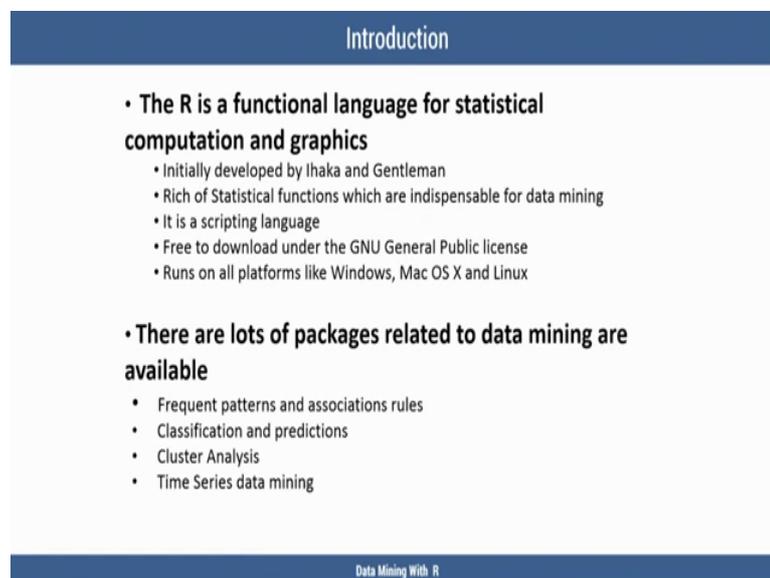


Data Mining
Prof. Arindam Dasgupta
Department of Computer Science and Engineering
Indian Institute of Technology, Kharagpur

Lecture – 43
Tutorial

Hello. Welcome to data mining with R course. I am going to present how to apply the R programming language in application of data mining, I am a research scholar in compressors department, my name is Arindam Dasgupta.

(Refer Slide Time: 00:40)



The slide is titled "Introduction" and contains the following text:

- **The R is a functional language for statistical computation and graphics**
 - Initially developed by Ihaka and Gentleman
 - Rich of Statistical functions which are indispensable for data mining
 - It is a scripting language
 - Free to download under the GNU General Public license
 - Runs on all platforms like Windows, Mac OS X and Linux
- **There are lots of packages related to data mining are available**
 - Frequent patterns and associations rules
 - Classification and predictions
 - Cluster Analysis
 - Time Series data mining

Data Mining With R

Now, one of what is R? Basically, R is a programming language, it is basically a functional language for statistical computation and graphics, initially, it is developed by Ihaka and gentleman, it is a riched with lots of statistical functions which are very essential for data mining applications. Basically, it is a scripting language and this tool can be downloaded freely from under the GNU license and you can execute run or run the R programming in any platform such as windows MAC, OS, X and Linux because it supports all the platforms. R is very beautiful language because it is basically functional language each functions are already implemented in the R library and we just use that functions for our pass passes.

There are lots of packages data mining available data mining packages are available these packages are frequent pattern analysis and rules as described in the theoretical class

and classification prediction algorithms cluster analysis packages, there are time series data mining packages there are lots of packages available just you have to download the package in real time and just apply the functions on that packages.

(Refer Slide Time: 02:23)

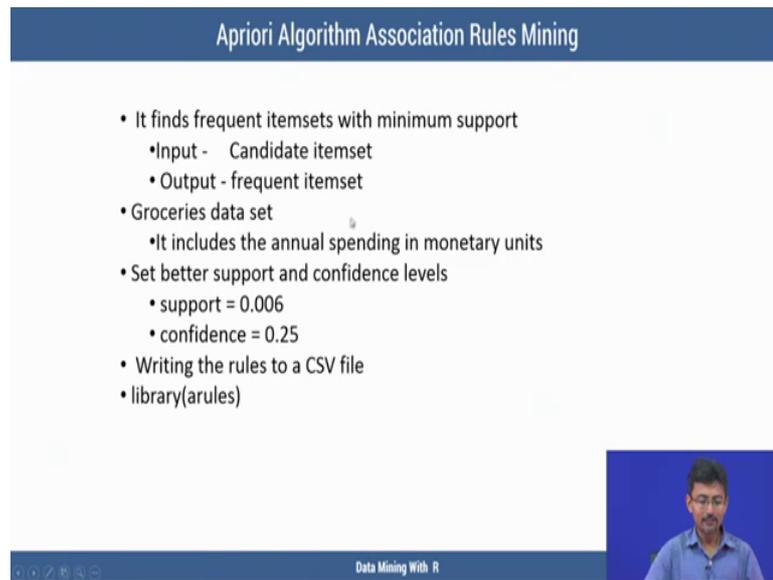
Software Installation

- Installing R in Windows Platform
 - Go to <https://cran.r-project.org/bin/windows/base/>
 - Under "Download and Install R", click on the "Windows" link
 - Now install R-2.10.1-win32.exe
- Installing R Studio
 - It is a free and open source IDE (integrated development environment) for R
 - It is available at <http://www.rstudio.com>

Data Mining With R

Now, how to get the R tool? To get the R tool, we have to go to the this website is called this website is called cran project org bin windows or base just go to this link at first, then download the R according to your platform in my case, these windows platform that is why I have downloaded R win 32 exe after downloading the software, we have to install the software, it is very easy to install just click on the R exe icon, I will show you the this is the this is the R software, you just double click on; this icon and just install by applying the next I have already installed here.

(Refer Slide Time: 03:48)



Apriori Algorithm Association Rules Mining

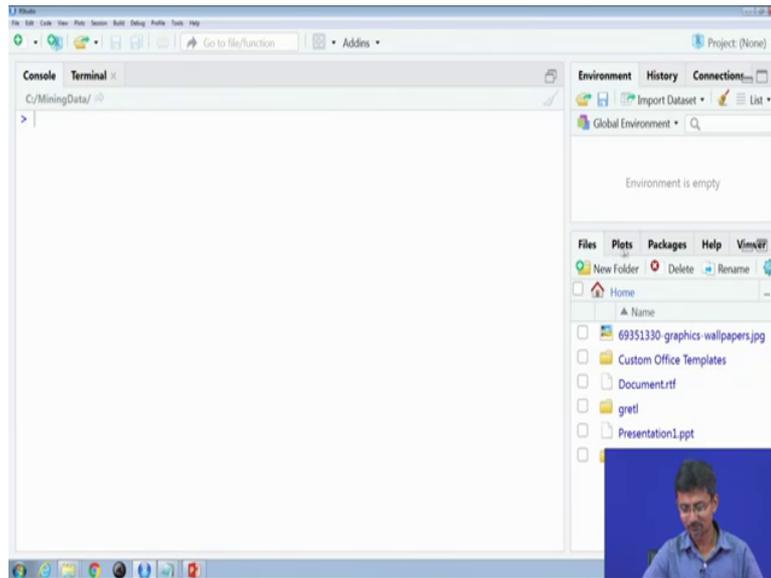
- It finds frequent itemsets with minimum support
 - Input - Candidate itemset
 - Output - frequent itemset
- Groceries data set
 - It includes the annual spending in monetary units
- Set better support and confidence levels
 - support = 0.006
 - confidence = 0.25
- Writing the rules to a CSV file
- library(arules)

Data Mining With R

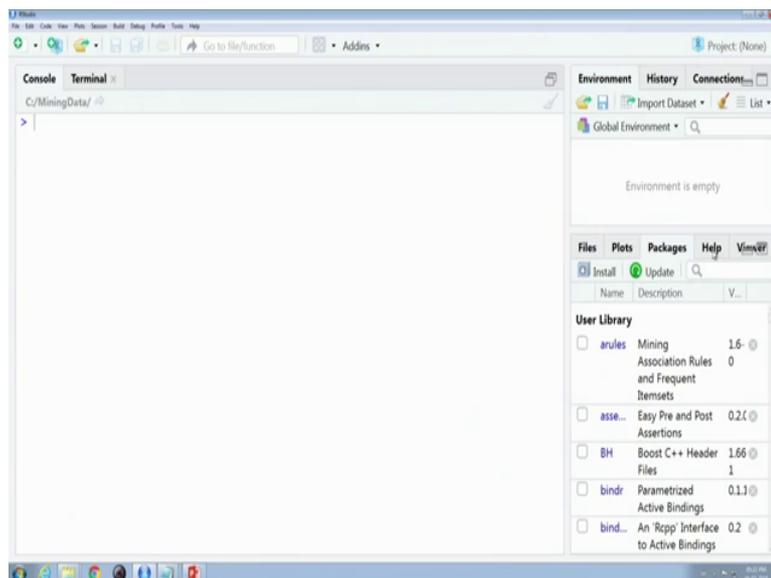
After that after installing the R studio, we have to install another software application that is basically intergrated development environment it is also free and open source id, it is called R studio, it is very user friendly, it provides an user friendly interface to execute your R programming language. I will show you the how to what is the interface of the R studio. These are interface of R studio just I am basically, this part of the top left one part of this area is used to write the programming language R programming language the bottom part is the basically the console of R console of R executing the R script.

Basically, this is the editor and these are R console and right hand side will display, this part top right hand side this part will display the local variables the which are just defined in the R terminal and the bottom of portion you can get the what is this is the files which are in this computer and any plot will be displayed here.

(Refer Slide Time: 05:19)



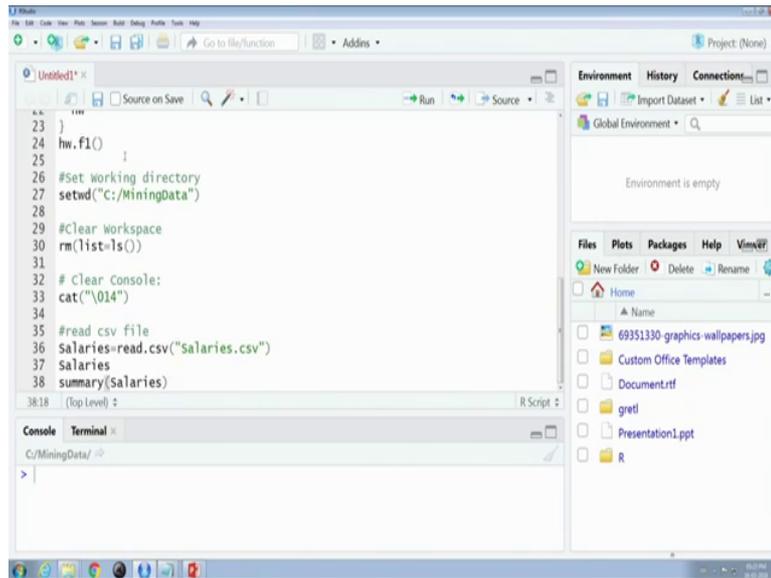
(Refer Slide Time: 05:28)



These are package list; these are which can be used in my program in our program and if any R resource is needed then you can help get help from this link.

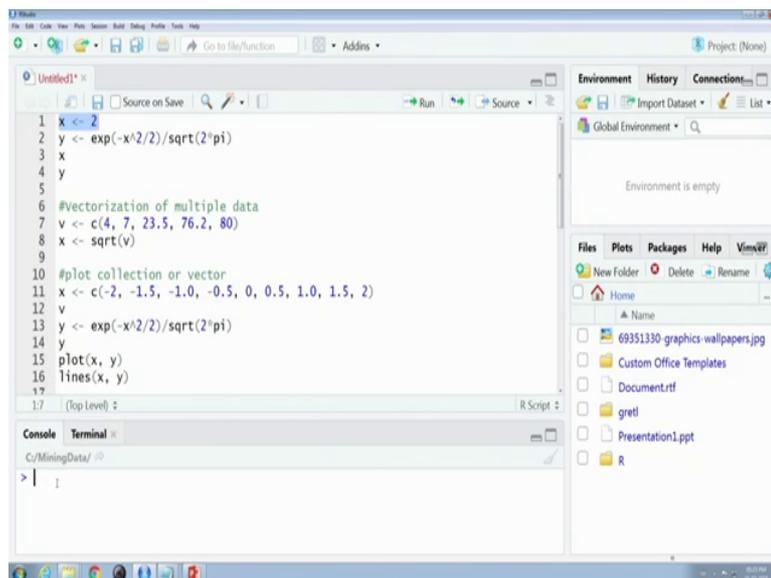
And now I want to show you some basic R program at first to click file create a 1 R script new file, I have created one new file and I have just put the script here each scripts are suppose first line in first line of the script, it is.

(Refer Slide Time: 06:31)



```
23 }
24 hw.fl()
25
26 #Set working directory
27 setwd("C:/MiningData")
28
29 #Clear workspace
30 rm(list=ls())
31
32 # Clear Console:
33 cat("\014")
34
35 #read csv file
36 salaries=read.csv("Salaries.csv")
37 salaries
38 summary(salaries)
```

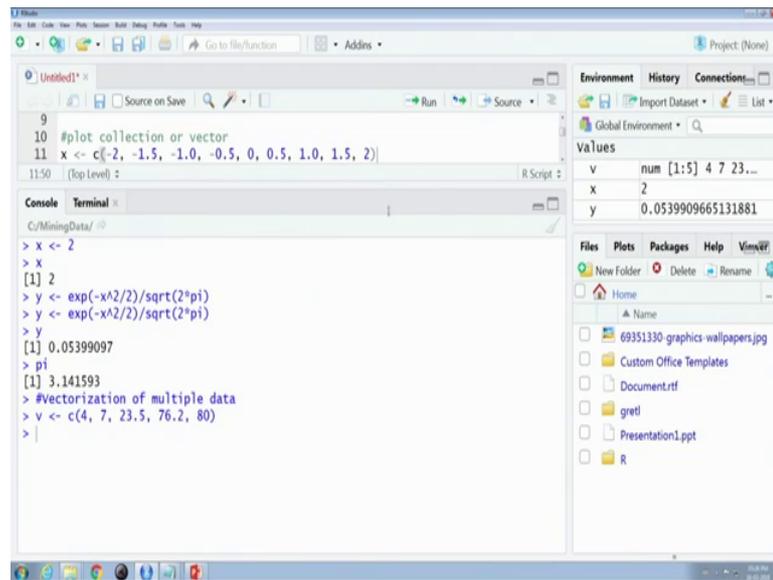
(Refer Slide Time: 06:33)



```
1 x <- 2
2 y <- exp(-x^2/2)/sqrt(2*pi)
3 x
4 y
5
6 #Vectorization of multiple data
7 v <- c(4, 7, 23.5, 76.2, 80)
8 x <- sqrt(v)
9
10 #plot collection or vector
11 x <- c(-2, -1.5, -1.0, -0.5, 0, 0.5, 1.0, 1.5, 2)
12 v
13 y <- exp(-x^2/2)/sqrt(2*pi)
14 y
15 plot(x, y)
16 lines(x, y)
17
```

It is showing x left arrow two here x is the object basically, not is exactly like very well, basically, it is a object it stores the very value 2 in the x variable.

(Refer Slide Time: 07:16)

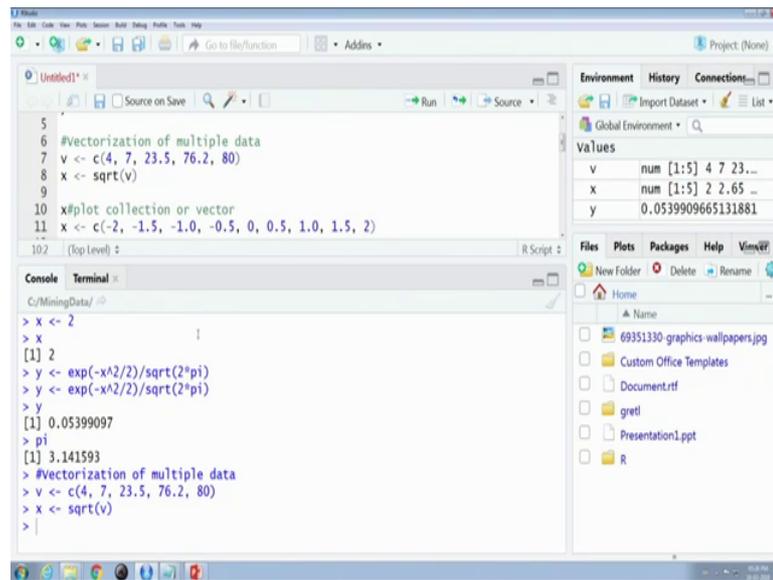


Suppose I am to execute this in R console. Now, I want to execute this after executing this value two is has been assigned to object x, I want to view this value just type x here and write which showing the value of x is 2 ok.

