# R Lang:

**Creating a file:** using file.create() function, a new file can be created from console or truncates if already exist.

**Syntax:**
file.create(" ")

**Example:** # create a file
# The file created can be seen
# In your working directory.

to create the csv file. The file gets created in the working directory.

# create a data frame
data_read.csv("input.csv")
> as.Date(2004-10-01"))

# write filtered data into a new file.
write.csv(retval," output.csv")
print(newdata)

---

**Reading a file:** using read.table() function in R, files can be read and output is shown below.

**Syntax:** read.table (file)

**Ex:** # Reading txt file
new.fata=read.table(file= "nFa.txt")

# print
print(new.fata)

**Writing into a csv file**
R can create csv file form existing data frame. The write.csv() function is used

# create a data frame
retval-Subset(data, as.Date(startdate)

# write filtered data into a new file.
newdata <- read.csv(" output.csv")

---

{Data Manipulation}

round(x,n)  # round the values of x to n decimal places

ceiling(x)  # vector x of smallest integers > x.

floor(x)  # vector x of largest integer < x.

# x truncates real x to integers
as.integer #
(compare to round (x,0)).

---

var( ) → # produces the variance covariance matrix.

sd( ) → #standard deviation.

{Transformation}

fivenum( ) → # Tukey five numbers: min, lowerhinge, 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=T) → # centers around the mean and scales by the (sd)

# input and display
read.table(filehome, header=True)→
# read files with tables in first row
# read a tab 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 a n x 2 matrix.

---

{Statistics}

min() → lowest value from given data

mean() → Average value.

median() → Middle value q1,q2,q3

sum() → Total

Write down a list of 6 numbers where the mean is 0, the mode is 1 and the range is 3.

Write down 8 numbers, 1 to 8, and demonstrate how to calculate the interquartile range

Here are ten numbers.

| 7 | 6 | 8 | 4 | 5 | 9 | 7 | 3 | 6 | 7 |
|---|---|---|---|---|---|---|---|---|---|

a) Range
b) Mode
c) Median
d) Mean
e) Interquartile range

7. Find the interquartile range of the following data.: 12, 17, 34, 56, 32, 23, 83, 31, 9, 36

8. The points scored by a football team during each game from last season are: 14, 21, 18, 21, 14, 35, 42, 42, 21, 14 Calculate the mean, median, and mode: Mean: Median: Mode:

9. Use the data provided to answer the questions below: 14, 9, 3, 11, 13, 16, 4, 11, 19, 10, 8, 20, 13, 7, 12

Find the median of the data set. This value is known as Quartile 2 or Q2 for short

Find the median of the first half of the data set. This value is known as Quartile 1 or Q1.

. Find the median of the second half of the data set. This value is known as Quartile 3 or Q3.

Subtract the Q3 value from Q1. This is the Interquartile Range or IQR.

1.9 people take a test. Their scores out of 100100 are:

56, 79, 77, 48, 90, 68, 79, 92, 7156,79,77,48,90,68,79,92,71

Calculating the Range

2.Find the range of 12, 8, 4, 16, 15, 15, 5, 15, 10, 812,8,4,16,15,15,5,15,10,8

3.There were 55 members of a basketball team who had a mean points score of 1212 points each per game.

4.One of the team members left, causing the mean point score to reduce to 1010 points each per game. What was the mean score of the player that left?: Find the mode and range of the list of numbers below: 280, 350, 320, 400, 350,280,350,320,400,350, 490, 590, 470, 280, 410, 350490,590,470,280,410,350

5.Dani recorded the heights of members of her extended family, to the nearest cm. Find the median of their heights (listed below):

163, 185, 164, 170, 188, 154, 168, 179163,185,164,170,188,154,168,179

6.
For the following list of numbers, $1, 2, -2, 0, 1, 8, 3, -3, 2, 4, -2, 2$ find the:
a) Range
b) Mode
c) Median
d) Mean
e) Interquartile range

