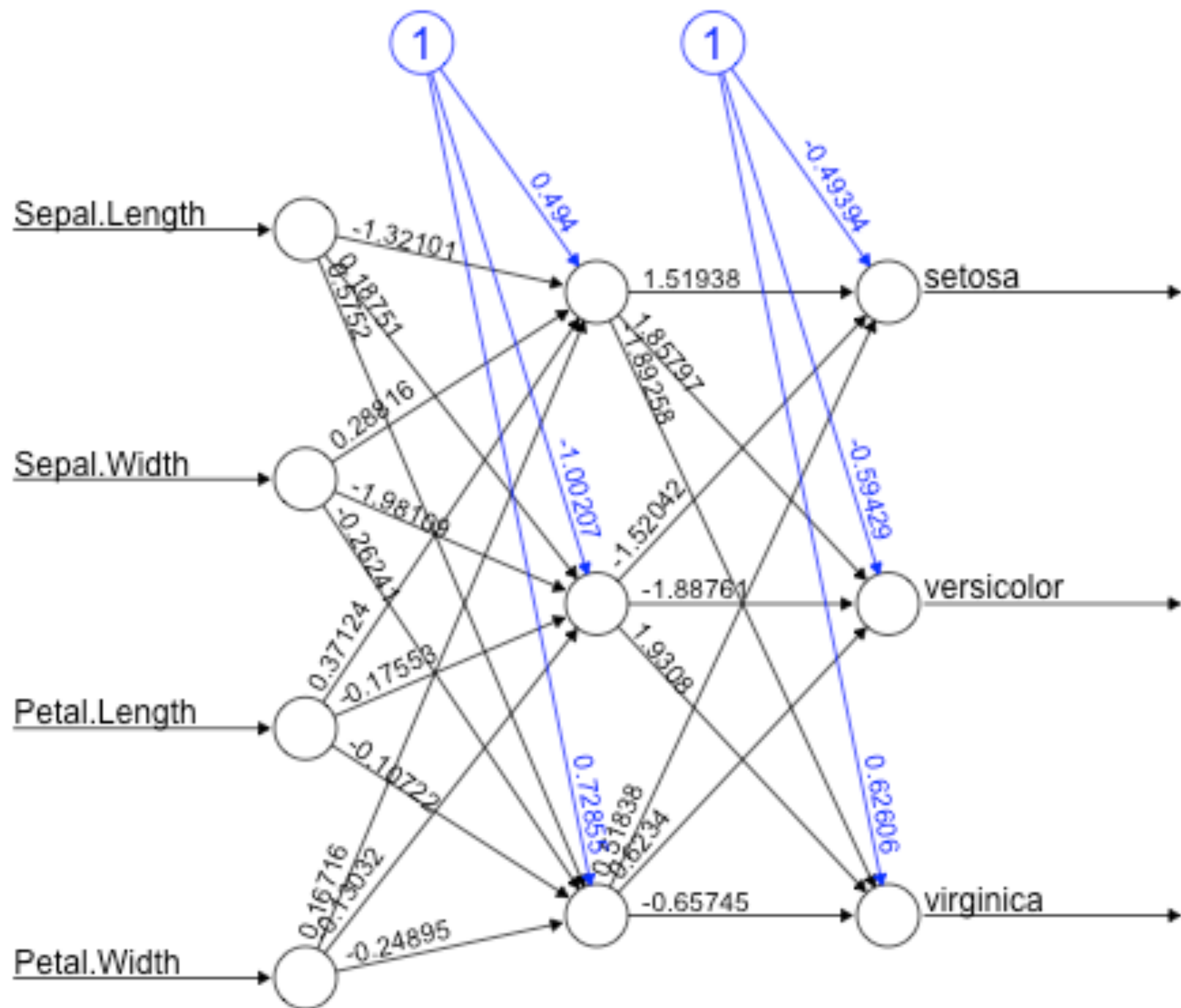
```
> names(nnet_istrain) [6] <- "versicolor"
```

```
> names(nnet_istrain) [6] <- "verginiaca"
```

```
# Run Neural Net
```

```
> nn <- neuralnet(setosa+versicolor+verginiaca~Sepal.Length+Sepal.Width+  
                  Petal.Length+Petal.Width
```

```
> plot(nn)
```



Error: 0.001277 Steps: 281

Code (continue)

```
> mypredict <- compute(nn, testdata[-5])$net.result  
  
> # Convert binary output into categorical output  
  
> maxidx <- function(arr) {  
    return (which(arr == max(arr)))  
}  
  
> idx <- apply(mypredict, c(1), maxidx)  
  
> prediction <- c("setosa", "versicolor", "virginica")[idx]  
  
> table(prediction, testdata$Species)
```

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

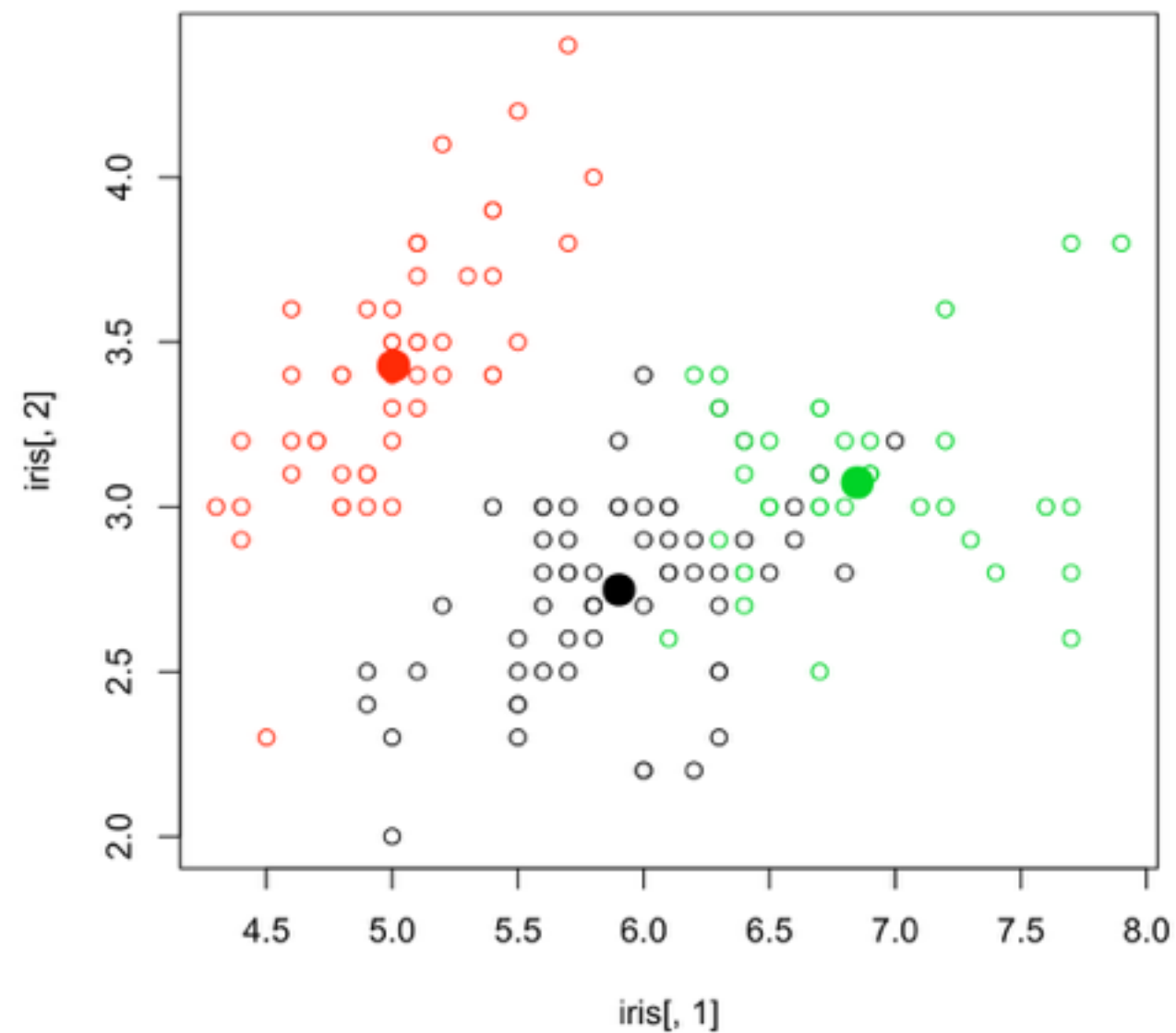
Clustering

- K-Means Clustering
- Hierarchical Clustering

K-Means Clustering

1. Pick an initial set of K centroids (this can be random or any other means)
2. For each data point, assign it to the member of the closest centroid according to the given distance function
3. Adjust the centroid position as the mean of all its assigned member data points. Go back to (2) until the membership isn't change and centroid position is stable.
4. Output the centroids.