Then I want to assign one expression, I am executing directly from my source script by clicking here just you have to just put the cursor here and click run, it will execute a particular command you see these executing this particular command is written here from this portion to this portion and now I want to view; what is the value of y, this is the value of y, you see, this is a basically expression e to the power minus x square by 2 and by sqrt 2 star pi, if I want to know that what is the value of pi you just type the pi value it is constant; it is already defined in R and then what is the;

Now, I have to show how to store the multiple datas into a particular object suppose I am to executes this v these vectorization, here see first bracket start within the bracket, there are lots of values 4 7; these are the values different types of values here; basically it is a it sees for collection it is a collection of datas; this data has been stored into the object v.

(Refer Slide Time: 09:54)



```
5
6 #Vectorization of multiple data
7 v <- c(4, 7, 23.5, 76.2, 80)
8 x <- sqrt(v)
9
10 #plot collection or vector
11 x <- c(-2, -1.5, -1.0, -0.5, 0, 0.5, 1.0, 1.5, 2)
12
13
```

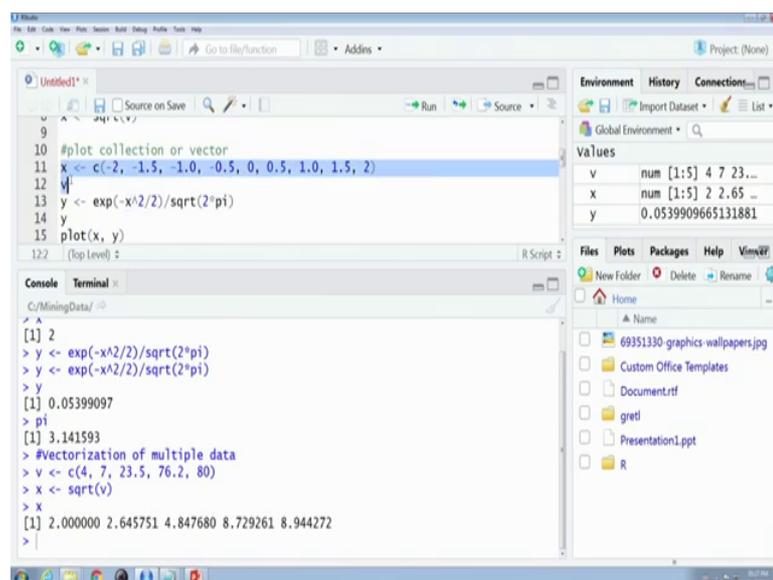
Console Terminal

```
C:/MiningData/ >>
> x <- 2
> x
[1] 2
> y <- exp(-x^2/2)/sqrt(2*pi)
> y <- exp(-x^2/2)/sqrt(2*pi)
> y
[1] 0.05399097
> pi
[1] 3.141593
> #Vectorization of multiple data
> v <- c(4, 7, 23.5, 76.2, 80)
> x <- sqrt(v)
> |
```

Variable	Value
v	num [1:5] 4 7 23.5 76.2 80
x	num [1:5] 2 2.65 4.85 8.72 8.94
y	0.0539909665131881

Then after then after storing the value into the v, I want to display the value of v into the x into x. Now see now each of the value has been square rooted. Now check the value of x here, it is showing the each value is square rooted; that means, the vectorization on multiple data and the square root function has been applied in each of the data into the collections.

(Refer Slide Time: 10:49)



```
9
10 #plot collection or vector
11 x <- c(-2, -1.5, -1.0, -0.5, 0, 0.5, 1.0, 1.5, 2)
12 y <- exp(-x^2/2)/sqrt(2*pi)
13 y <- exp(-x^2/2)/sqrt(2*pi)
14 y
15 plot(x, y)
16
```

Console Terminal

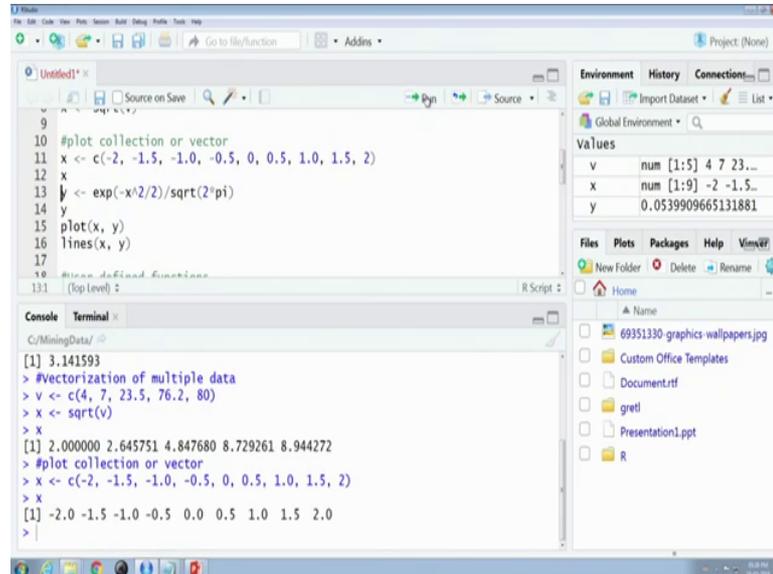
```
C:/MiningData/ >>
> x
[1] 2
> y <- exp(-x^2/2)/sqrt(2*pi)
> y <- exp(-x^2/2)/sqrt(2*pi)
> y
[1] 0.05399097
> pi
[1] 3.141593
> #Vectorization of multiple data
> v <- c(4, 7, 23.5, 76.2, 80)
> x <- sqrt(v)
> x
[1] 2.000000 2.645751 4.847680 8.729261 8.944272
> |
```

Variable	Value
v	num [1:5] 4 7 23.5 76.2 80
x	num [1:5] 2 2.65 4.85 8.72 8.94
y	0.0539909665131881

Now, I want to plot how to plot multiple vector data; suppose, first I store a collection of data into the x variable these the collection of data, then I will display the x variable, then

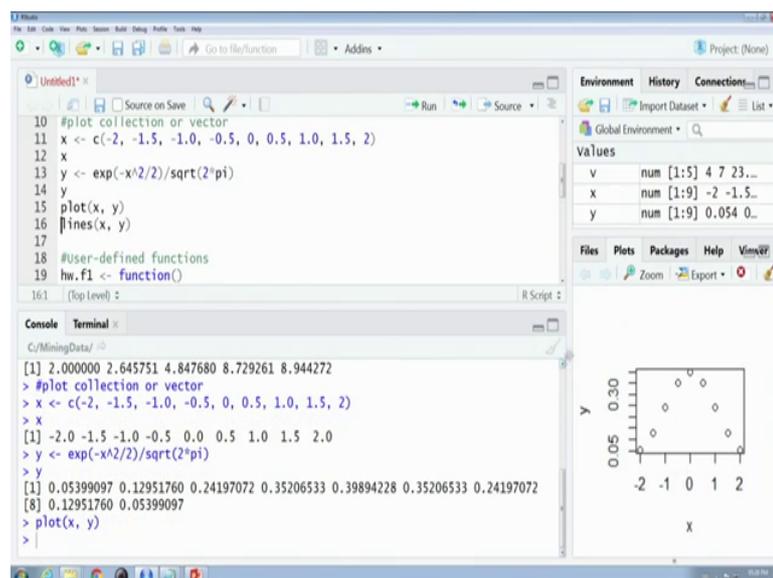
I apply this formula into each of data set here in x then plot, then I will display the value of y, then I plot the data set.

(Refer Slide Time: 11:47)



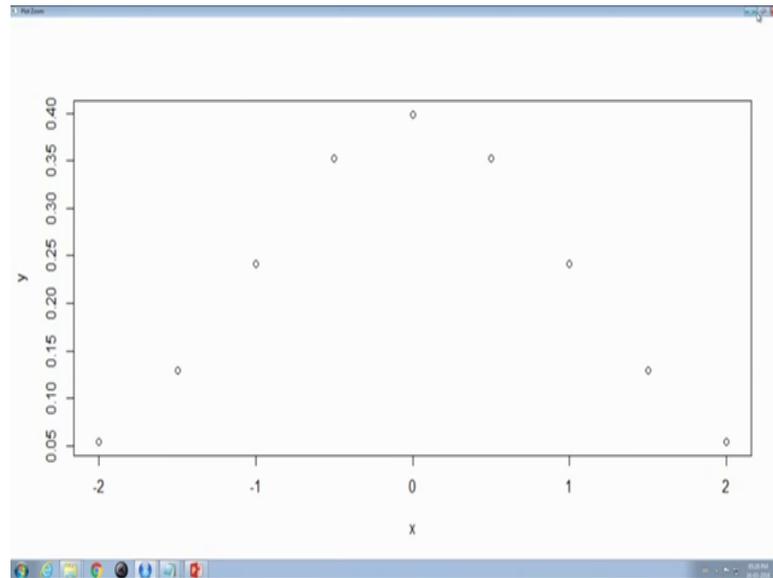
Now, start now it is already stored, then it is displaying the value of x object. Now, it is displaying, it is computing the values of y object, then it is displaying the values of y objects, I want to plot x data of x and y just for plotting the data of x and y just use plot within x axis and y axis.

(Refer Slide Time: 12:23)



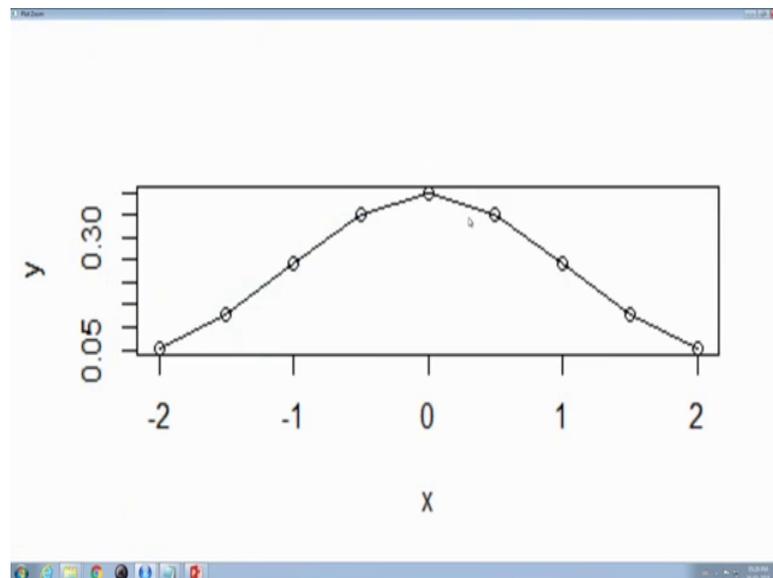
Now, it is plotting here.

(Refer Slide Time: 12:30)



It is showing; these are x axis and these are the y axis, I want to know I want to connect these points with line just typed lines there is a command lines. Now it is connected.

(Refer Slide Time: 13:00)

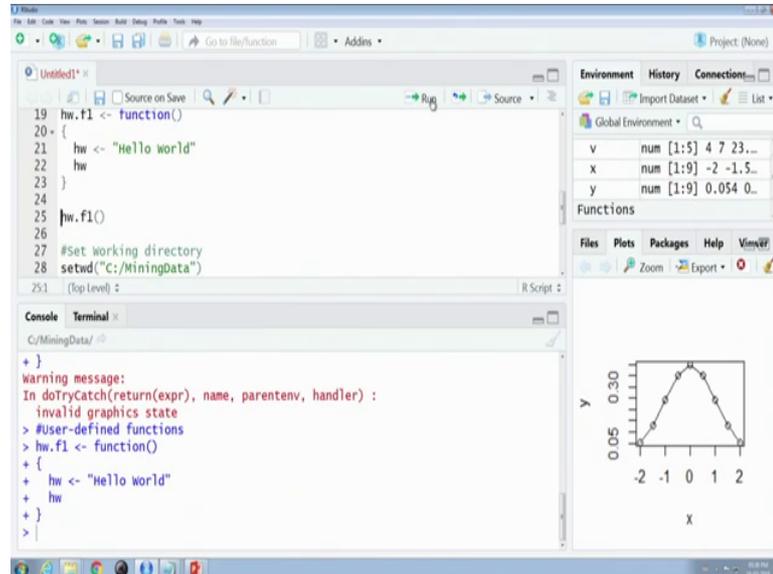


These are the very basic implementation of plotting which will be used in data mining application.

Now, we can define functions in R programming; these are the structure just write a function name here and then define it is a function the function name is this and function definition, I just generate use one object hw, then type insert one string here in this object

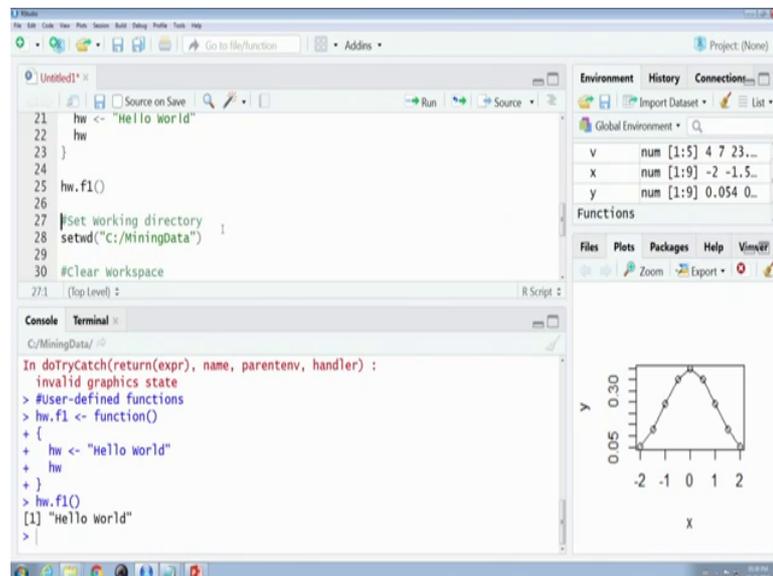
and display the string and this is the this is the calling of function. Now at first I have to define the function. Now function is defined. Now I want to call the function.

(Refer Slide Time: 14:08)



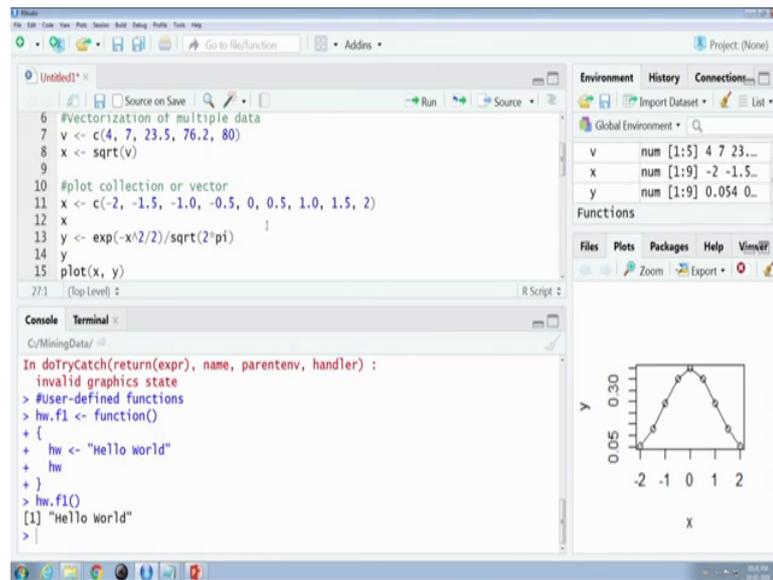
Here it is showing the hello world after by calling the function.

(Refer Slide Time: 14:21)



Now, we have learned how to store the values into the objects.