**Central Tendency:-** central around Statistics.

**Mean:-** The average of the number data set into Quartiles.

a Calculated.

$$\text{Mean:-} \frac{\text{Sum of all data values}}{\text{Total no of data values (n)}}$$

**Mean (census$ Total. Males):**

**Median:-** Middle Value (sorted)

Median (census $ Total. Males)

**Mode:-** MOST frequently occuring values.

Mode (census $ Total. Males)

**Measures of dispersion**

Interquartile Range → Quartiles, Mid range, Outlier

**Example of central Tendency.**

Ex:- 30,30, 32,38, 60

$$\text{Mean} = \frac{30+30+32+38+60}{5} = \frac{190}{5} = 38$$

Median = 32

Mode = 30.

---

**Inter Quartile Range:-** Measures of Variability, based on dividing a range ( census $ Total. Males)

**Example:-**

⟶ 2,3, 5, 7, 11,13,17,19,23, 29

Ten prime numbers are (Increasing order)

$Q_1 \to$ Lower Q.P (Quartile part)

$Q_2 \to$ Median

$Q_3 \to$ Upper Q.P

No of values =10 ⟶ Even number

∴ Median is Mean of 11 & 13

$$Q_2 = \frac{11+13}{2} = \frac{24}{2} = 12$$

$Q_1$ Part:- 2,3,5,7,11 ⟶ 5 ⟶ odd no.

Central value ⟶ 5 $\boxed{Q_1 = 5}$

$Q_3$ Part:- 13,17,19,23,29 ⟶ 5 (N⟶ 19.

The Subtraction $\boxed{Q_1 \leftarrow Q_3 \to 19-5 =1}$

$\boxed{11 \text{ is inter quartile Range Value.}}$

**Quartile** ( census $ Total. Males, 0.25)

**Quartiles:-** Compare set of observations

Ex:- 4,6,7,8,10, 23,24    n=7

Lower:- $Q_1 = \frac{(n+1)}{4} = \frac{7+1}{4} = 2^{nd}$ item ⟶6

Median: $Q_2 = \frac{n+1}{2} = \frac{7+1}{2} = 4^{th}$ item ⟶ 8

Upper:- $Q_3 = \frac{3(n+1)}{4} = 6^{th}$ item ⟶ 23

---

**Mid Range:-** The Arithmetic mean of largest & Smallest Value.

$$\text{Mid range} = \frac{(\text{Maximum Value} + \text{Minimum value})}{2}$$

Ex:- 2, 5,6, 7,8,9, 4 ⟶ $\frac{9+2}{2} = \frac{11}{2}$

**PLOTTING GRAPHS USING R-TOOL**

**Bar Plot:-** Stata:- barplot (data, xlab, ylab)

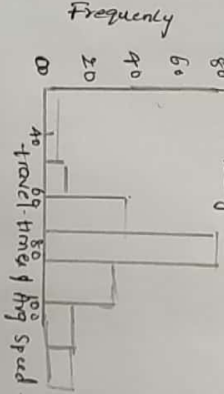Ex:-travel <- travel.times [which (travel $ Day of Week => Friday).

names (travel-times) <- in 1. c ("Day of Week", "Avg Speed")

**HISTOGRAM**

which displays of statistical information.

Hist (travel .times & Avg. Speed)
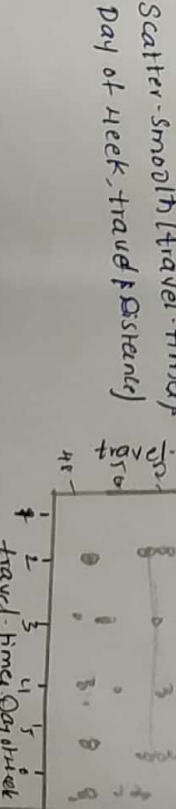
Histogram of travel.