- `library(stats)`
- `set.seed(101)`
- `km <- kmeans(iris[,1:4], 3)`
- `plot(iris[,1], iris[,2], col=km$cluster)`
- `points(km$centers[,c(1,2)], col=1:3, pch=19, cex=2)`

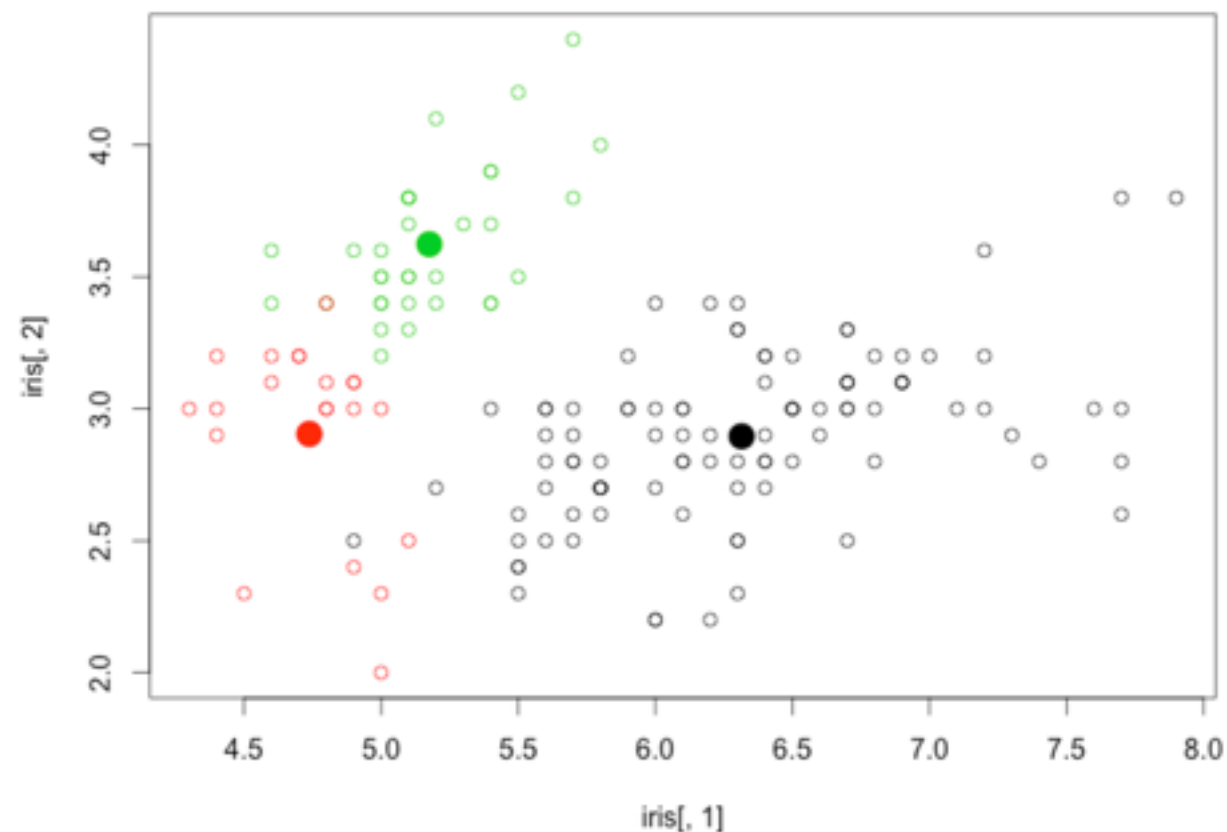


- `table(km$cluster, iris$Species)`

```
      setosa versicolor virginica  
1         0         48         14  
2        50          0          0  
3         0          2         36  
>
```

Another round

- `set.seed(900)`
- `km <- kmeans(iris[,1:4], 3)`
- `plot(iris[,1], iris[,2], col=km$cluster)`
- `points(km$centers[,c(1,2)], col=1:3, pch=19, cex=2)`



- `table(km$cluster, iris$Species)`

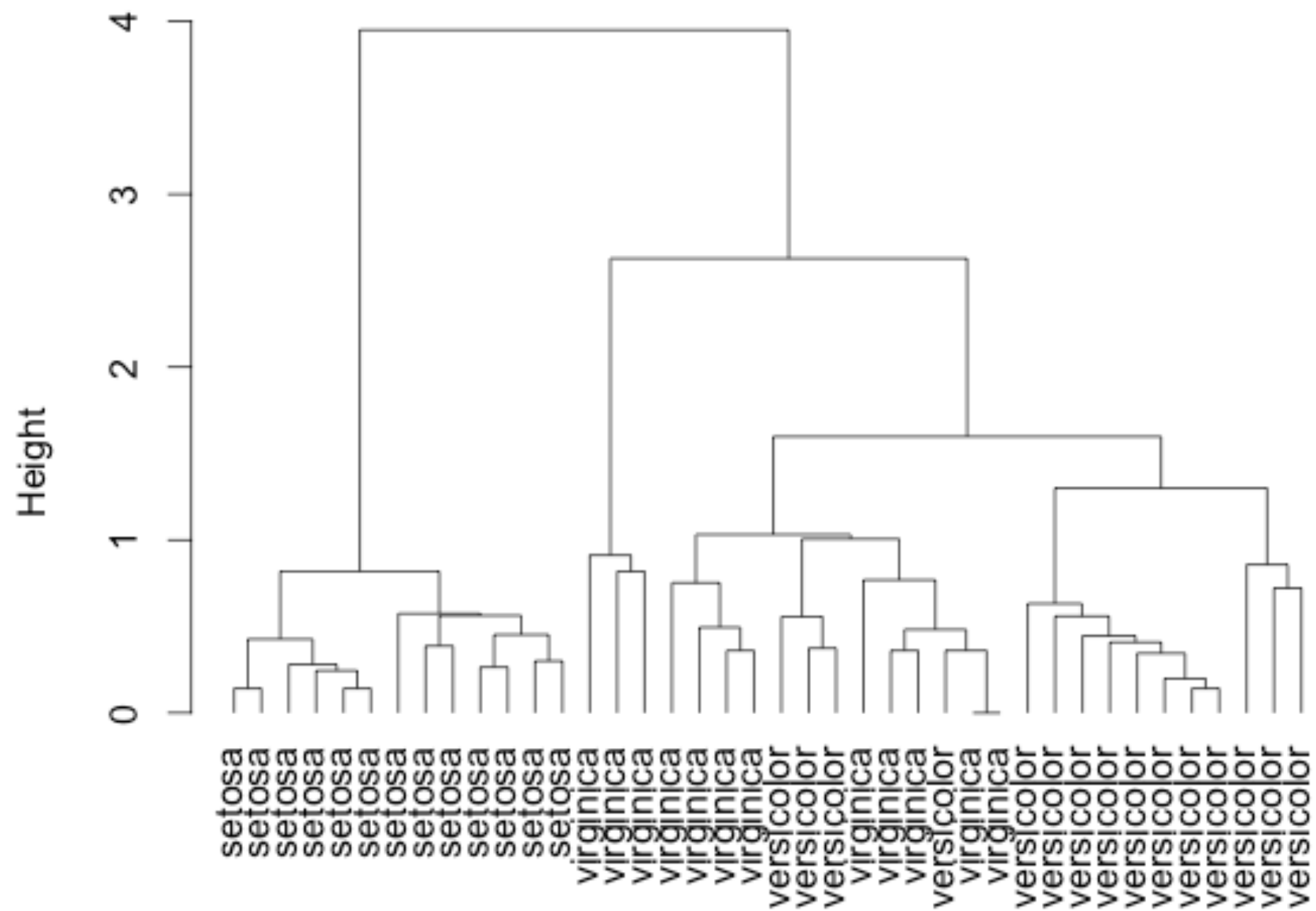
```
      setosa versicolor virginica  
1         0         46         50  
2        17          4          0  
3        33          0          0  
>
```