(Refer Slide Time: 14:26)



And how to plot data, then I want to access one data file for data mining, you have to access multiple data files or a single file for the analysis purpose for that we have to set one working directory in warren in in my working directory I have stored multiple data files in c drive, I have created one folder data mining data and within these folder I have stored lots of csv files for analysis these are the csv files. Now, I want to retrieve the csv files through R programming.

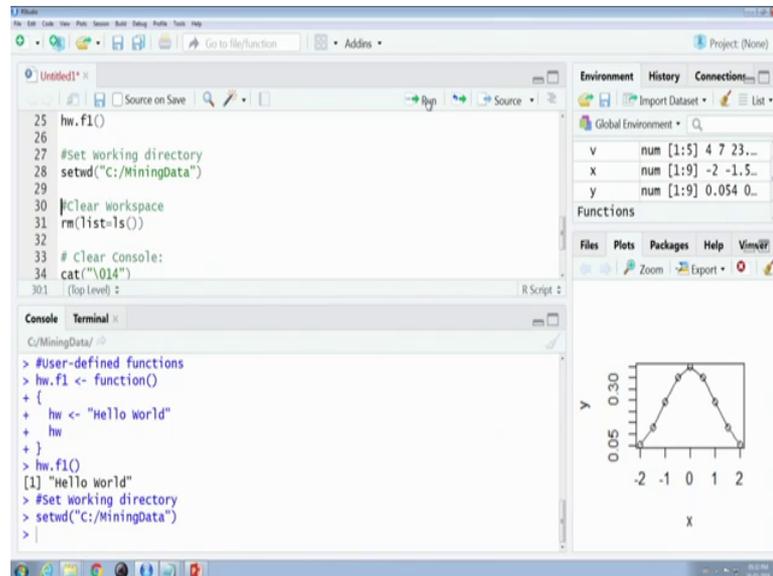
(Refer Slide Time: 15:28)

The screenshot shows an Excel spreadsheet with the following data:

Rank	Discipline	Yrs	Score	Salary
1	Prof	25	38	53750
2	Prof	20	38	57000
3	AsstProf	4	3	7670
4	Prof	45	39	52000
5	Prof	40	43	54200
6	AsstProf	6	4	8700
7	Prof	30	23	57000
8	Prof	45	45	54700
9	Prof	35	20	53700
10	Prof	38	38	52000
11	Prof	38	38	52000
12	AsstProf	13	8	13800
13	AsstProf	7	2	7600
14	AsstProf	5	3	7700
15	AsstProf	2	0	7600
16	Prof	30	10	54000
17	Prof	12	3	51700
18	Prof	19	20	50000
19	Prof	38	38	53000
20	Prof	17	23	54700
21	Prof	39	38	53700
22	Prof	35	20	53000
23	Prof	36	35	50500
24	Prof	34	30	50000
25	Prof	34	19	53000
26	AsstProf	13	8	14000
27	Prof	25	8	50000
28	Prof	30	20	54000
29	AsstProf	5	3	8270
30	AsstProf	11	0	7700
31	Prof	12	0	50000
32	Prof	20	4	50000
33	AsstProf	7	2	7600
34	Prof	18	9	51700
35	AsstProf	4	2	8020
36	Prof	4	2	8020
37	AsstProf	3	0	7700
38	Prof	22	25	53700
39	AsstProf	7	4	8070
40	Prof	40	10	51700
41	AsstProf	9	9	50000
42	Prof	23	2	54000
43	AsstProf	19	10	9400

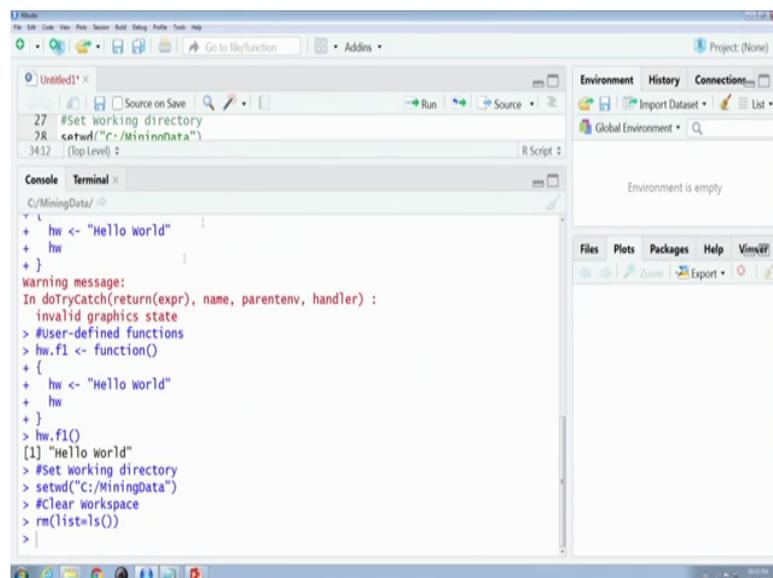
At first, I have to set the working directory because my data files are stored in mining data folder under c drive. Now I execute this instruction.

(Refer Slide Time: 15:59)



Then I have to clear my workspace these workspace these clear these variables and for that you have to use rm list equals to ls now execute this query now my environment is cleared after that clear my this console these R console.

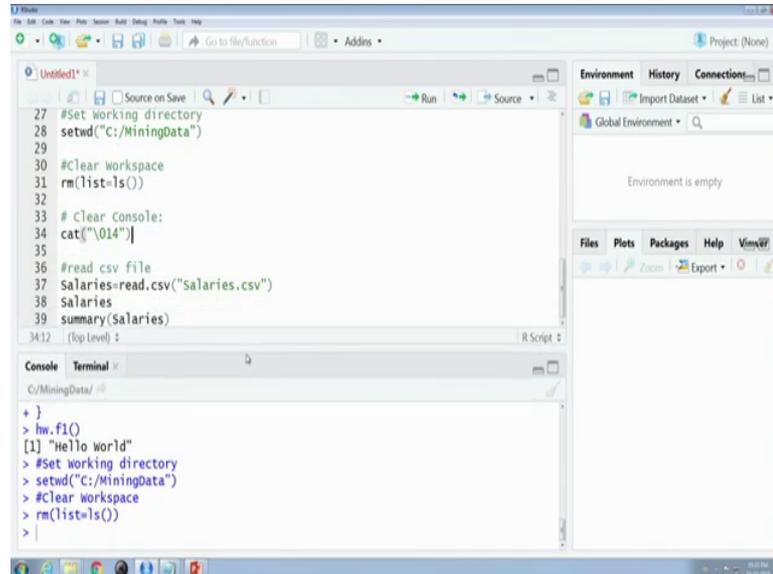
(Refer Slide Time: 16:28)



Clear this console, you have to use this command is cat slash 0 1 4 within string, just execute this command for clearing your console. Now my console is cleared now.

Now, second type we have to read there is a file salaries dot csv in my data mining directory we have to read this file.

(Refer Slide Time: 16:31)



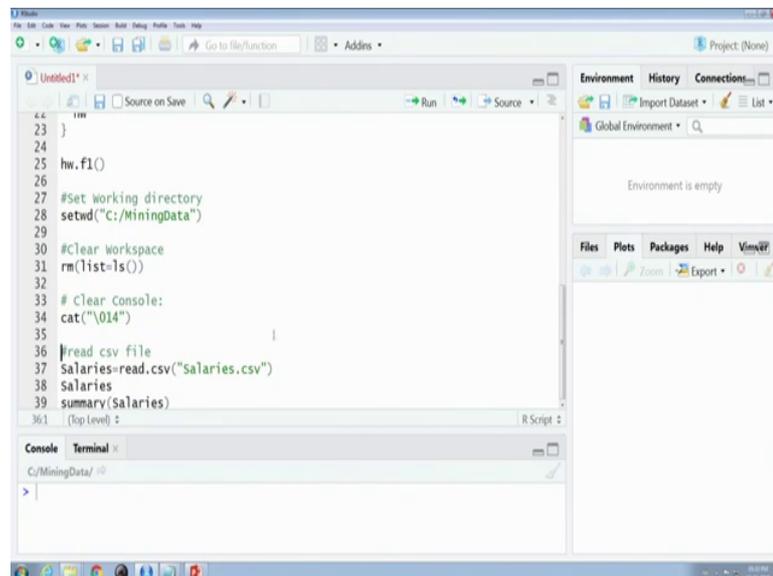
```
27 #Set working directory
28 setwd("C:/MiningData")
29
30 #Clear workspace
31 rm(list=ls())
32
33 # Clear Console:
34 cat("\014")
35
36 #read csv file
37 salaries=read.csv("Salaries.csv")
38 salaries
39 summary(salaries)
```

34:12 (Top Level) R Script

```
C:/MiningData/ >
> }
> hw.fl()
[1] "Hello world"
> #Set working directory
> setwd("C:/MiningData")
> #clear workspace
> rm(list=ls())
>
```

Through this command there is inbuilt function read dot csv.

(Refer Slide Time: 16:44)



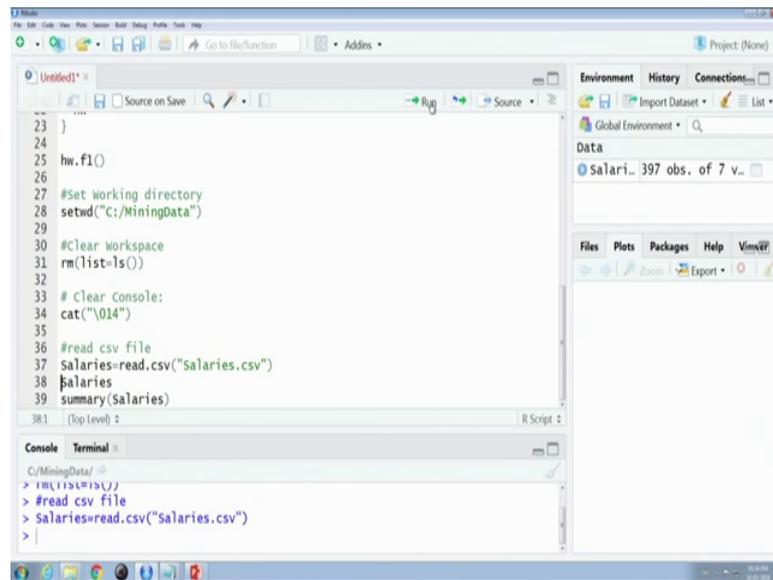
```
23 }
24
25 hw.fl()
26
27 #Set working directory
28 setwd("C:/MiningData")
29
30 #Clear workspace
31 rm(list=ls())
32
33 # Clear Console:
34 cat("\014")
35
36 #read csv file
37 salaries=read.csv("Salaries.csv")
38 salaries
39 summary(salaries)
```

36:1 (Top Level) R Script

```
C:/MiningData/ >
```

It is meant for accessing the csv file and after executing this command entire data will be stored in salary objects this salary object now salary objects is.

(Refer Slide Time: 17:21)



The screenshot shows the RStudio interface. The main editor window contains the following R script:

```
23 }
24
25 hw.fl()
26
27 #Set working directory
28 setwd("c:/MiningData")
29
30 #Clear workspace
31 rm(list=ls())
32
33 # Clear console:
34 cat("\014")
35
36 #read csv file
37 salaries=read.csv("Salaries.csv")
38 summary(salaries)
39
```

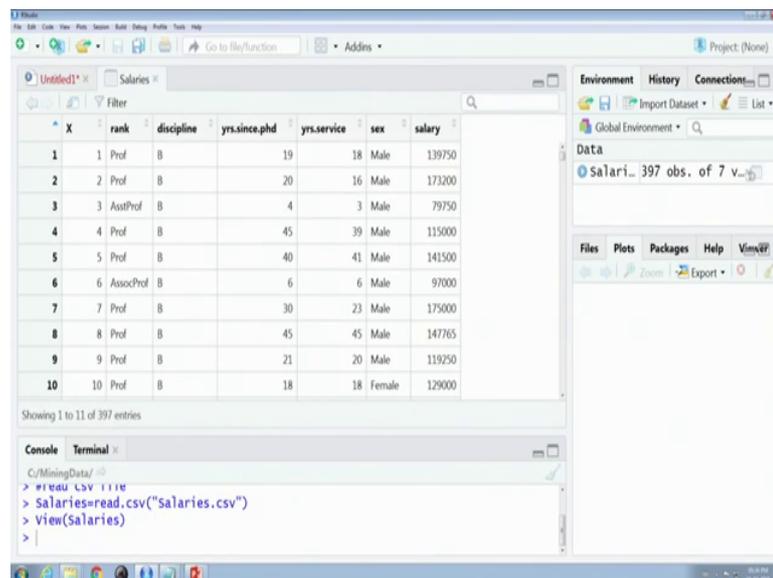
The console window shows the execution of the script:

```
C:/MiningData/ > rm(list=ls())
> #read csv file
> salaries=read.csv("Salaries.csv")
>
```

The Environment pane on the right shows the Global Environment with a data object named 'Salari...' containing 397 observations of 7 variables.

These salary objects you see destroying the tables.

(Refer Slide Time: 17:25)



The screenshot shows the RStudio interface with the 'Salaries' data frame loaded. The main editor window displays a table with the following columns: X, rank, discipline, yrs.since.phd, yrs.service, sex, and salary. The first 10 rows are shown:

X	rank	discipline	yrs.since.phd	yrs.service	sex	salary	
1	1	Prof	B	19	18	Male	139750
2	2	Prof	B	20	16	Male	173200
3	3	AsstProf	B	4	3	Male	79750
4	4	Prof	B	45	39	Male	115000
5	5	Prof	B	40	41	Male	141500
6	6	AssocProf	B	6	6	Male	97000
7	7	Prof	B	30	23	Male	175000
8	8	Prof	B	45	45	Male	147765
9	9	Prof	B	21	20	Male	119250
10	10	Prof	B	18	18	Female	129000

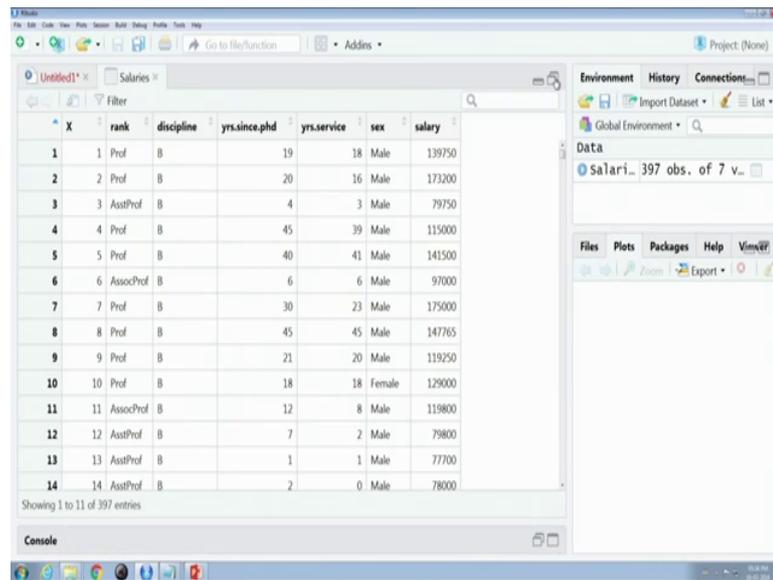
The console window shows the execution of the script:

```
C:/MiningData/ > #read csv file
> salaries=read.csv("Salaries.csv")
> View(salaries)
>
```

The Environment pane on the right shows the Global Environment with a data object named 'Salari...' containing 397 observations of 7 variables.

Now, I want to display the salary objects in my R console.

