

dotplot(mtcars\$mpg, labels = row.names(mtcars),  
groups = cyl, gcolor = "black", lcolor = "black",  
ptch = 19, main = "Dot plot example",  
xlab = "Miles per Gallon")

## UNIT - 3

### Probability distribution:

Probability distribution lie at heart of statistics,  
so naturally R provides numerous functions for making  
use of function like generating random numbers &  
calculating the distribution & quantile.

\* To draw random numbers from normal distribution  
use rnorm (n) function.

Ex: rnorm (n=10)

Or else, we can also specify values

Ex: rnorm (n=10, mean = 100, sd = 20)

\* dnorm is used to calculate density (probability of a particular value) for the normal distribution.

Ex: x = rnorm (10)

dnorm (x)

\* dnorm return the probability of a specific number occurring.

Like with rnorm, a mean and standard deviation can be specified for dnorm.

Eg: `randNorm <- rnorm(10)`  
`randNorm`

`dnorm(randNorm)`

`dnorm(c(-1, 0, 1))`

Ex: `r <- rnorm(80000)`

`rdensity <- dnorm(r)`

`require(ggplot2)`

constructs aesthetic mapping.

`ggplot(data.frame(x=r, y=rdensity) + aes(x=x, y=`

`+ geom_point()) + labs(x="Random Normal Variable"`

`y="Density")`

Normal distribution: Is also known as the Gaussian distribution

Is the probability distribution that plots all of its values

In a symmetrical fashion, and most of the results

are situated around the probability's mean.

Many things follow a Normal Distribution:

(i) heights of a people.

(ii) size of things produced by machines

(iii) errors in measurement.

(iv) blood pressure.

(v) marks on a test.

The normal distribution has

Mean = median = mode.

Standard deviation.

The standard deviation is measure of how spread numbers are (or) Standard deviation

Variance: The average of squared differences from the mean.

Standard scores:

The no. of standard deviations from the mean is also called the "standard score", "sigma" or "z-score".

→ The opposite of Pnorm is znorm giving cumulative (property) Probability & from the quantile

Q: How many high must an individual

Binomial distribution.

$$P(x; n, p) = \binom{n}{x} p^x (1-p)^{n-x}$$

$$\binom{n}{x} = \frac{n!}{x! (n-x)!}$$

n - no. of trials      p - probability of success of trial

`rbinom(n=1, size = 10, prob=0.4)`

`rbinom(n=1, size = 10, prob=0.4)` ~~Output~~ 3

`rbinom(n=5, size = 10, prob=0.4)` [1] 5 3654

`rbinom(n=10, size = 10, prob=0.4)` [1] 5344453333

`binomData <- data.frame(Success = rbinom(n=100, size = 10, prob=0.4))`

`ggplot(binomData)`

- 1) find out the no. of possible outcomes when the coin is rolled (or) flipped for 10 times and also calculate the no. of success (no. of occurrences of heads)

$$P(X=0) = \frac{10!}{0!} = 1$$

- 2) find out the possible outcomes when a die is rolled for 5 times and calculate the no. of times the face 4 and no. of times face 6 occurs.

$$P(X=0) = \frac{10!}{0!} = 1$$

$$P(X=1) = \frac{10!}{1!} = \frac{10!}{9!} = 10$$

$$P(X=2) = \frac{10!}{2!} = \frac{10!}{2! 8!} = \frac{10 \times 9 \times 8!}{2 \times 8!} = 45$$

$$P(X=3) = \frac{10!}{3!} = \frac{10!}{3! 7!} = \frac{10 \times 9 \times 8 \times 7!}{3 \times 2 \times 7!} = 120$$

$$P(X=4) = {}^{10}C_4 = \frac{10 \times 9 \times 8 \times 7 \times 6!}{4 \times 3 \times 2 \times 1 \times 6!} = 210$$

$$P(X=5) = {}^{10}C_5 = \frac{10 \times 9 \times 8 \times 7 \times 6 \times 5!}{5 \times 4 \times 3 \times 2 \times 1 \times 5!} = 252$$