Hierarchical Clustering

- Compute distance between every pairs of point/cluster.
 - (a) Distance between point is just using the distance function.
 - (b) Compute distance between pointA to clusterB may involve many choices (such as the min/max/avg distance between the pointA and points in the clusterB).
 - (c) Compute distance between clusterA to clusterB may first compute distance of all points pairs (one from clusterA and the other from clusterB) and then pick either min/max/avg of these pairs.
- Combine the two closest point/cluster into a cluster. Go back to (1) until only one big cluster remains

- `set.seed(101)`
- `sampleiris <- iris[sample(1:150, 40),] # get samples from iris dataset`
- `# each observation has 4 variables, ie, they are interpreted as 4-D points`
- `distance <- dist(sampleiris[,-5], method="euclidean")`
- `cluster <- hclust(distance, method="average")`
- `plot(cluster, hang=-1, label=sampleiris$Species)`

Cluster Dendrogram



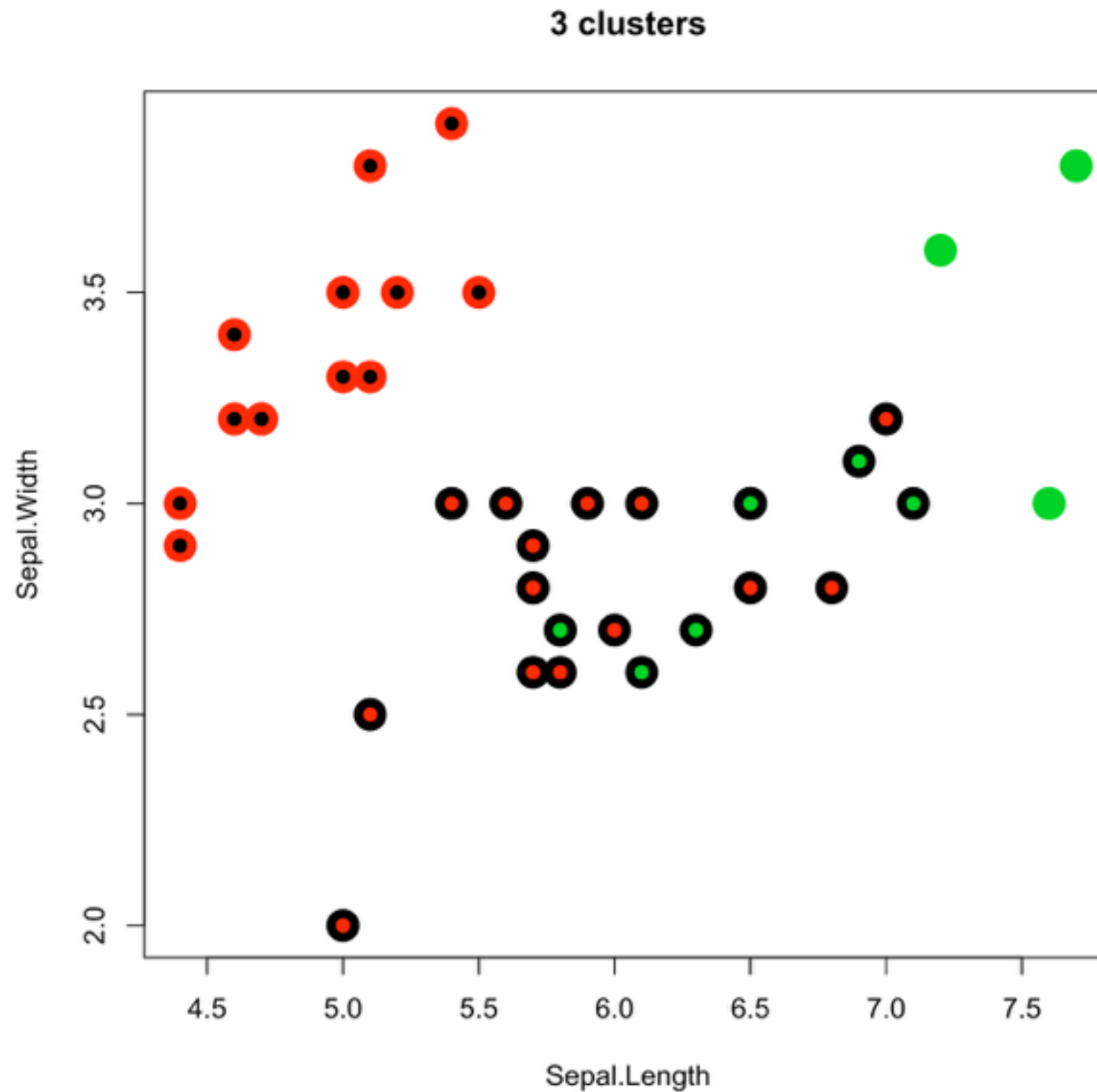
distance
hclust (*, "average")

It's possible to prune the result tree.

- `par(mfrow=c(1,2))`
- `group.3 <- cutree(cluster, k = 3) # prune the tree by 3 clusters`
- `table(group.3, sampleiris$Species) # compare with known classes`

```
group.3 setosa versicolor virginica
      1      0          15          9
      2     13          0          0
      3      0          0          3
> |
```

- `plot(sampleiris[,c(1,2)], col=group.3, pch=19, cex=2.5, main="3 clusters")`
- `points(sampleiris[,c(1,2)], col=sampleiris$Species, pch=19, cex=1)`



Association Rules (Market Basket Analysis)

- **Support:** The rule holds with support sup in T (the transaction data set) if $\text{sup}\%$ of transactions contain $X \cup Y$.
- $\text{sup} = \Pr(X \cup Y)$.
- **Confidence:** The rule holds in T with confidence conf if $\text{conf}\%$ of transactions that contain X also contain Y .
- $\text{conf} = \Pr(Y | X)$
- **Lift:** The Lift of the rule is $X \Rightarrow Y$ is the confidence of the rule divided by the expected confidence, assuming that the item sets are independent.



Apriori Algorithm

- ?apriori

Usage

```
apriori(data, parameter = NULL, appearance = NULL, control = NULL)
```

Arguments

data

object of class [transactions](#) or any data structure which can be coerced into [transactions](#) (e.g., a binary matrix or data.frame).

parameter

object of class [APparameter](#) or named list. The default behavior is to mine rules with support 0.1, confidence 0.8, and maxlen 10.

appearance

object of class [APappearance](#) or named list. With this argument item appearance can be restricted. By default all items can appear unrestricted.

control

object of class [APcontrol](#) or named list. Controls the performance of the mining algorithm (item sorting, etc.)