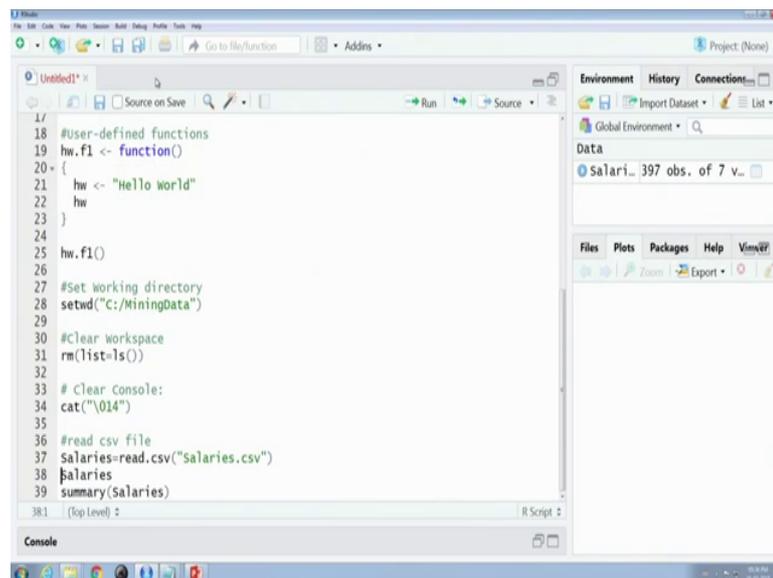
(Refer Slide Time: 17:28)



The screenshot shows the RStudio interface with a data table loaded. The table has 14 rows and 7 columns. The columns are labeled X, rank, discipline, yrs.since.phd, yrs.service, sex, and salary. The data is as follows:

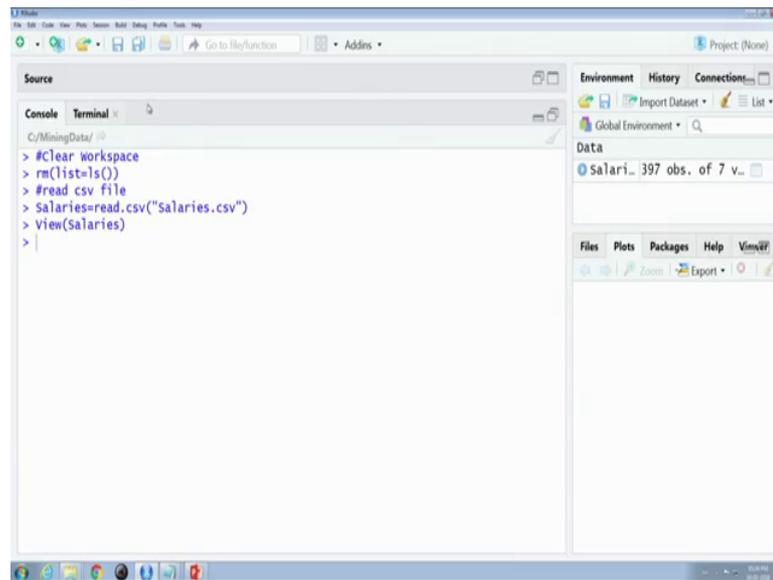
X	rank	discipline	yrs.since.phd	yrs.service	sex	salary	
1	1	Prof	B	19	18	Male	139750
2	2	Prof	B	20	16	Male	173200
3	3	AsstProf	B	4	3	Male	79750
4	4	Prof	B	45	39	Male	115000
5	5	Prof	B	40	41	Male	141500
6	6	AssocProf	B	6	6	Male	97000
7	7	Prof	B	30	23	Male	175000
8	8	Prof	B	45	45	Male	147765
9	9	Prof	B	21	20	Male	119250
10	10	Prof	B	18	18	Female	129000
11	11	AssocProf	B	12	8	Male	119800
12	12	AsstProf	B	7	2	Male	79800
13	13	AsstProf	B	1	1	Male	77700
14	14	AsstProf	B	2	0	Male	78000

(Refer Slide Time: 17:48)



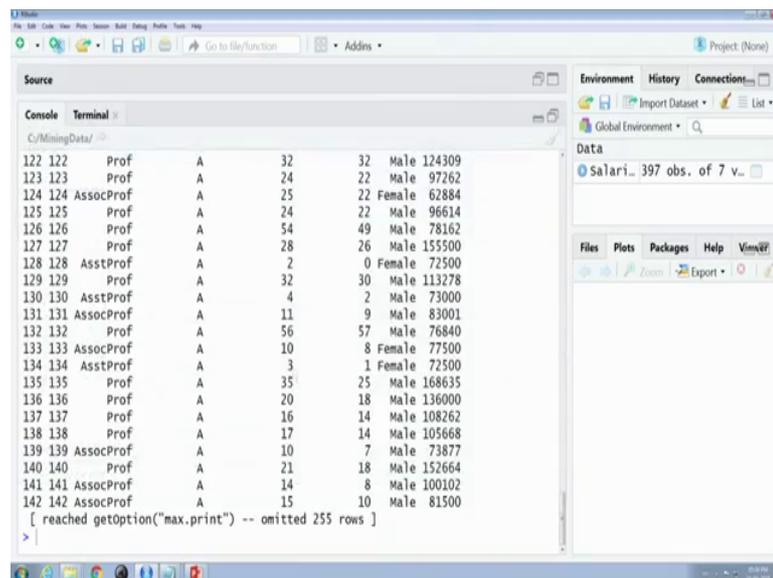
```
17
18 #User-defined functions
19 hw.f1 <- function()
20 {
21   hw <- "Hello world"
22   hw
23 }
24
25 hw.f1()
26
27 #Set working directory
28 setwd("c:/MiningData")
29
30 #Clear workspace
31 rm(list=ls())
32
33 # Clear Console:
34 cat("\n014")
35
36 #read csv file
37 salaries=read.csv("Salaries.csv")
38 salaries
39 summary(salaries)
```

(Refer Slide Time: 18:04)



```
C:/MiningData/ > #Clear workspace
> rm(list=ls())
> #read csv file
> Salaries=read.csv("Salaries.csv")
> View(Salaries)
>
```

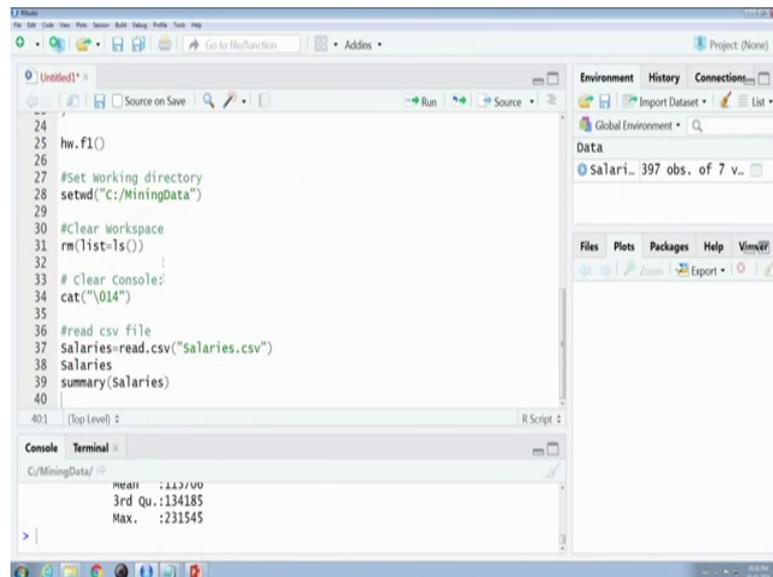
(Refer Slide Time: 18:13)



```
C:/MiningData/ 122 122 Prof A 32 32 Male 124309
123 123 Prof A 24 22 Male 97262
124 124 AssocProf A 25 22 Female 62884
125 125 Prof A 24 22 Male 96614
126 126 Prof A 54 49 Male 78162
127 127 Prof A 28 26 Male 155500
128 128 AsstProf A 2 0 Female 72500
129 129 Prof A 32 30 Male 113278
130 130 AsstProf A 4 2 Male 73000
131 131 AssocProf A 11 9 Male 83001
132 132 Prof A 56 57 Male 76840
133 133 AssocProf A 10 8 Female 77500
134 134 AsstProf A 3 1 Female 72500
135 135 Prof A 35 25 Male 168635
136 136 Prof A 20 18 Male 136000
137 137 Prof A 16 14 Male 108262
138 138 Prof A 17 14 Male 105668
139 139 AssocProf A 10 7 Male 73877
140 140 Prof A 21 18 Male 152664
141 141 AssocProf A 14 8 Male 100102
142 142 AssocProf A 15 10 Male 81500
[ reached getOption("max.print") -- omitted 255 rows ]
>
```

Now, it is showing in R console suppose, I want to summarize this salary summarize means, it will display the; suppose, there is a very well called discipline in this very well, there are 2 values A and B, it contains 181 values of A and 2 16 values of B. Similarly, these are basically nominal value and the numeric values in numeric values, it is showing the what is the minimum value what is the maximum value what is a quartile what is the median value? Mean value of each attributes in salary the minimum salary is this maximum salary is this and mean salary is this and it is it is it is showing all the summary information summary information of each attribute.

(Refer Slide Time: 19:38)



```
24  
25 hw.f1()  
26  
27 #Set working directory  
28 setwd("c:/MiningData")  
29  
30 #Clear workspace  
31 rm(list=ls())  
32  
33 # Clear console:  
34 cat("\014")  
35  
36 #read csv file  
37 salaries=read.csv("salaries.csv")  
38 salaries  
39 summary(salaries)  
40
```

Environment History Connections
Global Environment
Data
Salari_ 397 obs. of 7 v...

Files Plots Packages Help View
Zoom Export

401 (Top Level) R Script

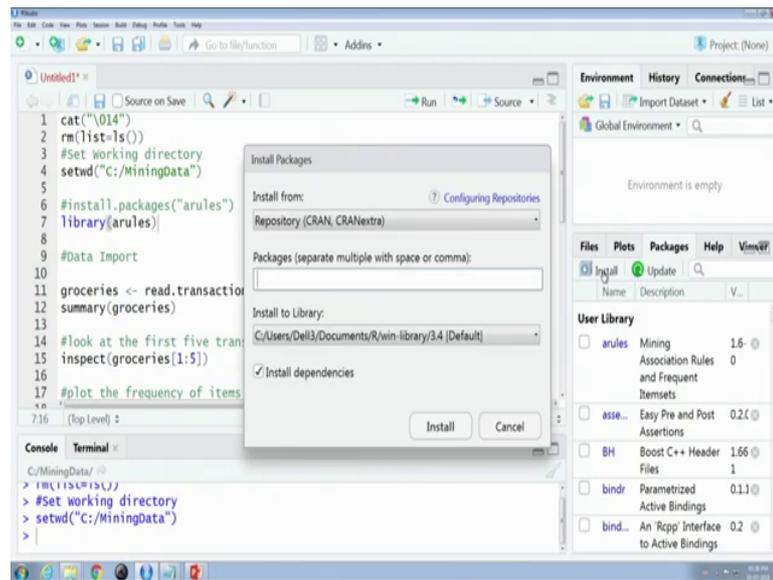
Console Terminal
C:/MiningData/
mean : 112700
3rd Qu.: 134185
Max. : 231545

Now, I will leave the basics of R programming which are needed to needed in any time any type of determining algorithms enough.

Now, at first I will discuss about how to use apriori algorithm for accessing association rules from a data set basically apriori algorithm, we have learned it finds the frequent data items with minimum support in I will show, I have a grocery data set and from that grocery data set, we just access the rules with support 0.006 and confidence 2.5 after generating the rules, we will store the rules into a csv file and for apriori algorithm there are the input is one file it is called candidates itemset and output is the frequent itemset ok.

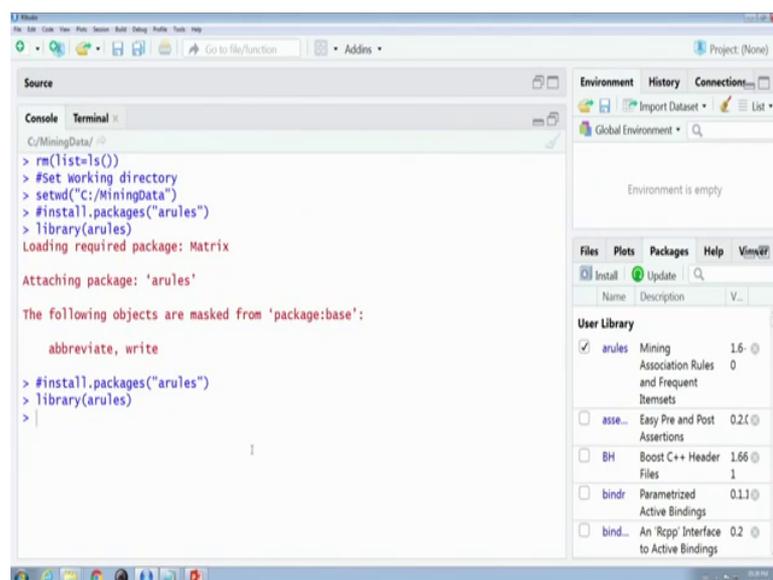
To execute apriori algorithm we have to use one package it is very essential it is called arules to download this package you just go to at right hand side click on install.

(Refer Slide Time: 21:24)



And type arules here arules here, it is displaying I selected the arules here. Now in install the arules packets from this terminal, after it is already installed in my system after installing the library installing the package, sorry, we have to include the library in R program for including library, just write library and type the package name and now I am including the library.

(Refer Slide Time: 22:24)



It is included; it is included smoothly

(Refer Slide Time: 22:52)

The screenshot shows an Excel spreadsheet with a single column of grocery items. The items listed include: citrus fruit, semi finished bread, tropical food, yogurt, coffee, whole milk, cream cheese, pip fruit, cream cheese, whole milk, butter, rice, abrasive cleaner, milk/burn, soft milk, milk/burn, bottled in liquid (appetizer), jarred plants, whole milk, tropical fruit, tropical fruit, whole milk, cereal, yogurt, flour, bottled in dishes, beer, frankfurter, milk/burn, chicken, tropical fruit, butter, sugar, fruit/veg newspapers, packaged fruit/vegetables, packaged fruit/vegetables, chocolate, specialty bar, other vegetables, butter milk, tropical fruit, cream cheese, detergent newspapers, tropical fruit, root veg, other veg, frozen, milk/burn, flour, sweet, soft, waffles, candy, bathroom cleaner, bottled in canned beer, yogurt, beverage, milk/burn, soda, chocolate, other vegetables, frozen bread, soda, fruit/veg, canned, bottled in shipping bags, yogurt, beverage, bottled in specialty bar, hamburger, meat, other veg, milk/burn, bottled in hygiene & napkins, root vegetables, other veg, whole milk, beverage, sugar, jam, berries, other veg, whole milk, waffles, soft, soda, abrasive cleaner, beer, grapes, detergent, pasta, soda, packaged fruit/vegetables, chocolate, canned beer, root vegetables, other veg, whole milk, dessert, citrus fruit, packaged newspapers, beverage, milk/burn, soda, canned, bottled in shipping bags, tropical fruit, root veg, whole milk, yogurt, domestic, frozen, bread, pasta, sugar, cereals, coffee, soda, waffles, candy, berries, yogurt.

Then I have to read the read the grocery csv file, the grocery csv file is like that these are grocery csv file.