**SCATTERPLOT:-** Extent of Correlation.

Scatter-smooth (travel-times, Day of Heek, travel $ distance)

# PERFORM CORRELATION ANALYSIS AND NORMALIZATION USING R-TOOL

## Correlation Analysis:

Steps involved :-

1. Create a new table with required dataframes

2. After that apply the formula or query for the chi-square test

Question:-

* diabetes1 ← table(diabetes $ Age . diabetes $ Insulin)
* diabetes1
* chi sq.test (diabetes1)

## Min Max Normalization formula

| Marks |
|---|
| 8 |
| 10 |
| 15 |
| 20 |

Min: The minimum value of the given attribute
Here Min is 8.

Max: The maximum value of given attribute.
Here Max is 20.

V: V is the representative value of attribute

for ex: $V_1 = 8$, $V_2 = 10$, $V_3 = 15$ & $V_4 = 20$
new Max: 1
new Min: 0

---

$$V' = \frac{V - min_A}{max_A - min_A}(new\_max_A - new\_min_A) + new\_min_A$$

for marks as 8:

$$Min\,Max = \frac{V - Min\,marks}{Max\,marks - Min\,marks}(new\,Max - new\,Min) + new\,Min$$

$$Min\,Max = \frac{(8-8)}{20-8} * (1-0) + 0$$

$$Min\,Max = \frac{0}{12} * 1$$

$$Min\,Max = 0$$

for marks as 10:

$$Min\,Max = \frac{(10-8)}{20-8} * (1-0) + 0$$

$$Min\,Max = \frac{(2)}{12} * 1$$

$$Min\,Max = 0.16$$

## Z-Score

standard deviation $= \sqrt{\dfrac{\Sigma\left(\begin{array}{c}\text{every individual value} \\ \text{of marks}\end{array} - \begin{array}{c}\text{mean of} \\ \text{marks}\end{array}\right)^2}{n}}$

$$S = \sqrt{\frac{\Sigma(x - \bar{x})^2}{n}}$$

mean of marks $= \dfrac{8 + 10 + 15 + 20}{4} = 13.25$

$$S = \sqrt{\frac{(8-13.25)^2 + (10-13.25)^2 + (15-13.25)^2 + (20-13.25)^2}{4}}$$

$$= \sqrt{\frac{27.56 + 10.56 + 3.06 + 45.56}{4}} = \sqrt{\frac{86.74}{4}} = \sqrt{21.6} = 4.6$$

---

Mean ($\mu$) = 13.25

standard deviation ($\sigma$) = 4.6

$$Z\,score = \frac{x - \mu}{\sigma} = \frac{8 - 13.25}{4.6} = -1.14$$

$$Z\,score = \frac{x - \mu}{\sigma} = \frac{10 - 13.25}{4.6} = -0.7$$

$$Z\,score = \frac{x - \mu}{\sigma} = \frac{15 - 13.25}{4.6} = 0.3$$

$$Z\,score = \frac{x - \mu}{\sigma} = \frac{20 - 13.25}{4.6} = 1.4$$

| marks | marks after z-score normalization |
|---|---|
| 8 | -1.14 |
| 10 | -0.7 |
| 15 | 0.3 |
| 20 | 1.4 |

Z-score Normalization:

* A ← C (diabetes $ Age)
* Mean ← mean(A)
* Std ← sd(A)
* Zscore ← (A-mean)|std
* Zscore

## Decimal Scaling formula

A value of attribute A is can be normalized by the formula

normalized value of attribute $= \dfrac{V'}{10^j}$

| Salary Bonus | formula | GPA Normalized after Decimal scaling |
|---|---|---|
| 400 | 400/1000 | 0.4 |
| 310 | 310/1000 | 0.31 |

Decimal scaling $= \dfrac{A}{100}$

# APRIORI ALGORITHM

| Transaction IP | ITEMS |
|---|---|
| T1 | Hot Dogs, Buns, Ketch up |
| T2 | Hot Dogs, Buns |
| T3 | Hot dogs, coke, chips |
| T4 | chips, coke |
| T5 | chips, Ketchup |
| T6 | Hot Dogs, Coke, chips |