Apriori Algorithm

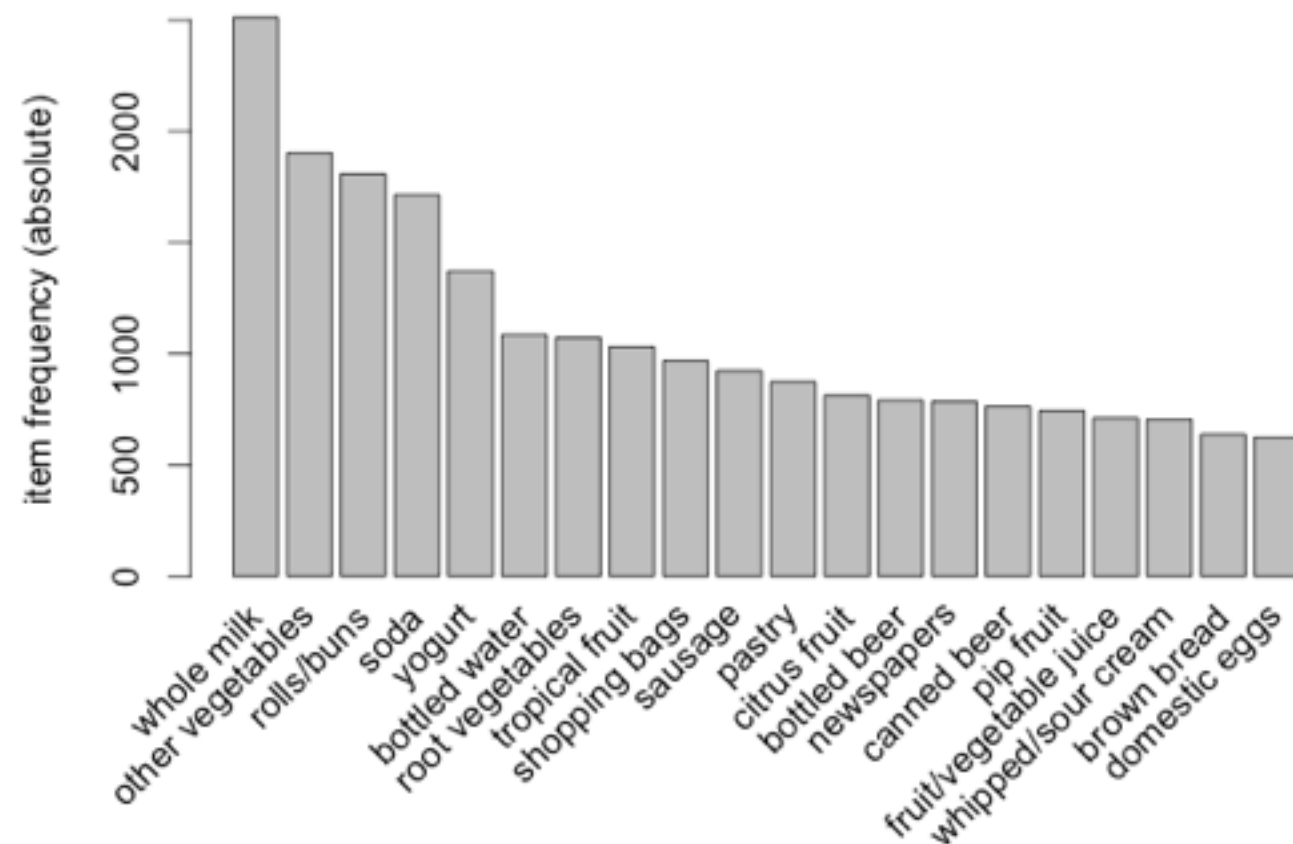
So lets get started by loading up our libraries and data set.

- `# Load the libraries`
- `library(arules)`
- `library(arulesViz)`
- `library(datasets)`
-
- `# Load the data set`
- `data(Groceries)`

Transaction ID	Milk	Bread	Butter	Beer
1	1	1	0	0
2	0	1	1	0
3	0	0	0	1
4	1	1	1	0
5	0	1	0	0

Explore Data

- # Create an item frequency plot for the top 20 items
- `itemFrequencyPlot(Groceries, topN=20,type="absolute")`



- **rules <- apriori(Groceries, parameter = list(supp = 0.001, conf = 0.8))**
- # Show the top 5 rules, but only 2 digits
- options(digits=2)
- inspect(rules[1:5])

```

  lhs                rhs      support confidence lift
1 {liquor,          => {bottled beer} 0.0019      0.90 11.2
   red/blush wine}
2 {curd,            => {whole milk}  0.0010      0.91  3.6
   cereals}
3 {yogurt,          => {whole milk}  0.0017      0.81  3.2
   cereals}
4 {butter,          => {whole milk}  0.0010      0.83  3.3
   jam}
5 {soups,           => {whole milk}  0.0011      0.92  3.6
   bottled beer}
> |

```

- summary(rules)

```

set of 410 rules

rule length distribution (lhs + rhs):sizes
  3   4   5   6
29 229 140  12

      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
      3.0   4.0   4.0   4.3   5.0   6.0

summary of quality measures:
      support      confidence      lift
Min.      :0.00102  Min.      :0.80  Min.      : 3.1
1st Qu.:0.00102  1st Qu.:0.83  1st Qu.: 3.3
Median :0.00122  Median :0.85  Median : 3.6
Mean    :0.00125  Mean    :0.87  Mean    : 4.0
3rd Qu.:0.00132  3rd Qu.:0.91  3rd Qu.: 4.3
Max.    :0.00315  Max.    :1.00  Max.    :11.2

mining info:
      data ntransactions support confidence
Groceries      9835      0.001      0.8