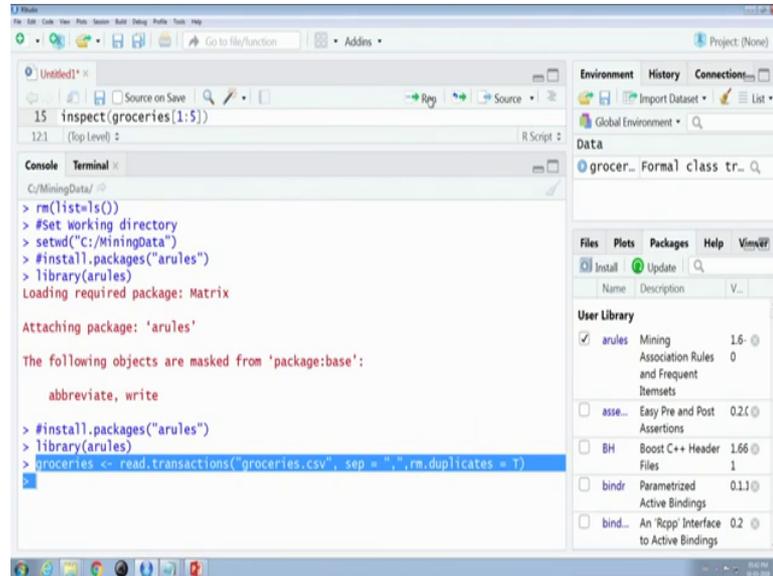
(Refer Slide Time: 22:53)

The screenshot shows an Excel spreadsheet with a single column of grocery items, identical to the one in the previous image. The items listed include: citrus fruit, semi finished bread, tropical food, yogurt, coffee, whole milk, cream cheese, pip fruit, cream cheese, whole milk, butter, rice, abrasive cleaner, milk/burn, soft milk, milk/burn, bottled in liquid (appetizer), jarred plants, whole milk, tropical fruit, tropical fruit, whole milk, cereal, yogurt, flour, bottled in dishes, beer, frankfurter, milk/burn, chicken, tropical fruit, butter, sugar, fruit/veg newspapers, packaged fruit/vegetables, packaged fruit/vegetables, chocolate, specialty bar, other vegetables, butter milk, tropical fruit, cream cheese, detergent newspapers, tropical fruit, root veg, other veg, frozen, milk/burn, flour, sweet, soft, waffles, candy, bathroom cleaner, bottled in canned beer, yogurt, beverage, milk/burn, soda, chocolate, other vegetables, frozen bread, soda, fruit/veg, canned, bottled in shipping bags, yogurt, beverage, bottled in specialty bar, hamburger, meat, other veg, milk/burn, bottled in hygiene & napkins, root vegetables, other veg, whole milk, beverage, sugar, jam, berries, other veg, whole milk, waffles, soft, soda, abrasive cleaner, beer, grapes, detergent, pasta, soda, packaged fruit/vegetables, chocolate, canned beer, root vegetables, other veg, whole milk, dessert, citrus fruit, packaged newspapers, beverage, milk/burn, soda, canned, bottled in shipping bags, tropical fruit, root veg, whole milk, yogurt, domestic, frozen, bread, pasta, sugar, cereals, coffee, soda, waffles, candy, berries, yogurt.

These are the items which are bought together these are transaction informations it contains the transactions of items this citrus fruit semi finished bread tropical food yogurt coffee are brought together in next transaction, only whole milk in next transaction pip fruit yogurt cream cheese these are the datasets my I have to find the item sets which are

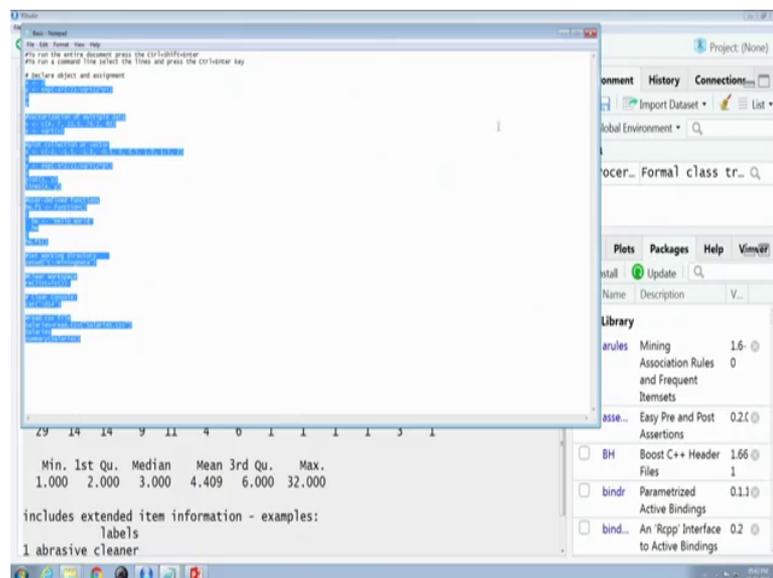
bought together next in most of time most frequently. Now I have to read the data set. Now dataset is readed the dataset is stored into the grocery.

(Refer Slide Time: 24:29)



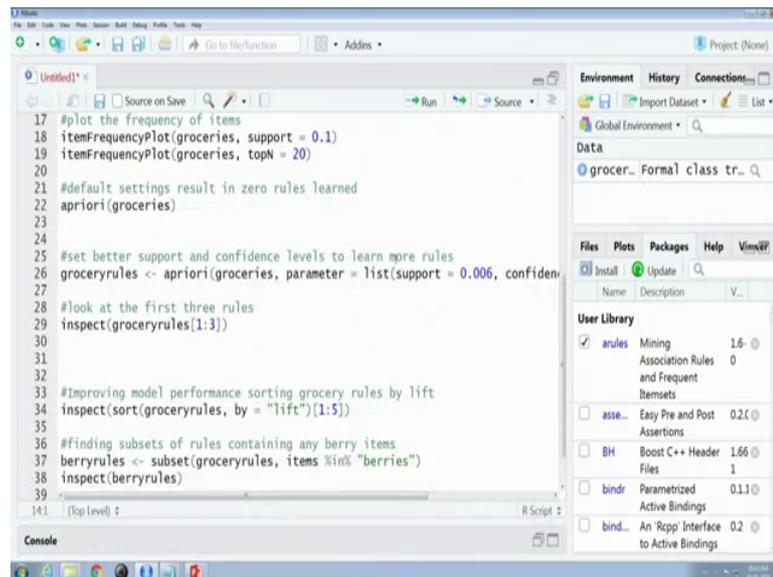
Then next it is showing the summary of the data set.

(Refer Slide Time: 24:53)



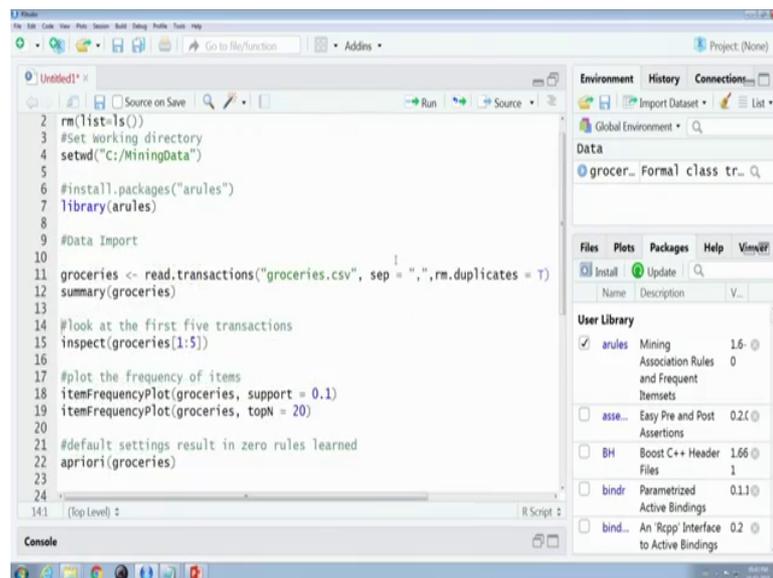
It contains 9835 rows and 169 columns because somebody in transaction, there are they have bought 169 items that is why it is showing here these are the most frequent items whole milk vegetables rolls and soda; these are more most frequently items and these are the size of distributions transactions.

(Refer Slide Time: 25:48)



```
17 #plot the frequency of items
18 itemFrequencyPlot(groceries, support = 0.1)
19 itemFrequencyPlot(groceries, topN = 20)
20
21 #default settings result in zero rules learned
22 apriori(groceries)
23
24
25 #set better support and confidence levels to learn more rules
26 groceryrules <- apriori(groceries, parameter = list(support = 0.006, confiden
27
28 #look at the first three rules
29 inspect(groceryrules[1:3])
30
31
32
33 #improving model performance sorting grocery rules by lift
34 inspect(sort(groceryrules, by = "lift")[1:5])
35
36 #finding subsets of rules containing any berry items
37 berryrules <- subset(groceryrules, items %in% "berries")
38 inspect(berryrules)
39
14.1 (Top Level) :
```

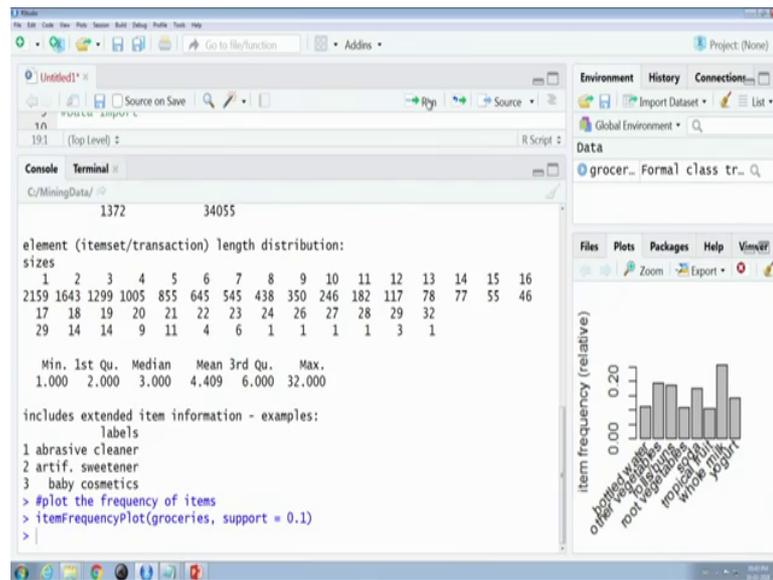
(Refer Slide Time: 25:49)



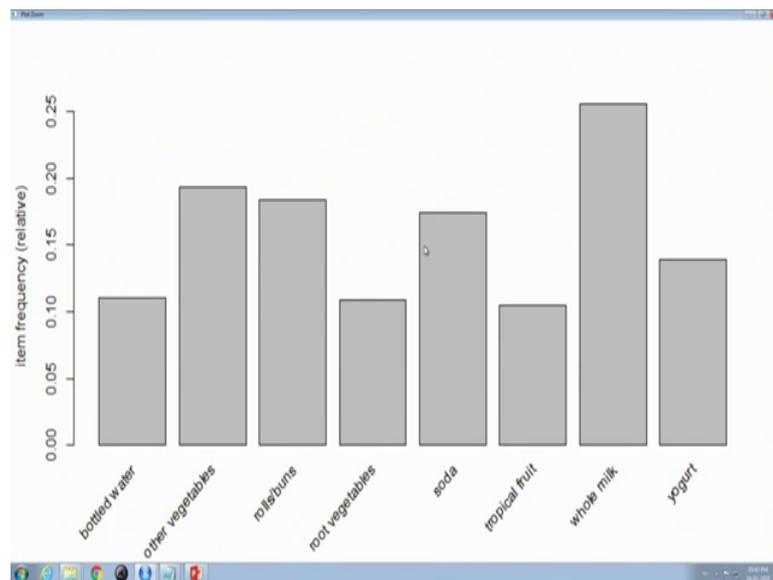
```
2 rm(list=ls())
3 #set working directory
4 setwd("C:/MiningData")
5
6 #install.packages("arules")
7 library(arules)
8
9 #Data Import
10
11 groceries <- read.transactions("groceries.csv", sep = ",", rm.duplicates = T)
12 summary(groceries)
13
14 #look at the first five transactions
15 inspect(groceries[1:5])
16
17 #plot the frequency of items
18 itemFrequencyPlot(groceries, support = 0.1)
19 itemFrequencyPlot(groceries, topN = 20)
20
21 #default settings result in zero rules learned
22 apriori(groceries)
23
24
14.1 (Top Level) :
```

Now, and to plot the frequency of items these are frequency of items.

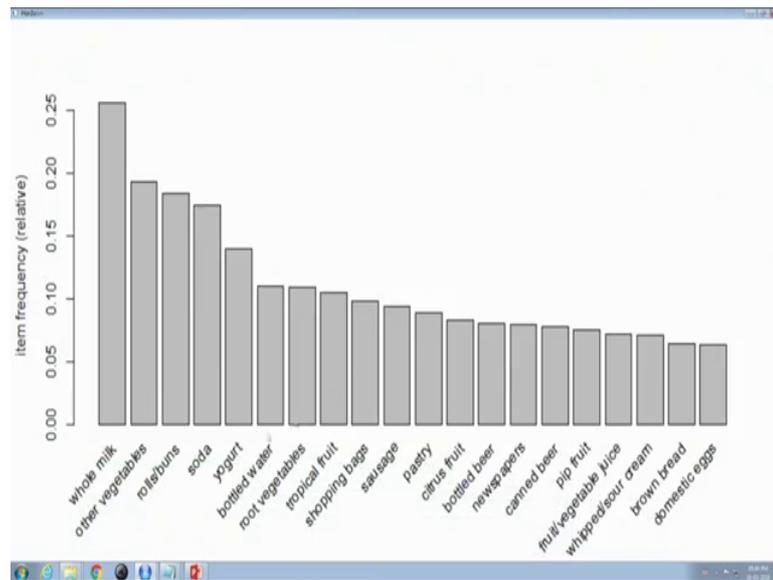
(Refer Slide Time: 26:01)



(Refer Slide Time: 26:05)



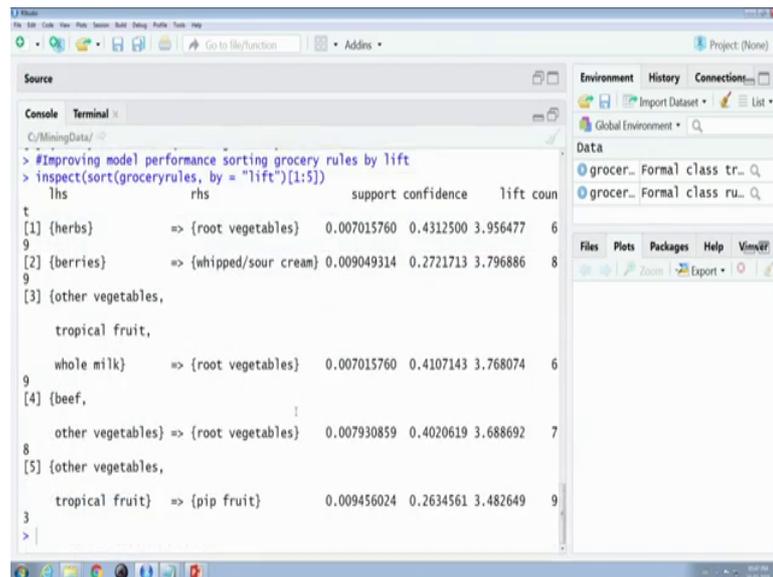
(Refer Slide Time: 27:20)



Item frequency bottled water other vegetables roll buns root vegetables these are the showing the frequency of items and to display frequency of plot this frequency of item top twenty items here is the top 20 items. Frequency of top 20 items, it is basically relative frequency in 0 to 1 scale, the most the whole milk is maximum frequency has a maximum frequency, the domestic is the minimum frequency. Now I apply the apriori. Now, apply now I have applied the apriori algorithm, it has confidence level initially the confidence is 0.8 minimum interval is 0.5 and now no rules has been generated in first iteration.

Now, we have to improve the grocery rules which support 0.6 and confidence 0.25 and minimum length is 2. Now the rules has been created, I have to inspect the rules. Now, it is showing the rules left hand side potted plants whole milk pasta depends on whole milk which support these are support. This is confidence and these are lift value, these are the rules. Now, this showing first 3 rules and to inspect by lift inspect the last last 5 rules by lift operation.

