

Date
2/10/2023

Dataware house and Data mining

CSA1675

M. Venkata Hari
19/10/2014-64
CSE

Assignment-1

Data Manipulation

Round(x,n):-

round the values of x to n decimal places

ceiling(x):-

vector x of smallest integer

floor(x):-

vector x of largest integer

as.integer:-

Truncates real x to integer (compared to round(x,n)).

Statistics

min() → lowest value from given data

mean() → Average value

median() → middle value Q_1, Q_2, Q_3

Sum() → Total

var() → # produces the variance covariance matrix

sd() → # standard deviation

writing into a csv file

R can create csv file from existing data-frame. The write.csv() function is used to create the csv file. This file gets created in the working directory.

```
# create a data frame  
data <- read.csv("input.csv")
```

```
retval <- subset(data, as.  
date(start.date)  
> as.date(2004-10-01))
```

write filtered data into a new file.

```
write.csv(retval, "output.csv")
```

```
newdata <- read.csv("output.csv")  
print(newdata)
```

Read:-

Creating a file:-

using file.create() function, a new file can be created from console or transcripts if already exist.

Syntax:-

```
file.create(" ")
```

Example:-

```
# create a file  
# the file created can be seen  
# in your working directory.
```

Reading a file:-

using read.table() function R, * files can be read and output is shown as data frame.

Syntax:-

```
read.table(file)
```

Example:-

```
# Reading txt file  
new.irisk <- read.table(file =  
"GFG.txt")  
  
# print  
print(new.irisk)
```

Transformation:-

five num:-

Turkey five numbers min, lower hinge, hinge, median, upper hinge max.

table() →

frequency counts of entries. Ideally the entries are factors (although it works with integers (or) even reals).

scale(data, scale = 1)

→ # centres around the mean and scales by the sd

Input and display:-

read.table(filename, header = TRUE)

read files with tabs in first row

read a table (or) space delimited file.

read.table(filename, header = TRUE, sep = ".")

read csv files

x = c(1:10) → # create a data vector with elements 1-10

vect = c(x,y) → # combine them into vector or length 2n.

mat = (bind(x,y)) → # combine them into an x2 matrix

② Suppose that the data for the analysis includes the attribute age. the age values for the data tuples are 18, 15, 16, 16, 19, 20, 20, 21, 22, 25, 25, 25, 30, 33, 33, 35, 35, 35, 36, 40, 45, 46, 52, 70.

* The first quartile (Q_1) is the 25th Percentile and the third quartile (Q_3) is the 75th percentile in a data set.

To find Q_1 and Q_3 , we first need to order the data set and find the medians. for odd number of elements in the data set.

Median = $(N+1)/2$ th element of the sorted dataset, where N is the number of elements in the dataset

for even number of elements in the dataset median = $(N/2$ th element + $(N/2 + 1$ th element)) / 2 of sorted dataset where N is the number of elements in the dataset.

Here we have 26 elements in the dataset, so the median is the average of the 13th and 14th elements which are 19 and 20 respectively

$$\text{therefore } Q_2 = (19 + 20) / 2 = 19.5$$

Now that we have Q_2 , we can find

Q_1 and Q_3 by finding the median of the lower and upper halves of the data set respectively

for the lower half of the dataset, we have the following values

18, 15, 16, 16, 19

the median of this set is 16,

$$\text{so } Q_1 = 16$$

for the upper half of the dataset, we have following values.

20, 20, 21, 22, 22, 25, 25, 30, 33, 33, 35,

35, 35, 35, 36, 40, 45, 46, 52, 70.

the median of this set is 35, so

$$Q_3 = 35$$

therefore the first quartile (Q_1) = 16 and the third quartile (Q_3) = 35