```

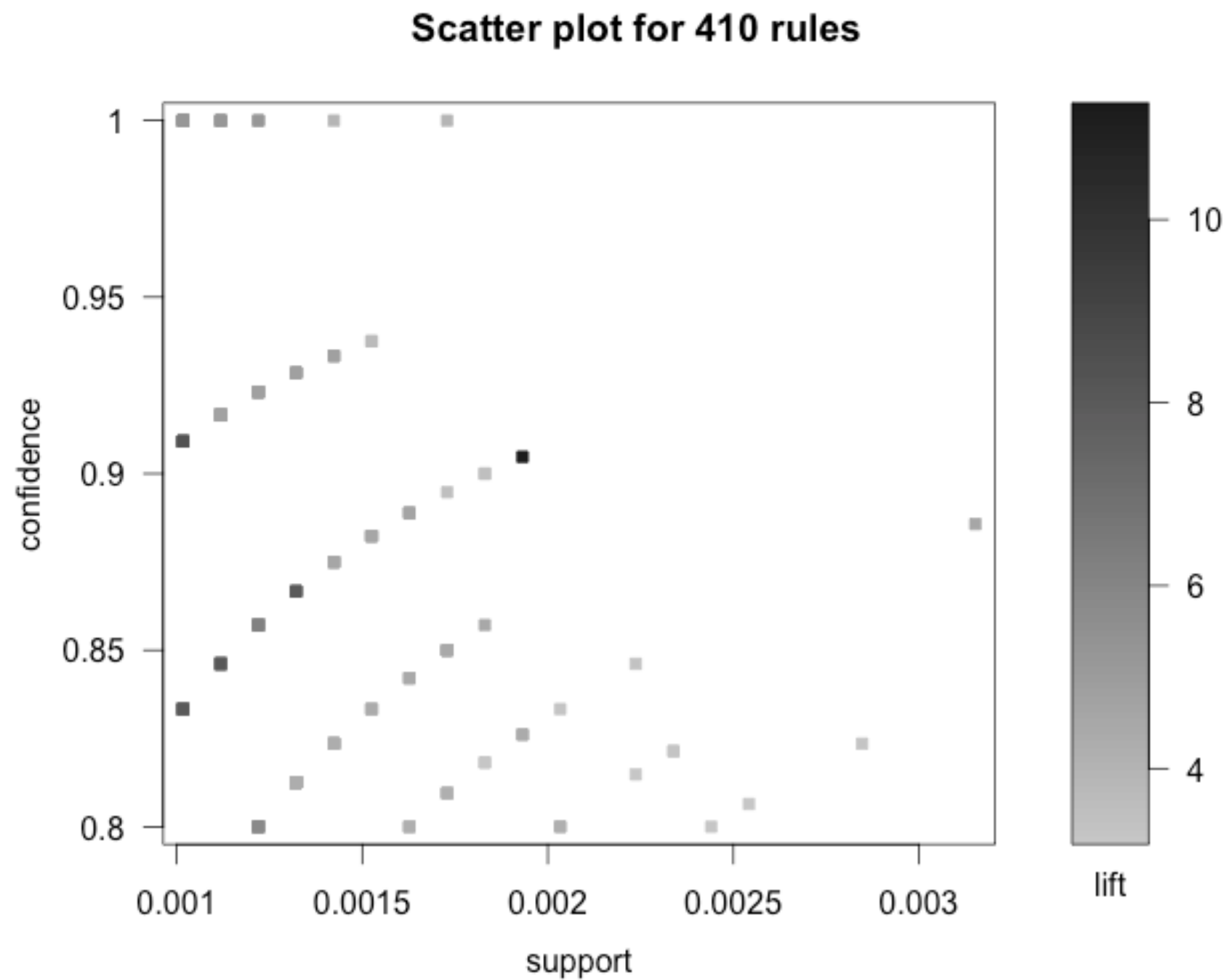
- # Sort Rules
- rules<-sort(rules, by="confidence", decreasing=TRUE)
- inspect(rules[1:5])

```

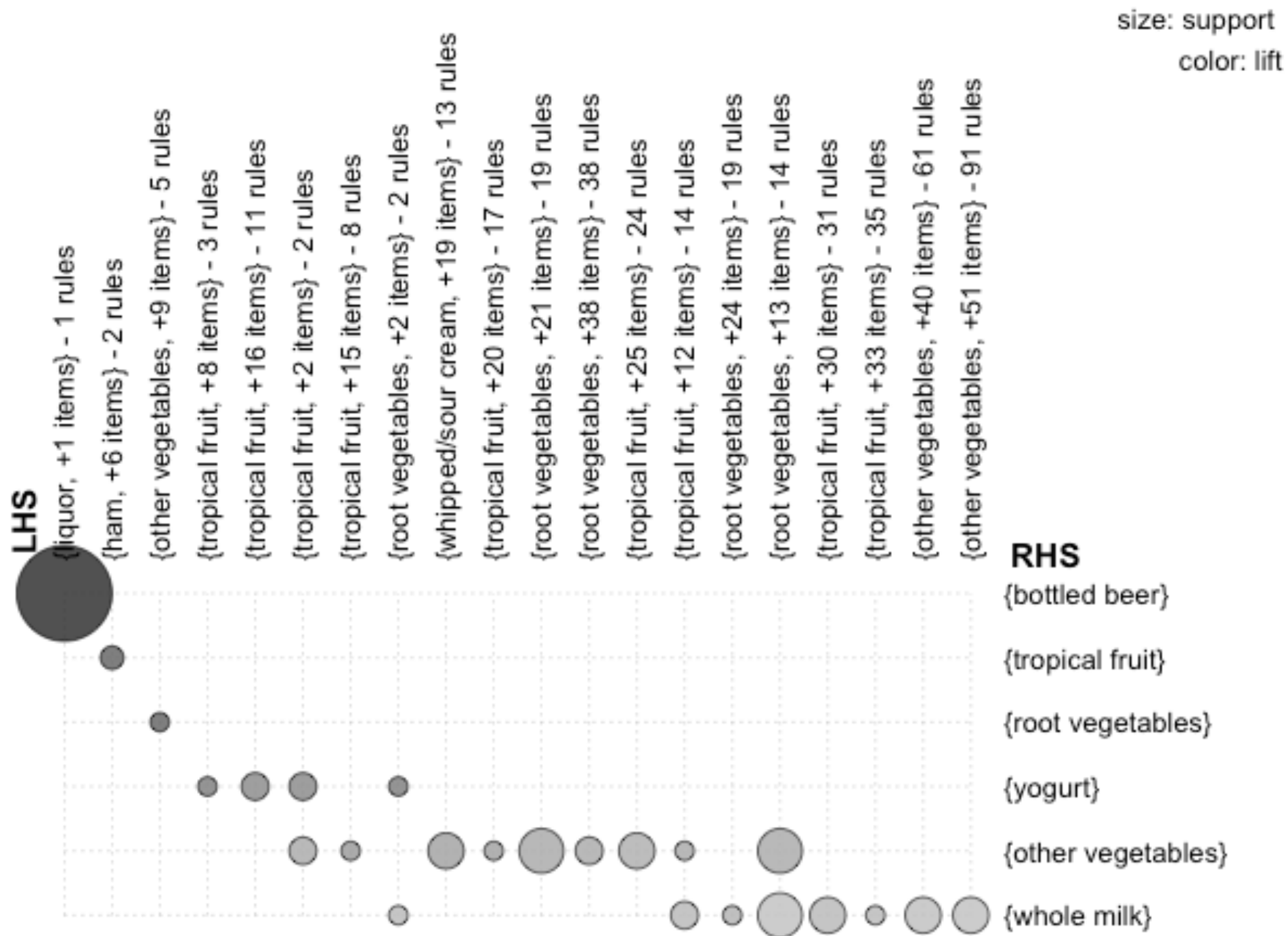
  lhs                rhs                support confidence lift
1 {rice,             => {whole milk}  0.0012             1  3.9
   sugar}
2 {canned fish,      => {whole milk}  0.0011             1  3.9
   hygiene articles}
3 {root vegetables,  => {whole milk}  0.0010             1  3.9
   butter,
   rice}
4 {root vegetables,  => {whole milk}  0.0017             1  3.9
   whipped/sour cream,
   flour}
5 {butter,           => {whole milk}  0.0010             1  3.9
   soft cheese,
   domestic eggs}
> |