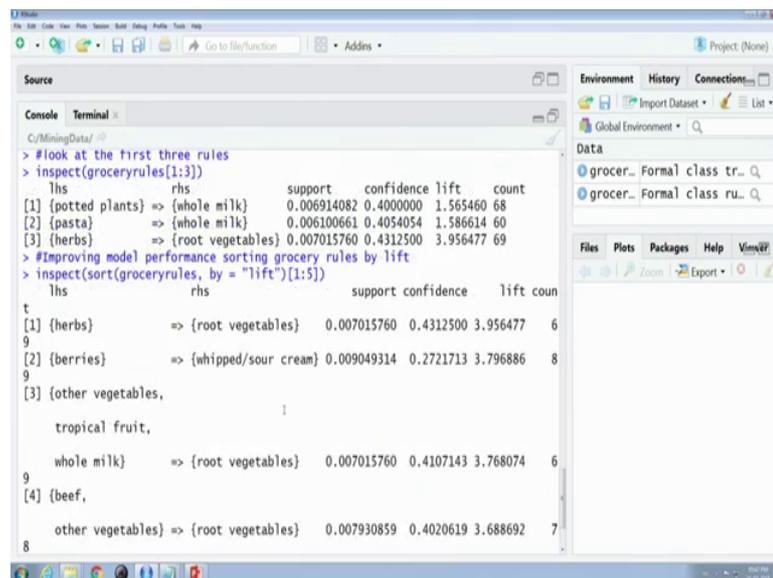
(Refer Slide Time: 29:50)



```
C:/MiningData/
> #Improving model performance sorting grocery rules by lift
> inspect(sort(groceryrules, by = "lift")[1:5])
```

	lhs	rhs	support	confidence	lift	count
[1]	{herbs}	=> {root vegetables}	0.007015760	0.4312500	3.956477	6
[2]	{berries}	=> {whipped/sour cream}	0.009049314	0.2721713	3.796886	8
[3]	{other vegetables, tropical fruit, whole milk}	=> {root vegetables}	0.007015760	0.4107143	3.768074	6
[4]	{beef, other vegetables}	=> {root vegetables}	0.007930859	0.4020619	3.688692	7
[5]	{other vegetables, tropical fruit}	=> {pip fruit}	0.009456024	0.2634561	3.482649	9

(Refer Slide Time: 29:51)



```
C:/MiningData/
> #look at the first three rules
> inspect(groceryrules[1:3])
```

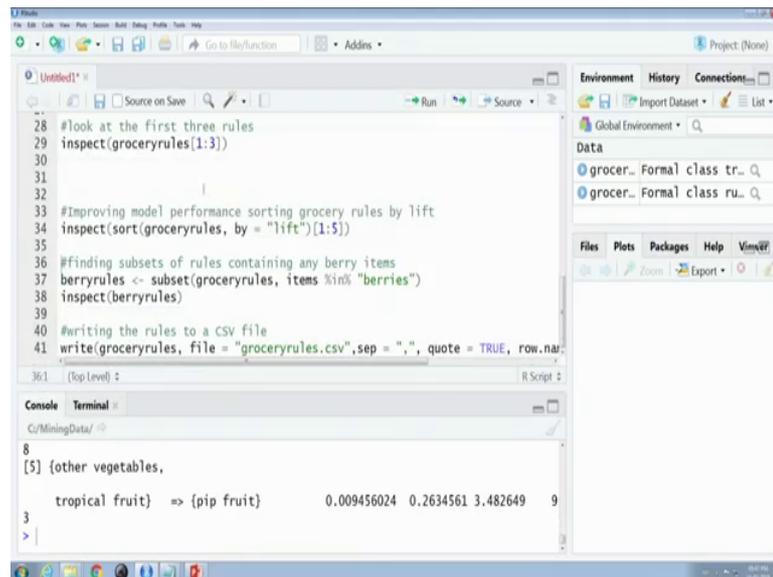
	lhs	rhs	support	confidence	lift	count
[1]	{potted plants}	=> {whole milk}	0.006914082	0.4000000	1.565460	68
[2]	{pasta}	=> {whole milk}	0.006100661	0.4054054	1.586614	60
[3]	{herbs}	=> {root vegetables}	0.007015760	0.4312500	3.956477	69

```
> #Improving model performance sorting grocery rules by lift
> inspect(sort(groceryrules, by = "lift")[1:5])
```

	lhs	rhs	support	confidence	lift	count
[1]	{herbs}	=> {root vegetables}	0.007015760	0.4312500	3.956477	6
[2]	{berries}	=> {whipped/sour cream}	0.009049314	0.2721713	3.796886	8
[3]	{other vegetables, tropical fruit, whole milk}	=> {root vegetables}	0.007015760	0.4107143	3.768074	6
[4]	{beef, other vegetables}	=> {root vegetables}	0.007930859	0.4020619	3.688692	7

Now, it is showing herbs; if it is herbs, then root vegetables will be in right hand side, it is berries whipped and sour cream and tropical other vegetables tropical fruit whole milk root vegetables and so on.

(Refer Slide Time: 30:18)

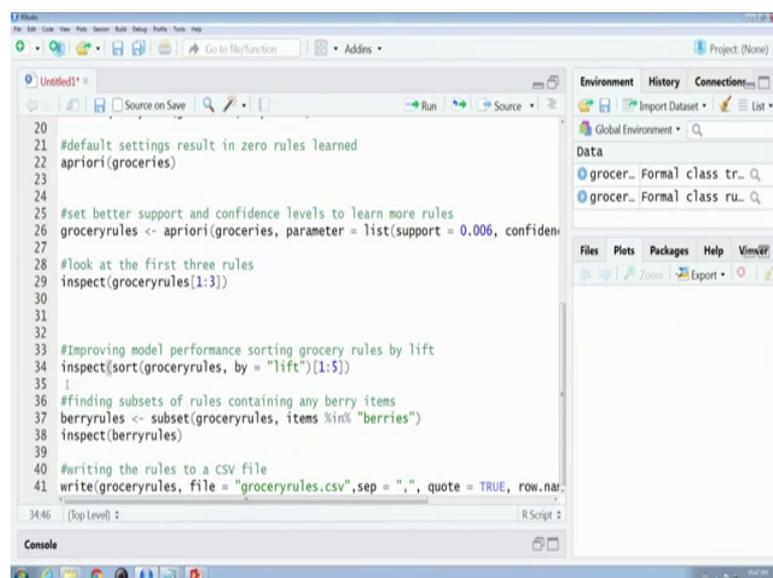


```
28 #look at the first three rules
29 inspect(groceryrules[1:3])
30
31
32
33 #Improving model performance sorting grocery rules by lift
34 inspect(sort(groceryrules, by = "lift")[1:5])
35
36 #finding subsets of rules containing any berry items
37 berryrules <- subset(groceryrules, items %in% "berries")
38 inspect(berryrules)
39
40 #writing the rules to a CSV file
41 write(groceryrules, file = "groceryrules.csv", sep = ",", quote = TRUE, row.nam
```

Console Terminal

```
C:/MiningData/ >>
8
[5] {other vegetables,
      tropical fruit} => {pip fruit}      0.009456024  0.2634561  3.482649  9
3
>
```

(Refer Slide Time: 30:25)

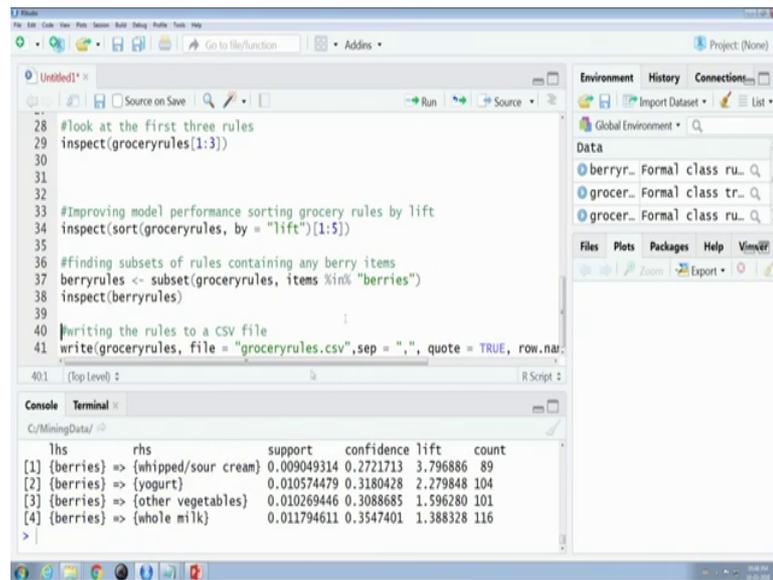


```
20
21 #default settings result in zero rules learned
22 apriori(groceries)
23
24
25 #set better support and confidence levels to learn more rules
26 groceryrules <- apriori(groceries, parameter = list(support = 0.006, confiden
27
28 #look at the first three rules
29 inspect(groceryrules[1:3])
30
31
32
33 #Improving model performance sorting grocery rules by lift
34 inspect(sort(groceryrules, by = "lift")[1:5])
35
36 #finding subsets of rules containing any berry items
37 berryrules <- subset(groceryrules, items %in% "berries")
38 inspect(berryrules)
39
40 #writing the rules to a CSV file
41 write(groceryrules, file = "groceryrules.csv", sep = ",", quote = TRUE, row.nam
```

Console

Now, I want to find a subset of rules containing any various items and how for that I have to use the subset function these is the it contains the rules and item is varies. Now I just execute this command after executing this command, I have to inspect this command.

(Refer Slide Time: 30:50)

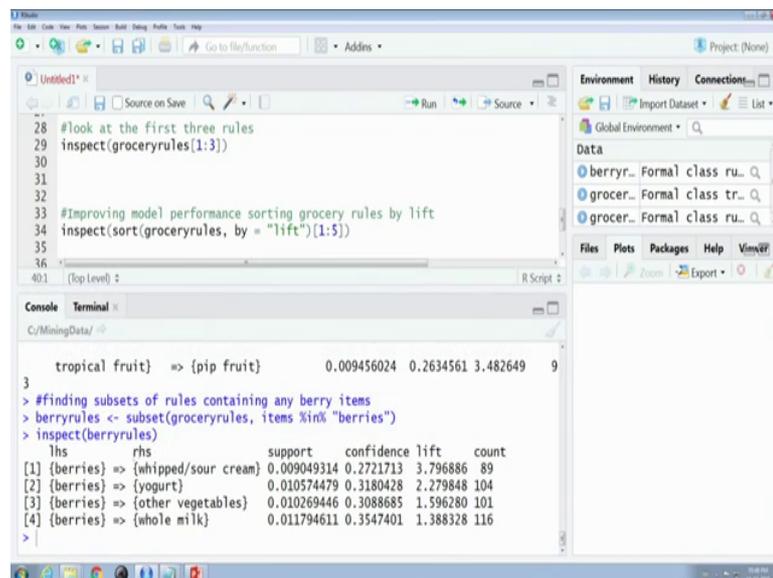


```
28 #look at the first three rules
29 inspect(groceryrules[1:3])
30
31
32
33 #Improving model performance sorting grocery rules by lift
34 inspect(sort(groceryrules, by = "lift")[1:5])
35
36 #finding subsets of rules containing any berry items
37 berryrules <- subset(groceryrules, items %in% "berries")
38 inspect(berryrules)
39
40 #writing the rules to a csv file
41 write(groceryrules, file = "groceryrules.csv", sep = ",", quote = TRUE, row.names = FALSE)
```

lhs	rhs	support	confidence	lift	count
[1] {berries}	=> {whipped/sour cream}	0.009049314	0.2721713	3.796886	89
[2] {berries}	=> {yogurt}	0.010574479	0.3180428	2.279848	104
[3] {berries}	=> {other vegetables}	0.010269446	0.3088685	1.596280	101
[4] {berries}	=> {whole milk}	0.011794611	0.3547401	1.388328	116

Now it is showing rules with berries.

(Refer Slide Time: 30:53)



```
28 #look at the first three rules
29 inspect(groceryrules[1:3])
30
31
32
33 #Improving model performance sorting grocery rules by lift
34 inspect(sort(groceryrules, by = "lift")[1:5])
35
36 #finding subsets of rules containing any berry items
37 berryrules <- subset(groceryrules, items %in% "berries")
38 inspect(berryrules)
39
40 #writing the rules to a csv file
41 write(groceryrules, file = "groceryrules.csv", sep = ",", quote = TRUE, row.names = FALSE)
```

lhs	rhs	support	confidence	lift	count
[1] {berries}	=> {whipped/sour cream}	0.009049314	0.2721713	3.796886	89
[2] {berries}	=> {yogurt}	0.010574479	0.3180428	2.279848	104
[3] {berries}	=> {other vegetables}	0.010269446	0.3088685	1.596280	101
[4] {berries}	=> {whole milk}	0.011794611	0.3547401	1.388328	116

Berries whipped and sour cream berries and yogurt with confidence level these berries and other vegetables and berries and whole milk, these are the most frequently bought items with berries. Now I have to write the all the rules into a particular file for that I have to use right command, right the grocery rules into the file grocery rules csv. Now I am now trying to execute this. Now the file has been created. Now I have to check the file this is the rule files here inter rule files has been saved.

Now, this is the procedure of how to generate association rules and save into the csv file the library is arules.

(Refer Slide Time: 32:53)

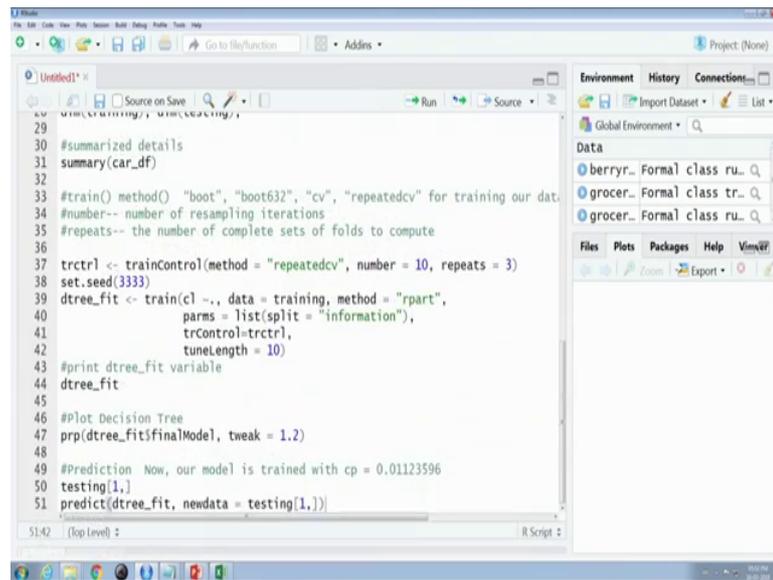
The slide is titled "Decision Trees" and contains the following content:

- The decision tree induction is learning of decision tree from class label training data
- The Cars Evaluation data set
 - #V1 buying price v-high, high, med, low
 - #V2 price of the maintenance v-high, high, med, low
 - #V3 doors 2, 3, 4, 5-more
 - #V4 persons 2, 4, more
 - #V5 luggage boot small, med, big
 - #V6 safety low, med, high
 - #V7 class unacc, acc, good, v-good
- Packages
 - caret
 - rpart.plot
 - Rpart
 - e1071

At the bottom of the slide, there is a navigation bar with icons and the text "Data Mining With R".

Now, I am going to show; how to generate the decision tree for creating decision tree we have to use 4 packages; one is caret another is Rpart dot plot and for decision tree generation e1071 and Rpart e1071 is the most important package for creating decision trees in I am using one data set; this name is car evaluation data sets, it contains 5 7 attributes price of the car maintenance value of the car doors of the car number of persons can sit into the car the booted the luggage area of the car the safety matters of the car and class; the class is either unsatisfactory, it is satisfactory good or very good these are the class now you want to make decision tree based on this information.