Find the "Frequent Item Sets" & generate "Association Rules" on this. Assume, Min Support threshold

$$S = 33.33\%$$

Min Confident Threshold

$$C = 60\%$$

**Solution:-**

Min Support Count = $\dfrac{33.33}{100} \times 6$

Min Support Count = 2

| Item Set | Sup count |
|---|---|
| Hot Dogs | 4 |
| Buns | 2 |
| Ketchup | 2 |
| Coke | 3 |
| chips | 4 |

← (arrow)

| Item Set | Sup·count |
|---|---|
| Hot Dogs, Buns | 2 |
| Hot Dogs, Ketchup | 1 |
| Hot Dogs, coke | 2 |
| Buns, Ketchup | 1 |
| Buns, Coke | 0 |
| Buns, chips | 0 |
| Ketchup coke | 0 |
| Ketchup chips | 1 |

← (arrow)

| Item Set | Sup·count |
|---|---|
| Hot Dogs, Buns | 2 |
| Hot Dogs, coke | 2 |
| Hot Dogs, chips | 2 |
| Coke, chips | 3 |

← (arrow)

| Item Set | Sup·count |
|---|---|
| Hot Dogs, Buns, coke | 0 |
| Hot Dogs, Buns, chips | 0 |
| Hot Dogs, coke, chips | 2 |

← (arrow)

| Item Set | Sup·count |
|---|---|
| Hot Dogs, coke, chips | 2 |

* There is only one item set with min support 2. So,only one itemset is frequent.

Frequent Itemset (I) = {Hot Dogs,coke,chips}

**Association Rules:-**

* (Hot Dogs ∧coke) ⇒ [chips]

Conf = $\dfrac{\text{Sup (Hot Dogs ∧coke ∧chips)}}{\text{Sup (Hot Dogs ∧coke)}}$

= 2/(2×100) = 100% [Selected]

* (Hot Dogs ∧chips) ⇒ [coke]

Conf = $\dfrac{\text{Sup(Hot Dogs ∧coke ∧chips)}}{\text{Sup(Hot Dogs ∧chips)}}$

=2/(2×100) = 100% [Selected]

* (coke ∧chips) ⇒ [Hot Dogs]

Conf = $\dfrac{\text{Sup[Hot dogs∧coke∧chips]}}{\text{Sup (coke ∧ chips)}}$

= 2/(3×100) = 66.67% [Selected]

* [Hot Dogs] ⇒ [coke ∧chips]

conf = $\dfrac{\text{Sup(Hot Dogs∧coke∧chips)}}{\text{Sup(Hot Dog)}}$

= 2/(4×100) = 50% [Rejected]

* [coke] ⇒ [Hot Dogs ∧chips]

Conf = $\dfrac{\text{Sup (Hot Dogs ∧coke∧chips)}}{\text{Sup (coke)}}$

=2/[3×100]=66.67% [Selected]

* [chips] ⇒[Hot Dogs∧coke]

conf = $\dfrac{\text{Sup[Hot Dogs∧coke]}}{\text{Sup(chips)}}$

= 2/(4×100) = 50% [Rejected]

there are 4-Strong results [min.conf >60%]

* * * * *

## FP-Growth Algorithm

* Frequent Pattern -Growth Alg

| TID | ITEMS BOUGHT | [ORDERED] Freq-Item |
|---|---|---|
| 100 | f,a,c,d,g,i,m,P | f,c,a,m,P |
| 200 | a,b,c,f,l,m,o | f,c,a,b,m |
| 300 | b,f,h,j,o | f,b |
| 400 | b,c,k,s,P | c,b,P |
| 500 | a,f,c,e,l,P,m,n | f,c,a,m,P,v√ |

* Frequent Pattern - Type

* Frequent Pattern



Fig Frequent Pattern Tree

* Conditional -Frequent Pattern Tree {Conditional}

| Item | Conditional Pattern Base | FP-Tree |
|---|---|---|
| P | {{f,c,a,m:2},{c,b:1}} | {c:3} |
| m | {{f,c,a:2},{f,c,a,b:1}} | {f,c,a:3} |
| b | {{f,c,a:1},{f:1},{c:1}} | ∅ |
| a | {{f,c:3}} | {f,c:3} |
| c | {{f:3}} | {f:3} |
| f | ∅ | ∅ |