```

- `plot(rules);`
- `plot(rules, method="grouped")`



Grouped matrix for 410 rules

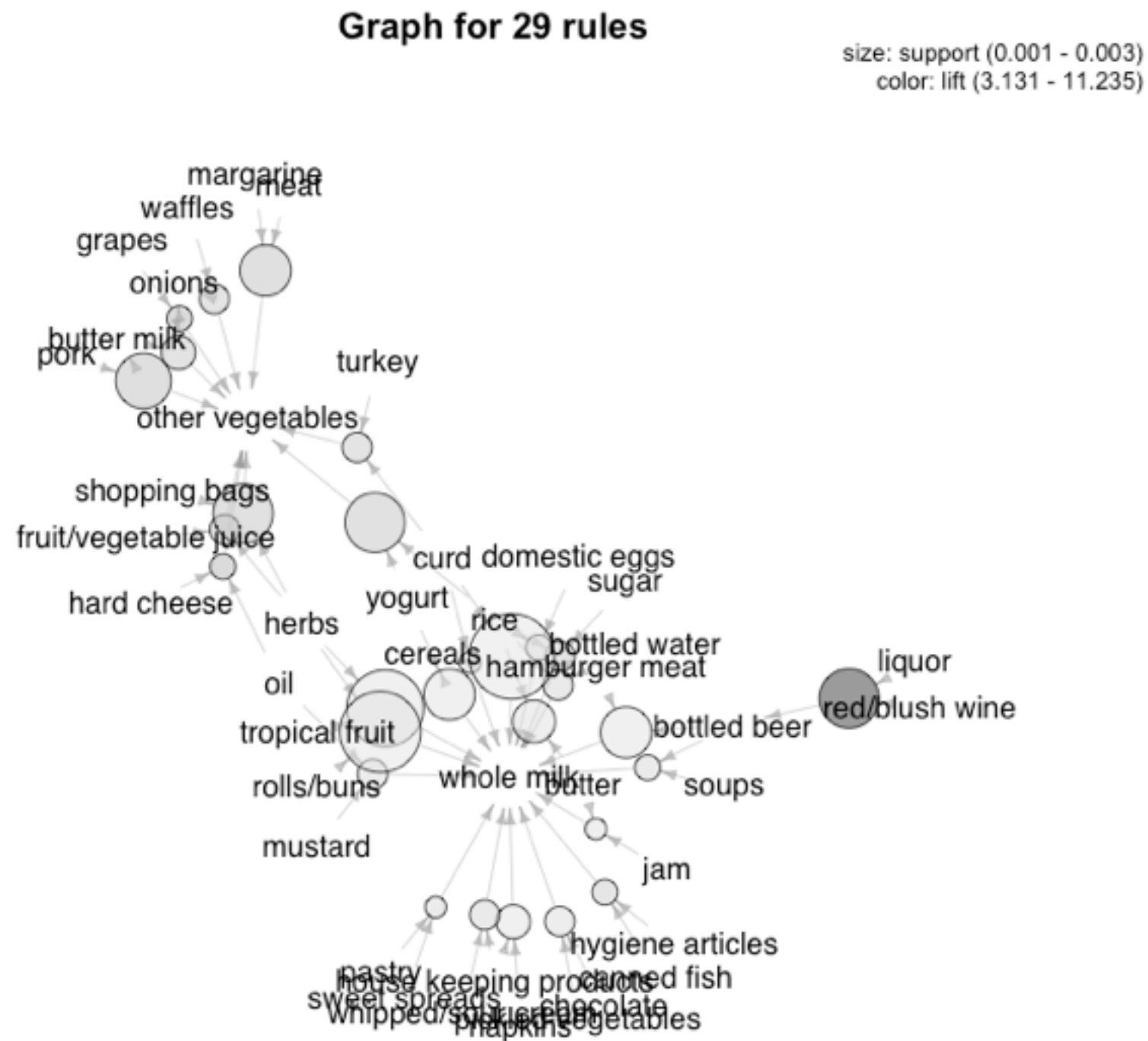


Change to have limit association in one rule

- # change to have maximum of 3
- `rules <- apriori(Groceries, parameter = list(supp = 0.001, conf = 0.8,maxlen=3))`
- `inspect(rules[1:5])`

	lhs	rhs	support	confidence	lift
1	{liquor, red/blush wine}	=> {bottled beer}	0.0019	0.90	11.2
2	{curd, cereals}	=> {whole milk}	0.0010	0.91	3.6
3	{yogurt, cereals}	=> {whole milk}	0.0017	0.81	3.2
4	{butter, jam}	=> {whole milk}	0.0010	0.83	3.3
5	{soups, bottled beer}	=> {whole milk}	0.0011	0.92	3.6

- `plot(rules, method="graph")`



Rules pruned

- `subset.matrix <- is.subset(rules, rules)`
- `subset.matrix[lower.tri(subset.matrix, diag=T)] <- NA`
- `redundant <- colSums(subset.matrix, na.rm=T) >= 1`
- `rules.pruned <- rules[!redundant]`
- `rules <- rules.pruned`
- `summary(rules)`

set of 330 rules

rule length distribution (lhs + rhs): sizes

3	4	5	6
29	216	84	1

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
3.0	4.0	4.0	4.2	5.0	6.0

summary of quality measures:

support	confidence	lift
Min. :0.00102	Min. :0.80	Min. : 3.1
1st Qu.:0.00102	1st Qu.:0.82	1st Qu.: 3.3
Median :0.00122	Median :0.85	Median : 3.6
Mean :0.00127	Mean :0.86	Mean : 3.8
3rd Qu.:0.00132	3rd Qu.:0.91	3rd Qu.: 4.3
Max. :0.00315	Max. :1.00	Max. :11.2

mining info:

data	ntransactions	support	confidence
Groceries	9835	0.001	0.8

Targeting Items

- What are customers likely to buy before buying whole milk?
- What are customers likely to buy if they purchase whole milk?
- This essentially means we want to set either the Left Hand Side and Right Hand Side. This is not difficult to do with R!

Find whole milk's antecedents

- `rules<-apriori(data=Groceries, parameter=list(supp=0.001,conf = 0.08), appearance = list(default="lhs",rhs="whole milk"), control = list(verbose=F))`
- `rules<-sort(rules, decreasing=TRUE,by="confidence")`
- `inspect(rules[1:5])`

```
  lhs                rhs      support confidence lift
1 {rice,              => {whole milk}  0.0012         1  3.9
   sugar}
2 {canned fish,       => {whole milk}  0.0011         1  3.9
   hygiene articles}
3 {root vegetables,   => {whole milk}  0.0010         1  3.9
   butter,
   rice}
4 {root vegetables,   => {whole milk}  0.0017         1  3.9
   whipped/sour cream,
   flour}
5 {butter,            => {whole milk}  0.0010         1  3.9
   soft cheese,
   domestic eggs}
> |
```

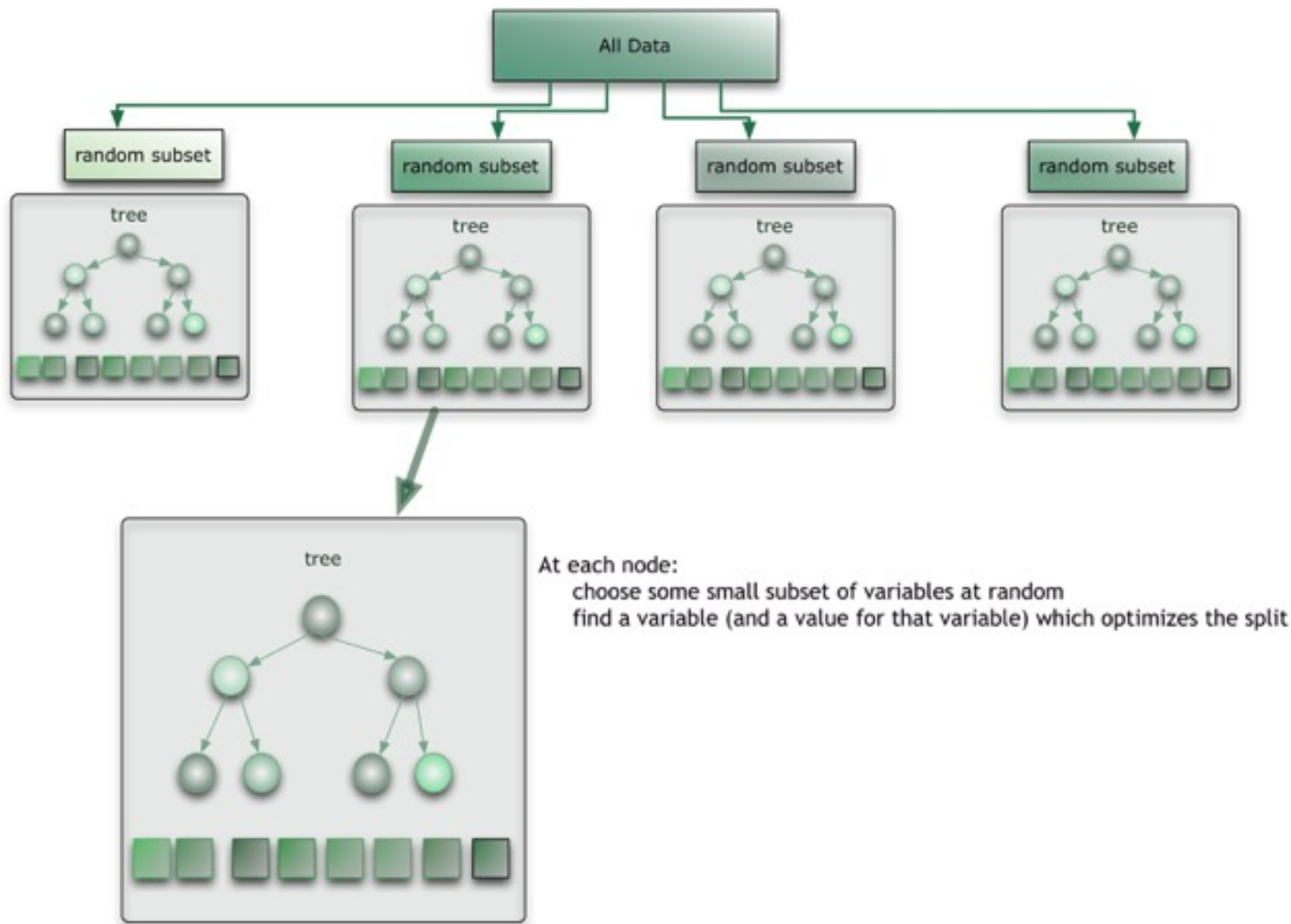
Likely to buy after buy whole milk

- `rules<-apriori(data=Groceries, parameter=list(supp=0.001,conf = 0.15,minlen=2), appearance = list(default="rhs",lhs="whole milk"), control = list(verbose=F))`
- `rules<-sort(rules, decreasing=TRUE,by="confidence")`
- `inspect(rules[1:5])`

```
  lhs                rhs      support confidence lift
1 {whole milk} => {other vegetables} 0.075      0.29  1.5
2 {whole milk} => {rolls/buns}      0.057      0.22  1.2
3 {whole milk} => {yogurt}          0.056      0.22  1.6
4 {whole milk} => {root vegetables} 0.049      0.19  1.8
5 {whole milk} => {tropical fruit}  0.042      0.17  1.6
> |
```

Bagging and Boosting using R

Ensemble : Bagging



Random Forest

- Here is how such a system is trained; for some number of trees T :
- 1) Sample N cases at random with replacement to create a subset of the data. The subset should be about 66% of the total set.
- 2) At each node:
 - a) For some number m (see below), m predictor variables are selected at random from all the predictor variables.
 - b) The predictor variable that provides the best split, according to some objective function, is used to do a binary split on that node.
 - c) At the next node, choose another m variables at random from all predictor variables and do the same.