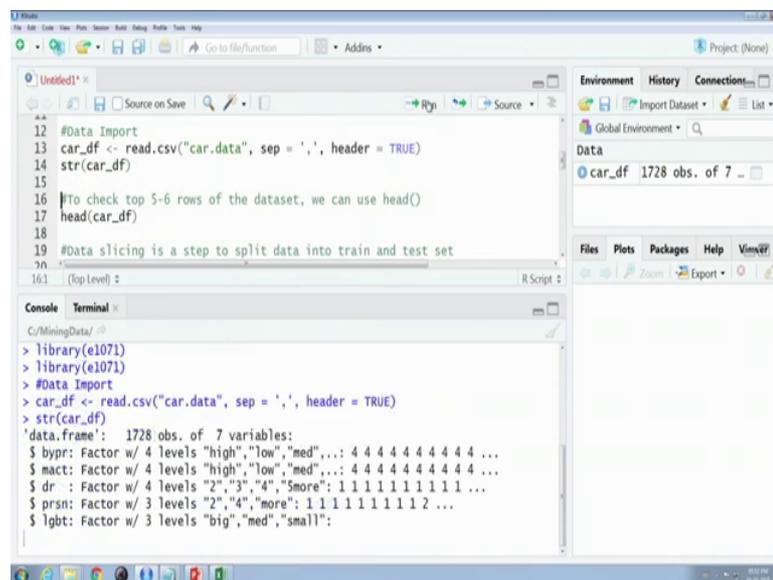
(Refer Slide Time: 34:42)



```
29 #summary
30 #summarized details
31 summary(car_df)
32
33 #train() method() "boot", "boot632", "cv", "repeatedcv" for training our data
34 #number-- number of resampling iterations
35 #repeats-- the number of complete sets of folds to compute
36
37 trctrl <- trainControl(method = "repeatedcv", number = 10, repeats = 3)
38 set.seed(3333)
39 dtree_fit <- train(cl ~., data = training, method = "rpart",
40                   parms = list(split = "information"),
41                   trControl = trctrl,
42                   tuneLength = 10)
43 #print dtree_fit variable
44 dtree_fit
45
46 #Plot Decision Tree
47 prp(dtree_fit$finalModel, tweak = 1.2)
48
49 #Prediction Now, our model is trained with cp = 0.01123596
50 testing[1,]
51 predict(dtree_fit, newdata = testing[1,])
```

Now, I want to add the library caret Rpart Rplot and e701 e1071. Now, I have to read the data from card data, then I have to view the structure of the data set.

(Refer Slide Time: 35:26)



```
12 #Data Import
13 car_df <- read.csv("car.data", sep = ',', header = TRUE)
14 str(car_df)
15
16 #To check top 5-6 rows of the dataset, we can use head()
17 head(car_df)
18
19 #Data slicing is a step to split data into train and test set
20
21
```

```
C:/MiningData/ >>
> library(e1071)
> library(e1071)
> #Data Import
> car_df <- read.csv("car.data", sep = ',', header = TRUE)
> str(car_df)
'data.frame': 1728 obs. of 7 variables:
 $ bypr: Factor w/ 4 levels "high","low","med",...: 4 4 4 4 4 4 4 4 4 4 ...
 $ mact: Factor w/ 4 levels "high","low","med",...: 4 4 4 4 4 4 4 4 4 4 ...
 $ dr : Factor w/ 4 levels "2","3","4","5more": 1 1 1 1 1 1 1 1 1 1 ...
 $ prsn: Factor w/ 3 levels "2","4","more": 1 1 1 1 1 1 1 1 1 2 ...
 $ lgbt: Factor w/ 3 levels "big","med","small":
```

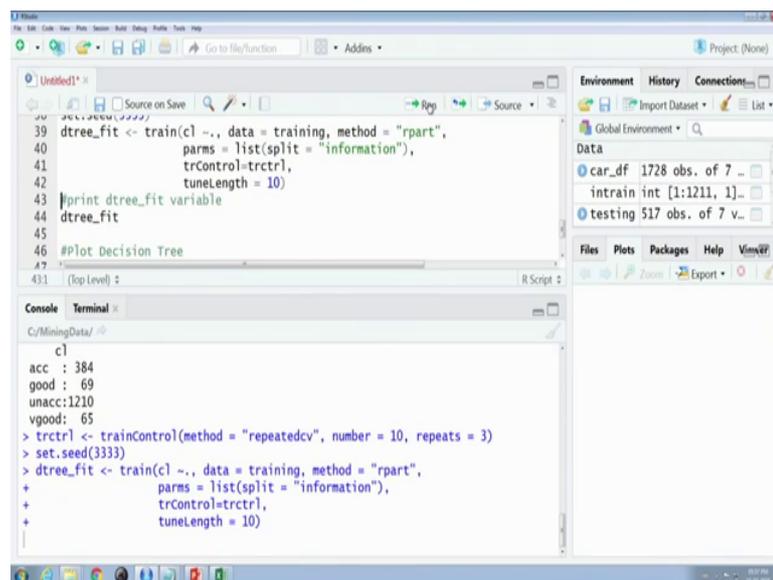
It is showing the structure of the data set it has 4 levels bypr means buying price mact means the maintaining cost, it has 4 levels door number of doors peers number of persons here these are the structure of the data set, then I am to display, the first few data set and it is displaying the data set 6 data sets, then I have to divides the data set into 2 parts 30 percent for testing the decision tree and 70 percent for training purpose that is

why I am creating a party one partition based on class column; cl is class based on class column and 0.7 means seventy percent of data will be used for training purpose list is false.

Now, at first, I have to make a seed on random value, then I have generated a training data set. Now, training data set has been created. Now test data set; these are the test data set. Now I have to check the dimension of the training data set. Now, it is show showing the dimension of the training data set and testing data set the training data set contains 1211 records and test data set has 5 517 records with 7 attributes.

Here is the summary of the car data sets. Now, I have to apply the training with number of 10 folds number of 10 folds number of resampling iterations, these are number and repeats means the number of complete sets of folds to compute I am here using repeated cv method for training, there are lots of method available boot method boot six 32 method cv method repeat cv method for training our data. Now execute training make one random variable, then then I have to create the tree these function training these are training function and it will create store the tree into the tree feet this is a class level data is training data I am using Rpart method is the transaction control tune length is 10; 10 is the number of resampling iterations.

(Refer Slide Time: 39:52)



```
39 dtree_fit <- train(c1 ~., data = training, method = "rpart",
40                    parms = list(split = "information"),
41                    trControl=trctrl,
42                    tuneLength = 10)
43 #print dtree_fit variable
44 dtree_fit
45
46 #Plot Decision Tree
47
481 (Top Level) :
```

```
C:/MiningData/
c1
acc : 384
good : 69
unacc:1210
vgood: 65
> trctrl <- trainControl(method = "repeatedcv", number = 10, repeats = 3)
> set.seed(3333)
> dtree_fit <- train(c1 ~., data = training, method = "rpart",
+                   parms = list(split = "information"),
+                   trControl=trctrl,
+                   tuneLength = 10)
```

The screenshot shows the R Studio environment. The script editor contains R code for training a decision tree using the rpart method with repeated cross-validation. The console shows the output of the training process, including the class distribution of the training data and the execution of the training function.

Now, you have to execute the train data to create the tree. Now, tree is created; now I have to get that tree information; these are the data of a this is the complexity parameters

this is the accuracy and kappa value and it is showing it is using cart algorithm, there are 6 predictors and 4 classes. Now I want to plot the tree plot the decision tree; these are decision tree, but it will be due to graphics problem in this machine, it will not displaying properly. Now we have to test one dataset. Suppose, these one testing data set the testing dataset showing buying price is very high maintenance cost is very high and class is unsatisfactory these my testing data whether these data is correctly predicted or not now.

(Refer Slide Time: 40:41)

```

44 dtree_fit
45
46 #Plot Decision Tree
47 prp(dtree_fit$finalModel, tweak = 1.2)
48
49 #Prediction Now, our model is trained with cp = 0.01123596
50 testing[1,]
51 predict(dtree_fit, newdata = testing[1,])

```

Console Terminal

```

C:/MiningData/ >
0.01648352 0.8065620 0.5775304
0.01831502 0.7979943 0.5577698
0.02060440 0.7941443 0.5489378
0.02197802 0.7922137 0.5480683
0.06668132 0.7883246 0.5726618
0.09340659 0.7233582 0.2223118

Accuracy was used to select the optimal model using the largest value.
The final value used for the model was cp = 0.01098901.
> #Plot Decision Tree
> prp(dtree_fit$finalModel, tweak = 1.2)
>

```

I am now I am predicting the same data here new data set predict dtree the data record is this now; now it is predicting correctly because this unsatisfactory and in my dataset it is also unsatisfactory here these both are same.

(Refer Slide Time: 42:31)

K Means Clustering

- The k-means algorithm generates the clusters where centre of each cluster denotes the mean value of the respective objects
- Input – k number of clusters and data set
- Output- set of k clusters
- The Whole Sales Customers data set
 - Channel
 - Region
 - Fresh
 - Milk,
 - Grocery
 - Frozen

Data Mining With R

Now, I want to show how to apply R program to use k means clustering the k means clustering generates cluster where the centre of each cluster denotes mean value of respective objects the input set of k means algorithm the number of clusters you want to generate and the data set and it will create k k clusters in my data set, I am using 1 wholesale customers data set in that data set, there are lots of attributes; I am I am trying to the attributes are channel regions fresh milk grocery I will generate the clusters on grocery attribute.

(Refer Slide Time: 43:59)

```
1 cat("\014")
2 rm(list=ls())
3 setwd("c:/MiningData")
4
5 #Data Import
6 x<- read.csv("wholeSaleCustomersData.csv",header=TRUE)
7
8 #kmeans cluster
9 km <- kmeans(x[,grocery],4,10)
10 print(km)
11 plot(x[,grocery, col = km$cluster)
```

Environment History Connections

Global Environment

Data

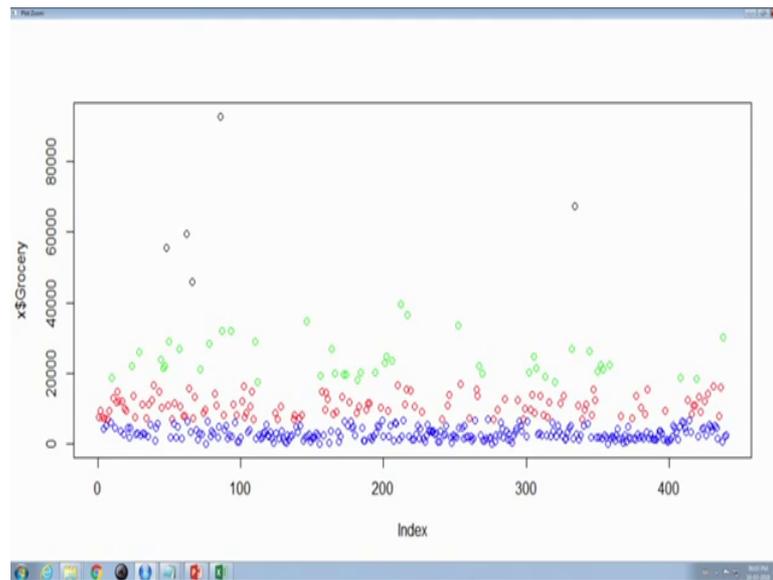
- car_df 1728 obs. of 7
- dtree... Large train (24
- ewdata 1 obs. of 7 var...

Files Plots Packages Help View

Zoom Export

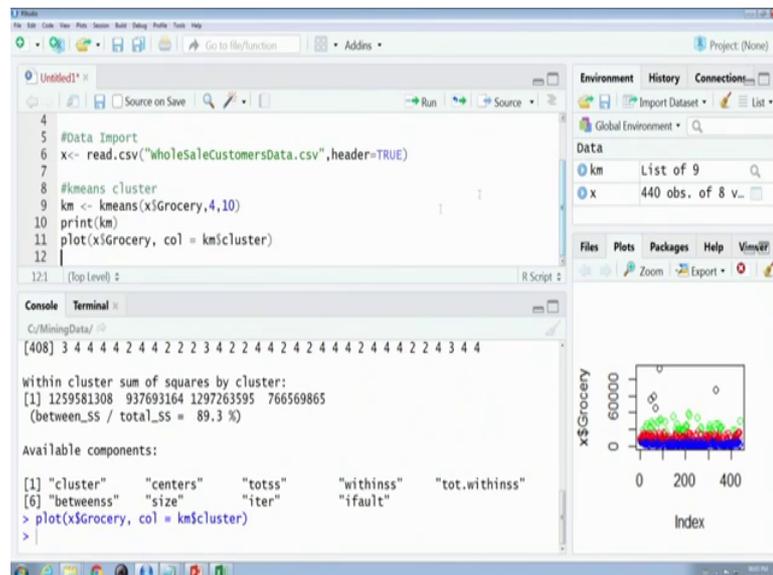
1:12 (Top Level) R Script

(Refer Slide Time: 45:54)



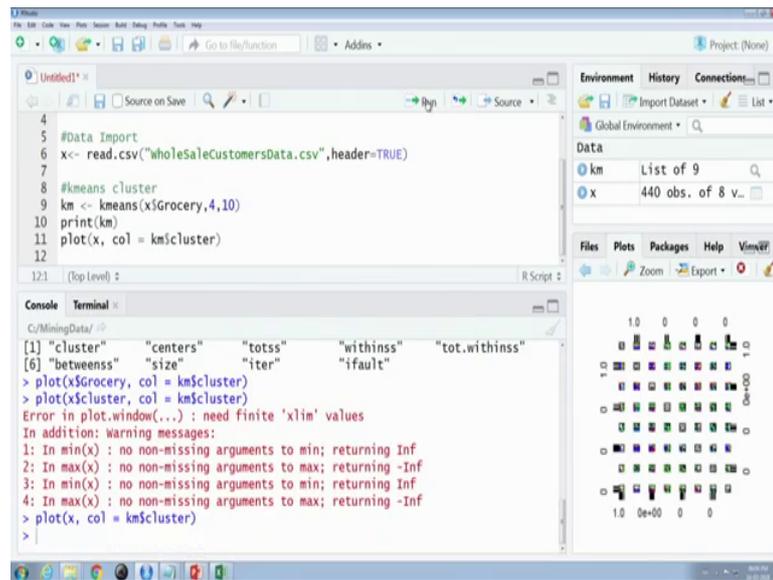
Now, we want to generate the cluster on grocery. Now it is showing the clusters on grocery; there are 4 clusters this blue, one is one cluster, red one is another cluster and green one is another cluster and very small; these are very small cluster, it scattered basically.

(Refer Slide Time: 46:17)

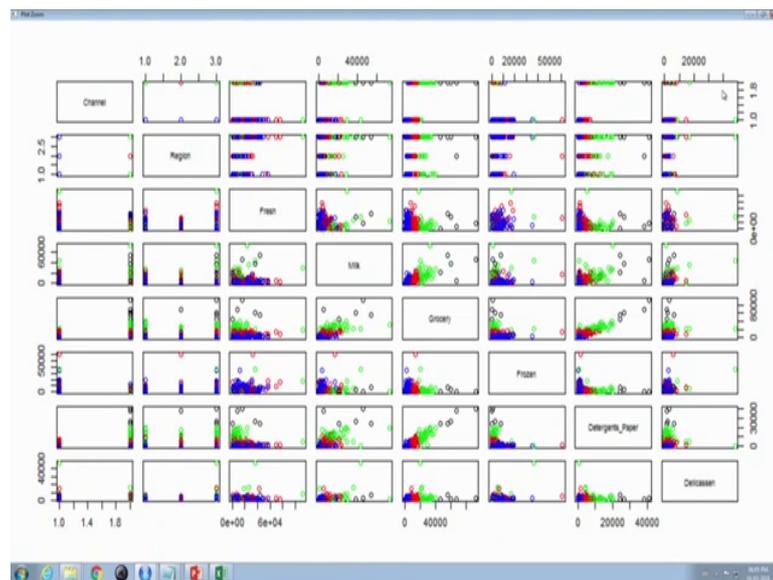


Suppose, I want to display inter cluster all the clusters of all attributes now it is showing all the cloud clusters of all attributes here these are clusters of channel attributes these are clusters of where is the grocery is the clusters of grocery attributes.