$$P(X=6) = {}^{10}C_6 = \frac{10 \times 9 \times 8 \times 7 \times 6 \times 5 \times 4!}{6 \times 5 \times 4 \times 3 \times 2 \times 1 \times 4!} = 210$$

$$P(X=7) = {}^{10}C_7 = \frac{10!}{7!(3!)!} = \frac{10 \times 9 \times 8 \times 7 \times 6 \times 5 \times 4 \times 3!}{7! \times 3 \times 2 \times 1 \times 3!} = 120$$

$$P(X=8) = {}^{10}C_8 = \frac{10!}{8!(2!)!} = \frac{10 \times 9 \times 8 \times 7}{8! \times 2} = 45$$

$$P(X=9) = {}^{10}C_9 = \frac{10!}{9! \times 1!} = \frac{10 \times 9!}{9! \times 1} = 10$$

$$P(X=10) = {}^{10}C_{10} = \frac{10!}{10! \times 0!} = \frac{1}{1} = 1$$

3. Write a function for calculating density binomial distribution for probability of the occurrence of 5 success out of 10 trials. The probability of success is 50%

## Basic Statistics:

Statistics: "a branch of mathematics used to summarize, analyze, and interpret a group of number of observations".

Population: any group of interest or any group that researchers

Two types of statistics

- Descriptive Statistics
- Inferential Statistics.

1. Descriptive: procedures used to summarize, organize set of measurements.

2. Inferential: procedures used that allow researchers to infer or generalize observations made with samples to the larger population from which they were selected.

Descriptive.

• use descriptive statistics to communicate with other researchers and the public.

• Descriptive statistics: central tendency and dispersion.

Measures of central tendency: we use statistical measures to locate a single score that is most representative of all scores in a distribution.

- Mean, Median, Mode.

- The notations used

Population:  $N$

Sample:  $n$

Mean for Sample -  $\bar{x}$

\* for population mean -  $\mu$

- Median:

- Data: 2, 3, 4, 5, 7, 10, 80 mean of these scores is 15.86

- 80 is an outlier (out of boundary)

- Mean fails to reflect most of the data, we use median instead of mean to remove the influence of an outlier.

- Median is the middle value in a distribution of data listed in a numeric order.

Mode

- The value in a dataset that occurs most often or most frequently.

Ex: 2, 3, 3, 3, 4, 4, 4, 4, 7, 8, 8, 8

Mode = 4

- Dispersion (variability): A measure of the spread of scores in a distribution.

- Range, Variance and standard deviation.

Range: difference b/w largest and smallest value

Variance measures the averaged squared distance that scores deviate from the mean.

$$s^2 = \frac{\sum (x - \bar{x})^2}{n-1} \text{ or } \frac{SS}{n-1}$$

↓  
degree of freedom: the no. of

scores in a sample that are free to vary.

Standard deviation: square root of variance.

### Correlation:

- The correlation is one of the most common and most useful statistics.
- A correlation is a single number that describes the degree of relationship between two variables.
- It is a statistical method which enables the researcher to find whether two variables are related and to what extent they are related.
- Correlation is considered as the syncretic movement of two or more variables.
- We can observe this when a change in one particular variable is accompanied by changes in other variables as well, as this happens either in the same or opposite direction, then the resultant variables are said to be correlated.
- The word correlation is made of co- (meaning "together") and Relation.

- correlation is positive when the value increase together, and
- correlation is negative when one value decrease as the other increases.

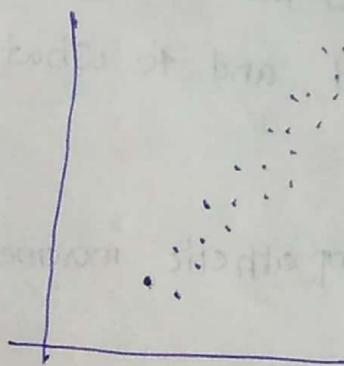
### correlation ~~is~~ of three types

- Positive
- Negative
- No correlation.