Bagging

- > library(randomforest)
- > # Train 500 trees, random selected attributes
- > model <- randomForest(Species~., data=traindata, nTree=500)
- > prediction <- predict(model, newdata=testdata, type='class')
- > table(prediction, testdata\$Species)

prediction	setosa	versicolor	virginica
setosa	15	0	0
versicolor	0	13	2
virginica	0	0	13

Boosting

```
> library(adabag)
```

```
> iris.adaboost <- boosting(Species~., data=traindata,  
boost=TRUE, mfinal=5)
```

```
> iris.adaboost
```

Get Text Mining Library

- `Needed <- c("tm", "SnowballCC", "RColorBrewer", "ggplot2", "wordcloud", "biclust", "cluster", "igraph", "fpc")`
- `install.packages(Needed, dependencies=TRUE)`

Load File

- http://www.cs.cornell.edu/people/pabo/movie-review-data/review_polarity.tar.gz
- # Load Text file
- `cname <- file.path("~", "Downloads", "text")`
- `cname`
- `dir(cname)`

Text Mining

- `library(tm)`
- `docs <- Corpus(DirSource(cname))`
- `summary(docs)`

Text Mining Basic

- **# Remove Punctuation**

- docs <- tm_map(docs, removePunctuation)
- for(j in seq(docs))
- {
- docs[[j]] <- gsub("/", " ", docs[[j]])
- docs[[j]] <- gsub("@", " ", docs[[j]])
- docs[[j]] <- gsub("\\\\", " ", docs[[j]])
- }

- **#remove Number**

- docs <- tm_map(docs, removeNumbers)
- docs <- tm_map(docs, tolower)

- **# remove stop words**

- docs <- tm_map(docs, removeWords, stopwords("english"))

Clean Data

- `# remove ing s, es`
- `library(SnowballC)`
- `docs <- tm_map(docs, stemDocument)`
- `docs <- tm_map(docs, stripWhitespace)`
- `# tells R to treat your preprocessed documents as text documents.`
- `docs <- tm_map(docs, PlainTextDocument)`

Step of Text Mining

- # Create Document Term Matrix
- `dtm <- DocumentTermMatrix(docs)`
- `dtm`
- # Create Term Document Matrix
- `tdm <- TermDocumentMatrix(docs)`
- `tdm`

Explore Data

- `# Explore Data`
- `freq <- colSums(as.matrix(dtm))`
- `length(freq)`
- `ord <- order(freq)`
- `# Start by removing sparse terms:`
- `dtms <- removeSparseTerms(dtm, 0.1) # This makes a matrix that is 10% empty space, maximum.`
- `inspect(dtms)`
- `freq[head(ord)]`
- `freq[tail(ord)]`
- `head(table(freq), 20)`

Explore Data

- `# we can view a table of the terms we selected when we removed sparse terms`
- `freq <- colSums(as.matrix(dtms))`
- `freq <- sort(colSums(as.matrix(dtm)), decreasing=TRUE)`
- `head(freq, 14)`
- `findFreqTerms(dtm, lowfreq=500)`
- `wf <- data.frame(word=names(freq), freq=freq)`
- `head(wf)`

Create Word Cloud

- `library(ggplot2)`
- `p <- ggplot(subset(wf, freq>500), aes(word, freq))`
- `p <- p + geom_bar(stat="identity")`
- `p <- p + theme(axis.text.x=element_text(angle=45, hjust=1))`
- `p`

Hadoop and Map-reduce Paradigm

Large-Scale Data Analytics

- MapReduce computing paradigm (E.g., Hadoop) vs. Traditional database systems



vs.

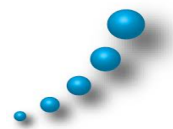


- Many enterprises are turning to Hadoop
 - Especially applications generating big data
 - Web applications, social networks, scientific applications

Why Hadoop is able to compete?



vs.



Scalability (petabytes of data, thousands of machines)



Flexibility in accepting all data formats (no schema)



Efficient and simple fault-tolerant mechanism



Commodity inexpensive hardware



Performance (tons of indexing, tuning, data organization tech.)



Features:

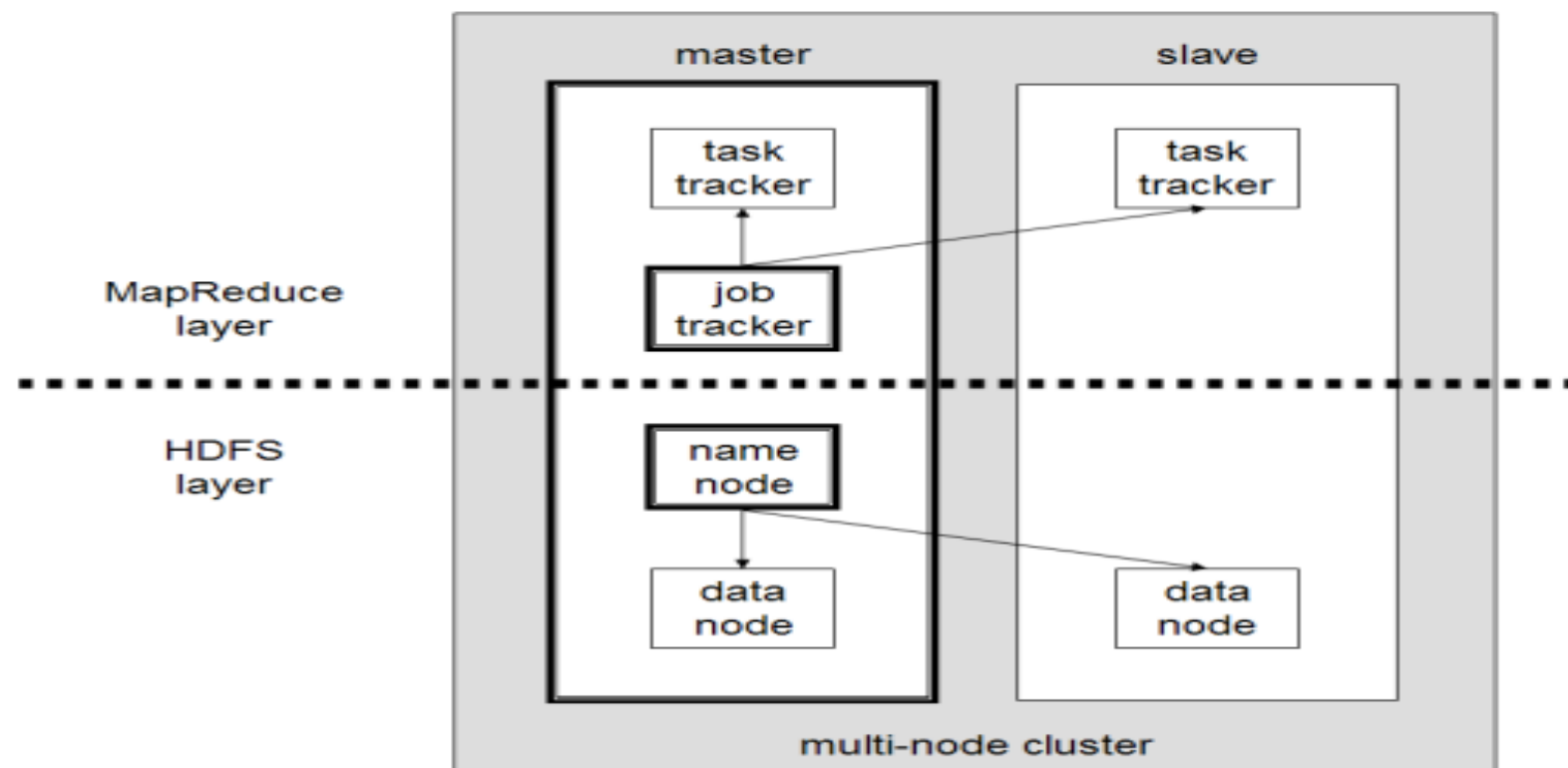
- Provenance tracking
- Annotation management
-

What is Hadoop?