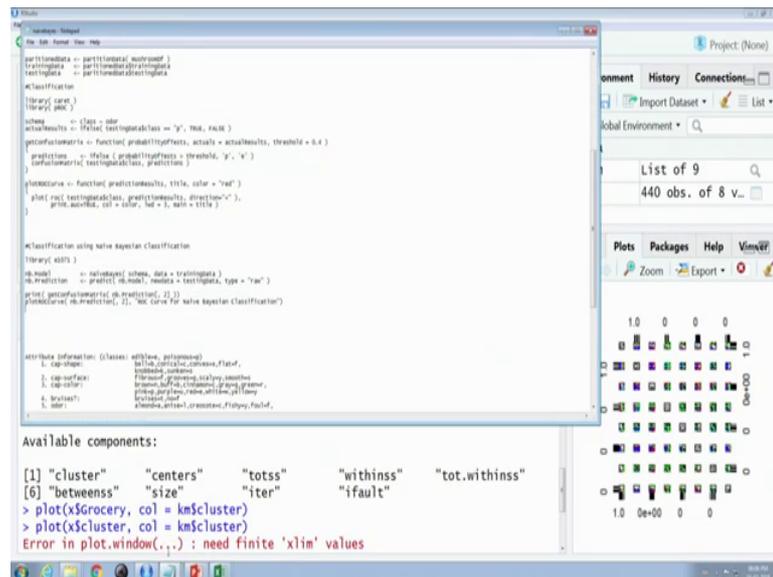
(Refer Slide Time: 47:27)



(Refer Slide Time: 47:31)



(Refer Slide Time: 49:17)



```
partitionsdata <- partition(mushroom)
trainingsdata <- partition(partitionsdata)
testingsdata <- partition(partitionsdata)

#NaiveBayes
library(car)
library(glm)

scheme <- list = odds
actualResults <- predict(testingsdata[1:5], model, FUN = )
getConfusionMatrix <- function(probabilityPress, actual = actualResults, threshold = 0.4)
{
  predictions <- ifelse(probabilityPress > threshold, "Y", "N")
  getConfusionMatrix(testingsdata$class, predictions)
}

plotAccuracy <- function(getConfusionMatrix, title, color = "red")
{
  plot(mtcars, las=1, col=getConfusionMatrix, direction="v")
  print(mtcars, col = color, las = 1, main = title)
}

#NaiveBayes using naive Bayes Classification
library(e1071)
library(ada)

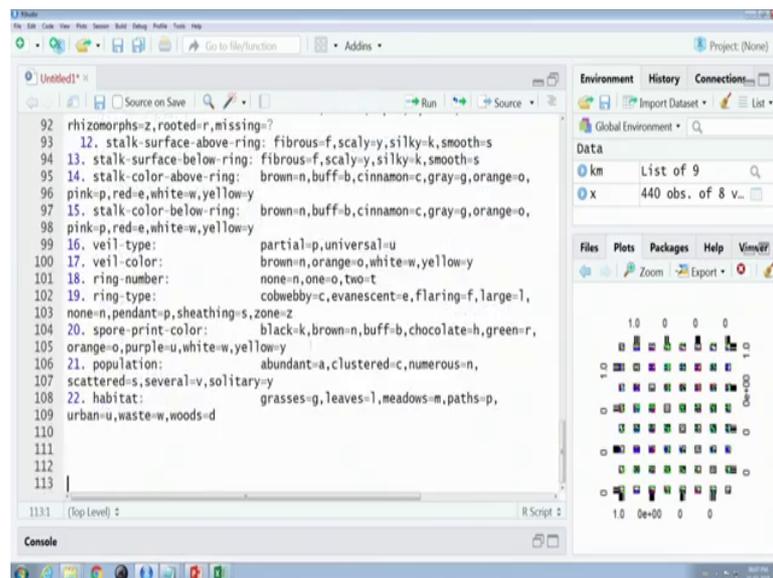
#m.model <- naiveBayes(scheme, data = trainingsdata)
#m.Prediction <- predict(m.model, testingsdata, type = "raw")
print(getConfusionMatrix(m.Prediction, 2)) #m.Curve for naive Bayes Classification

attributes(Information) #classes = e1071, partition
1. cap-shape: brown-n,conical-c,convex-f,flat-f,
convex-g,knobbed-k,
2. cap-surface: fibrous-f,smooth-s,silky-s,smooth-s
3. cap-color: brown-n,black-b,cinnamon-c,gray-g,orange-o,
pink-p,red-e,white-w,yellow-y
4. bracket: none-n,one-o,two-t
5. odor: anise-a,atlas-t,creosote-c,fishy-f,leaves-l,
```

Available components:

```
[1] "cluster" "centers" "totss" "withins" "tot.withins"
[6] "betweens" "size" "iter" "ifault"
> plot(x$Grocery, col = km$cluster)
> plot(x$cluster, col = km$cluster)
Error in plot.window(...): need finite 'xlim' values
```

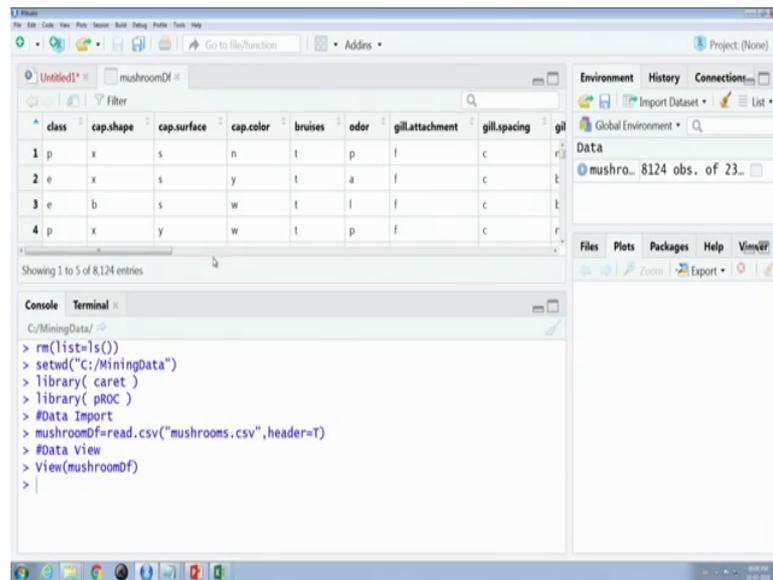
(Refer Slide Time: 49:32)



```
92 rhizomorphs=z,rooted-r,missing=?
93 12. stalk-surface-above-ring: fibrous=f,scaly=y,silky-k,smooth-s
94 13. stalk-surface-below-ring: fibrous=f,scaly=y,silky-k,smooth-s
95 14. stalk-color-above-ring: brown-n,buff-b,cinnamon-c,gray-g,orange-o,
96 pink-p,red-e,white-w,yellow-y
97 15. stalk-color-below-ring: brown-n,buff-b,cinnamon-c,gray-g,orange-o,
98 pink-p,red-e,white-w,yellow-y
99 16. veil-type: partial-p,universal-u
100 17. veil-color: brown-n,orange-o,white-w,yellow-y
101 18. ring-number: none-n,one-o,two-t
102 19. ring-type: cobwebby-c,evanescent-e,flaring-f,large-l,
103 none-n,pendant-p,sheathing-s,zone-z
104 20. spore-print-color: black-k,brown-n,buff-b,chocolate-h,green-r,
105 orange-o,purple-u,white-w,yellow-y
106 21. population: abundant-a,clustered-c,numerous-n,
107 scattered-s,several=v,solitary=y
108 22. habitat: grasses=g,leaves=l,meadows=m,paths=p,
109 urban-u,waste-w,woods=d
110
111
112
113
```

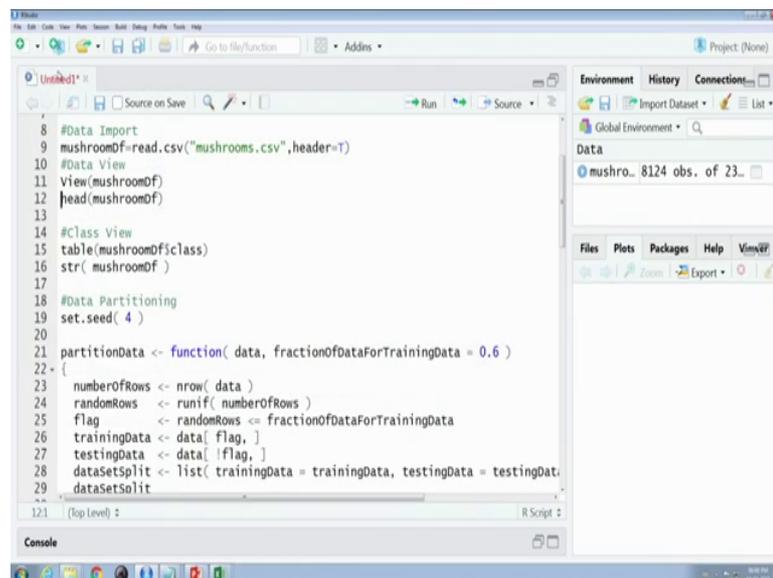
It contains lots of attributes cap shape cap surface cap colour, these are the abbreviations of the table, these are the abbreviations of the tables and two classes eatable or poisonous to execute Naive Bayesian algorithm, we need to include two libraries one is Carret and the other one is proc library, then I have to import data same way.

(Refer Slide Time: 50:45)



Similarly, then want to view the mushroom datasets this is the same data.

(Refer Slide Time: 50:59)



These are the first few first 6 data sets, it is very big table that is why it is displaying this way first six in second second set of attributes these are third set of attributes these are the fourth set of attributes this because there are lots of attributes that is why it is displaying this way.

Now, I want to view the class table; there are 2 classes edible and poisonous; there are 4000; around 4000 classes mushrooms are edible around 3000 classes are poisonous mushrooms.

(Refer Slide Time: 52:05)

```

11 View(mushroomDF)
12 head(mushroomDF)
13
14 #Class View
15 table(mushroomDF$class)
16 str(mushroomDF)
17
18 #Data Partitioning
19 set.seed(4)
20
21 partitionData <- function( data, fractionOfDataForTrainingData = 0.6 )
22
181 (Top Level)

```

The console output shows the structure of the data set:

```

C:/MiningData/ >
$ ring.type      : Factor w/ 5 levels "e","t","l","n",...: 5 5 5 1 5 5 5
$ ...
$ spore.print.color : Factor w/ 9 levels "b","h","k","n",...: 3 4 4 3 4 3 3 4
$ ...
$ population      : Factor w/ 6 levels "a","c","n","s",...: 4 3 3 4 1 3 3 4
$ ...
$ habitat         : Factor w/ 7 levels "d","g","l","m",...: 6 2 4 6 2 2 4 4
$ ...

```

what is the structure of the; these are these is the structure of the data sets there are twenty three variables in the data sets and around eight thousand observations, then for again have to partition the data for training set and test set.

(Refer Slide Time: 52:33)

```

24 nRows <- nrow(mushroomDF)
25 flag <- randomRows <- fractionOfDataForTrainingData
26 trainingData <- data[flag, ]
27 testingData <- data[!flag, ]
28 dataSetSplit <- list( trainingData = trainingData, testingData = testingData )
29 dataSetSplit
30 }
31
32 partitionedData <- partitionData( mushroomDF )
33 trainingData <- partitionedData$trainingData
34 testingData <- partitionedData$testingData
35
36
371 (Top Level)

```

The console output shows the partitioning process:

```

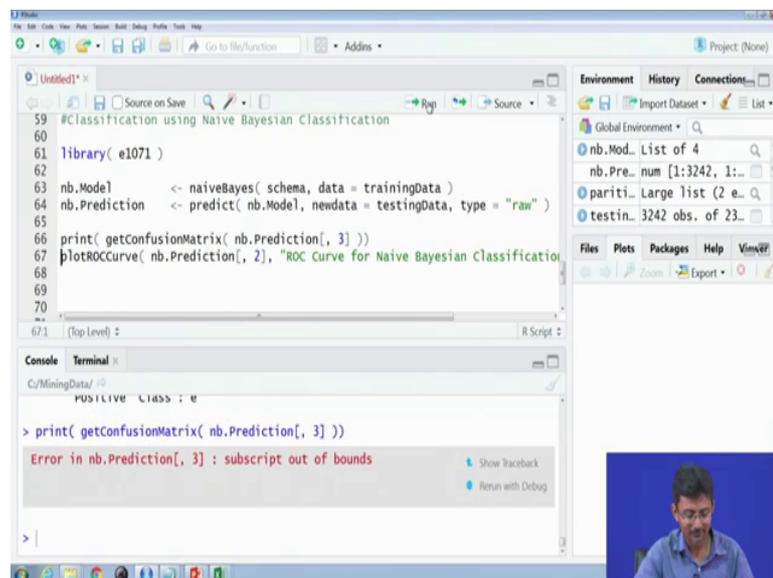
C:/MiningData/ >
+ nRows <- nrow( data )
+ randomRows <- runif( numberOfRows )
+ flag <- randomRows <- fractionOfDataForTrainingData
+ trainingData <- data[ flag, ]
+ testingData <- data[ !flag, ]
+ dataSetSplit <- list( trainingData = trainingData, testingData = testingData )
+ dataSetSplit
+ }
>

```

Now training set is being generated now test set is generated now to make a schema class and odor class versus odor schema. Now testing the data set on class equal to p actual results now have to develop the confusion matrix with the probability test with the function probability test. Now confusion matrix has been generated. Now I have to make prediction curve. Now, execute the naive Bayesian algorithm on training data and it will be stored in nb dot model and this will be the predicted data.

Now, prediction on any model testing data type is raw now and to print the confusion matrix the prediction two prediction data two see for prediction data 2, this is the confusion matrix references is e to e is and with accuracy 0.9851; the class is edible here.

(Refer Slide Time: 55:33)



The screenshot shows the R Studio environment. The script editor contains the following R code:

```
59 #Classification using Naive Bayesian Classification
60
61 library( e1071 )
62
63 nb.Model <- naiveBayes( schema, data = trainingData )
64 nb.Prediction <- predict( nb.Model, newdata = testingData, type = "raw" )
65
66 print( getConfusionMatrix( nb.Prediction[, 3] ) )
67 plotROCcurve( nb.Prediction[, 2], "ROC Curve for Naive Bayesian Classification" )
68
69
70
71 (Top Level) :
```

The Environment pane on the right shows the following objects:

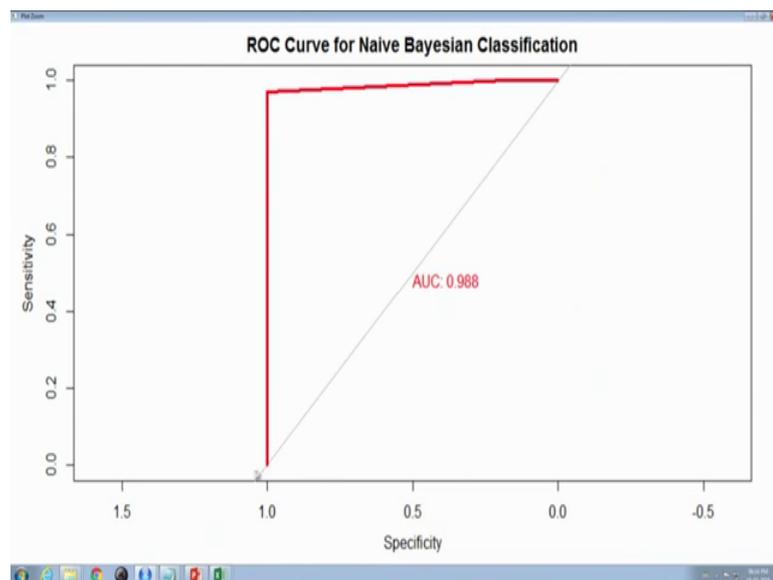
- nb.Mod... List of 4
- nb.Pre... num [1:3242, 1:...
- pariti... Large list (2 e... Q
- testin... 3242 obs. of 23...

The Console pane shows the following output and error:

```
> print( getConfusionMatrix( nb.Prediction[, 3] ) )
Error in nb.Prediction[, 3] : subscript out of bounds
```

Suppose, I want to check another data this is the first data no there is no; now I want to plot the ROC curve or Bayesian matrix prediction matrix. Now this is the curve these ROC curve for Bayesian classification is a specificity and sensitivity curve.

(Refer Slide Time: 56:32)



Therefore, for executing the Naive Bayesian algorithm at first, we have to partition the data into training set and test set then make a schema in this class with an particular attribute then the actual results based on poisonous here then generate the confusion matrix after that execute the Naive Bayesian algorithm these are procedure ok.

Thank you.