* Conditional Pattern Base (CPB)

| Item | conditional Pattern Base |
|---|---|
| P | {f,c,a,m:2},{c,b:1} |
| m | {{f,c,a:2},{f,c,a,b:1} |
| b | {f,c,a:1},{f:1},{c:1} |
| a | {f,c:3} |
| c | {f:3} |

(common in all Paths)

## Frequent Pattern rules

| Item | Conditional Pattern Base | Conditional FP-Tree | Freq- Pattern Generated |
|---|---|---|---|
| P | {{f,c,a,m:2}, {c,b:1}} | {c:3} | {<c,P:3>} |
| m | {{f,c,a:2},{f,c,a,b:1}} | {f,c,a:3} | {<f,m:3>,<c,m:3> <a,m:3>, <f,c,m:3> <f,a,m:3>,<c,a,m:3>} |
| b | {{f,c,a:1},{f:1},{c:1}} | ∅ | {} |
| a | {{f,c:3}} | {f;c:3} | {(f,a:3), <c,a:3> <f,c,a:3>} |
| c | {{f:3}} | {f:3} | {<f,c:3>} |
| f | ∅ | ∅ | {} |

Hence Frequent Pattern rules — generated —

| rec | Age | Income | Student | Credit _rating | Buys_computer |
|---|---|---|---|---|---|
| r1 | <=30 | High | No | Fair | No |
| r2 | <=30 | High | No | Excellent | No |
| r3 | 31..40 | High | No | Fair | Yes |
| r4 | >40 | Medium | No | Fair | Yes |
| r5 | >40 | Low | Yes | Fair | Yes |
| r6 | >40 | Low | Yes | Excellent | No |
| r7 | 31..40 | Low | Yes | Excellent | Yes |
| r8 | <=30 | Medium | No | Fair | No |
| r9 | <=30 | Low | Yes | Fair | Yes |
| r10 | >40 | Medium | Yes | Fair | Yes |
| r11 | <=30 | Medium | Yes | Excellent | Yes |
| r12 | 31..40 | Medium | No | Excellent | Yes |
| r13 | 31..40 | High | Yes | Fair | Yes |
| r14 | >40 | Medium | No | Excellent | No |

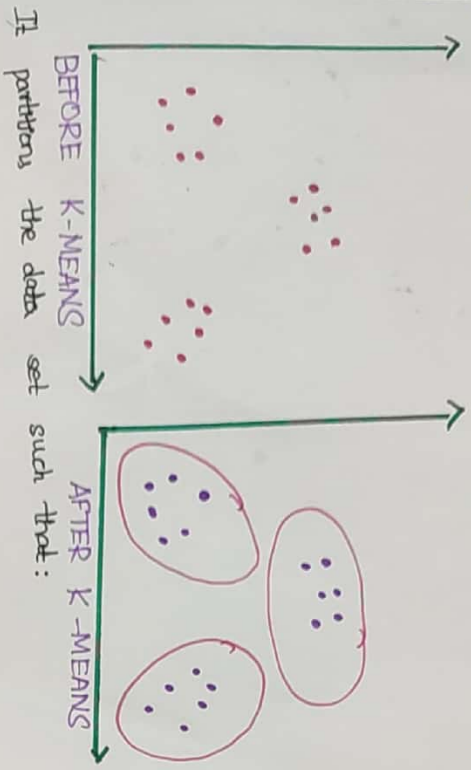| Day | outlook | temp | humidity | windy | play |
|---|---|---|---|---|---|
| 1 | sunny | hot | high | FALSE | no |
| 2 | sunny | hot | high | TRUE | no |
| 3 | overcast | hot | high | FALSE | yes |
| 4 | rainy | mild | high | FALSE | yes |
| 5 | rainy | cool | normal | FALSE | yes |
| 6 | rainy | cool | normal | TRUE | no |
| 7 | overcast | cool | normal | TRUE | yes |
| 8 | sunny | mild | high | FALSE | no |
| 9 | sunny | cool | normal | FALSE | yes |
| 10 | rainy | mild | normal | FALSE | yes |
| 11 | sunny | mild | normal | TRUE | yes |
| 12 | overcast | mild | high | TRUE | yes |
| 13 | overcast | hot | normal | FALSE | yes |
| 14 | rainy | mild | high | TRUE | no |