- Hadoop is a software framework for *distributed processing* of *large datasets* across *large clusters* of computers
 - *Large datasets* → Terabytes or petabytes of data
 - *Large clusters* → hundreds or thousands of nodes
- Hadoop is open-source implementation for Google **MapReduce**
- Hadoop is based on a simple programming model called *MapReduce*
- Hadoop is based on a simple data model, *any data will fit*

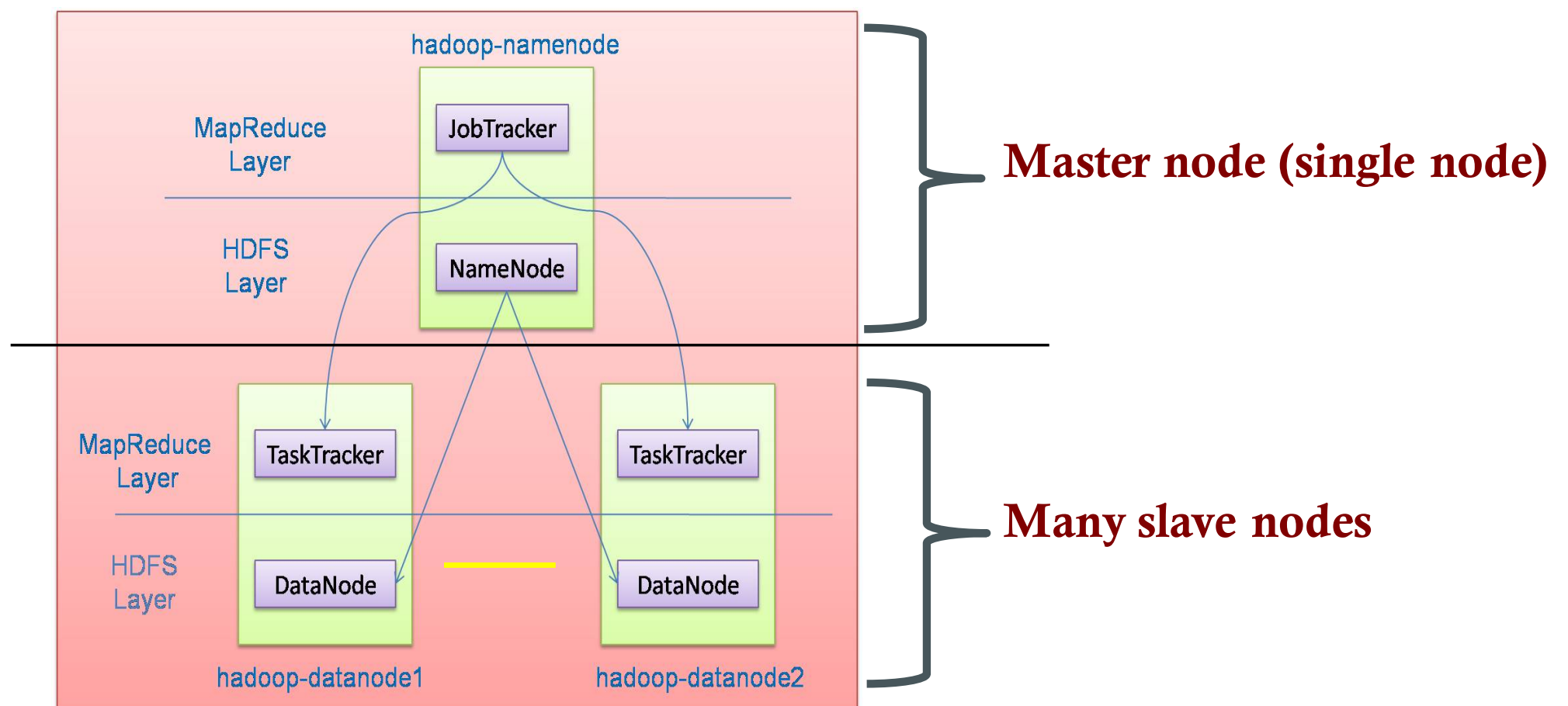
What is Hadoop?

- **Hadoop framework consists on two main layers**
 - Distributed file system (HDFS)
 - Execution engine (MapReduce)



Hadoop Architecture

- Hadoop is designed as a *master-slave shared-nothing* architecture



Design principles of Hadoop

- Need to process big data
- Need to parallelize computation across thousands of nodes
- **Commodity hardware**
 - Large number of low-end cheap machines working in parallel to solve a computing problem
- This is in contrast to **Parallel DBs**
 - Small number of high-end expensive machines

Design principles of Hadoop

- **Automatic parallelization & distribution**
 - Hidden from the end-user
- **Fault tolerance and automatic recovery**
 - Nodes/tasks will fail and will recover automatically
- **Clean and simple programming abstraction**
 - Users only provide two functions “map” and “reduce”

RHadoop

- `install.packages(c('rJava','RJSONIO', 'itertools',
'digest','Rcpp','httr','functional','devtools', 'plyr','reshape2'))`
- `Sys.setenv("HADOOP_CMD"="/usr/local/Cellar/hadoop/2.6.0/bin/
hadoop")`
- `Sys.setenv("HADOOP_STREAMING"="/usr/local/Cellar/hadoop/
2.6.0/libexec/share/hadoop/tools/lib/hadoop-streaming-2.6.0.jar")`
- `Sys.getenv("HADOOP_CMD")`
- `Sys.setenv("HADOOP_HOME"="/usr/local/Cellar/hadoop/2.6.0")`

Install RHadoop

- Installing RHadoop [rhdfs, rmr, rhbase]
 1. Download RHadoop packages from GitHub repository of Revolution Analytics: <https://github.com/RevolutionAnalytics/RHadoop>
 - rmr: [rmr-2.2.2.tar.gz]
 - rhdfs: [rhdfs-1.6.0.tar.gz]
 - rhbase: [rhbase-1.2.0.tar.gz]
 2. Installing packages.
 - For rmr we use: R CMD INSTALL rmr-2.2.2.tar.gz
 - For rhdfs we use: R CMD INSTALL rhdfs-1.6.0.tar.gz
 - For rhbase we use: R CMD INSTALL rhbase-1.2.0.tar.gz

gdp data

- `library(rmr2)`
- `library(rhdfs)`
- `gdp <- NA`
- `gdp <- read.csv("~/Downloads/GDP.csv")`
- `gdp <- gdp[,1:4]`
- `gdp$GDP <- as.double(gsub(",", "", gdp$GDP))`
- `head(gdp)`

Setup Map-Reduce Function

- `hdfs.init()`
- `gdp.values <- to.dfs(gdp)`
- `aaplRevenue = 181890`
- `gdp.map.fn <- function(k,v) {`
- `key <- ifelse(v[4] < aaplRevenue, "less", "greater")`
- `keyval(key, 1)`
- `}`
- `count.reduce.fn <- function(k,v) {`
- `keyval(k, length(v))`
- `}`

Run Map-Reduce Function

- `count <- mapreduce(input=gdp.values, map = gdp.map.fn, reduce = count.reduce.fn)`
- `from.dfs(count)`

Q&A