# K-Means Clustering:

* K-Means clustering is an unsupervised iterative clustering technique.

* It partitions the given data set into k predefined distinct clusters.

* A cluster is defined as a collection of data points exhibiting certain similarities.

→



→

BEFORE K-MEANS

↓

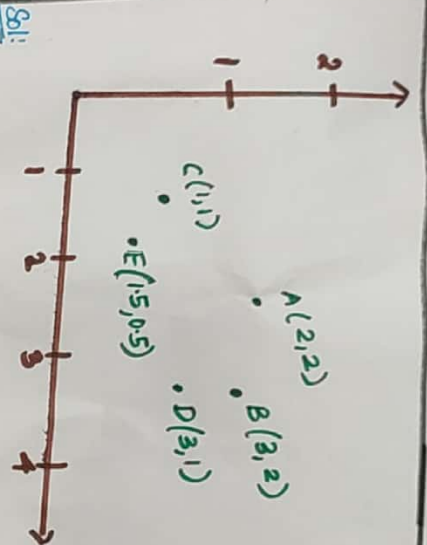It partitions the data set such that:

* Each data point belongs to a cluster with the nearest mean.

* Data points belonging to one cluster have high degree of similarity.

* Data points belonging to different clusters have high degree of dissimilarity.

**PROBLEM:**
Use K-Means Algorithm to create two clusters.



AFTER K-MEANS

→

---

Cluster-01: contains points A (2,2)
B (3,2)
E (1·5,0·5)

Cluster-02: contains points C (1,1)
D (3,1)

$A(2,2)$
$B(3,2)$
$C(1,1)$
$D(3,1)$
$E(1·5,0·5)$



**Sol!:**

**Iteration 01:**

We calculate the distance of each point from each of the center of the two clusters.

Distance is calculated by using euclidean distance formula.

↳ Calculating Distance Between $A(2,2)$ and $C_1(2,2)$:

$P(A, C1)$
$= \text{sqrt}\left[(x2-x_1)^2 + (y_2-y_1)^2\right]$
$= \text{sqrt}\left[(2-2)^2 + (2-2)^2\right] = 0$

↳ Calculating Distance Between $A(2,2)$ and $C_2(1,1)$:

$P(A, C2)$
$= \text{sqrt}\left[(x2-x_1)^2 + (y2-y_1)^2\right]$
$= \text{sqrt}\left[(1-2)^2 + (1-2)^2\right] = 1·41$

---

→ Now we re-compute the new cluster centers.

→ The new cluster center is computed by taking mean of all the points contained. In that cluster.

**For Cluster-01:**
Center of Cluster - 01
$= ((2+3+3)/3 , (2+2+1)/3)$
$= (2·67, 1·67)$

**For Cluster-02:**
Center of cluster - 02
$= ((1+1·5/2 , (1+0·5)/2)$
$= (1·25, 0·75)$

* Next, we go to iteration-02
iteration-03 and so on until the centers do not change anymore.

* This is completion of iteration-01.

| Given points | Distance from Center (2,2) | Distance from Center (1,1) | Point belongs to cluster |
|---|---|---|---|
| A(2,2) | 0 | 1·41 | C1 |
| B(3,2) | 1 | 2·24 | C1 |
| C(1,1) | 1·41 | 0 | C2 |
| D(3,1) | 1·41 | 2 | C1 |
| E(1·5,0·5) | 1·58 | 0·71 | C2 |