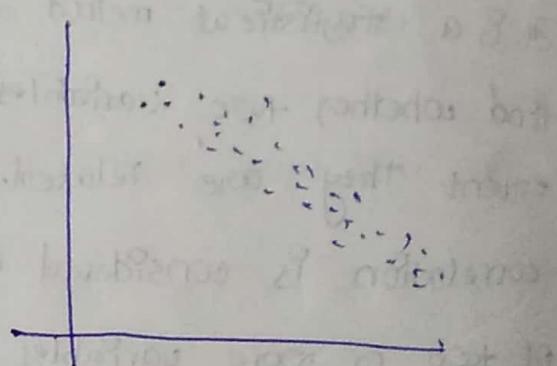
There is no change in a variable with any change in other variable.

### Correlation can have a values

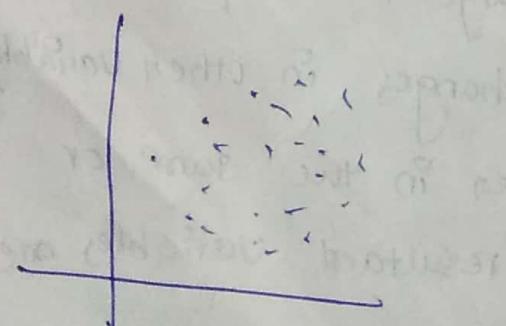
- 1 is a perfect positive correlation
- 0 is no correlation
- 1 is a perfect negative correlation.



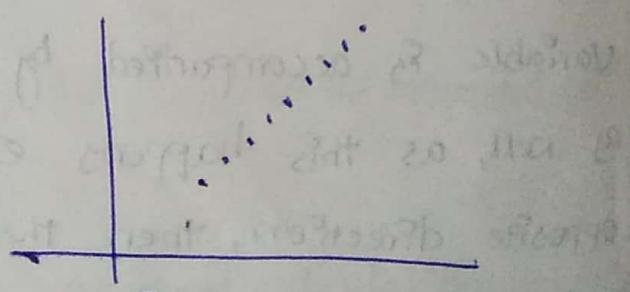
Positive correlation



Negative correlation



No correlation



Perfect positive correlation

If the points are close to each other than it is said to be strong correlation if the points are far from each other then it is negative correlation.

$$\text{correlation } r = \frac{n \sum xy - \sum x \sum y}{\sqrt{n \sum x^2 - (\sum x)^2} \sqrt{n \sum y^2 - (\sum y)^2}}$$

Ex:  
Let's assume that we want to look the relationship exist b/w Height and self-esteem.

Hypothesis Statement: correlation coefficient ( $r$ ) is equal to zero

Alternate Hypothesis: correlation coefficient ( $r$ ) is not equal to zero.

Person	Height (x)	Self Esteem (y)	$n=10$
1	68	4.1	
2	71	4.6	
3	62	3.8	
4	75	4.4	
5	58	3.2	
6	60	3.1	
7	67	3.8	
8	68	4.1	
9	71	4.3	
10	69	3.7	

$$\sum xy = 447561 \quad \sum x^2 = 447561$$

There are mainly three coefficients of correlation.

1. Karl Pearson's coefficient of correlation.
2. Pearson's rank
3. concurrent correlation

$$r = \frac{n \sum xy - \sum x \sum y}{\sqrt{n \sum x^2 - (\sum x)^2} \sqrt{n \sum y^2 - (\sum y)^2}}$$

$r = -1$  We say there is perfect negative correlation.

$r$  between  $-1$  &  $0$  We say that negative correlation.

$r = 0$  We say that it is no correlation.

$r = 1$  We say perfect positive correlation

$r = 0 & 1$  positive correlation.

Finding correlation coefficient can be done in 3 methods.

1. Pearson's method
2. Spearman's method.
3. Kendall's method

If we wanted to find in other method, the command is

`cor(x, y, method = "method-name")`

What is a t-test?

A t-test is a statistic

$$t = \frac{\text{variance between groups}}{\text{variance within groups}}$$

T test is often called Student's T test in the name of its founder "student".

- T test is used to compare two different set of values.
- Generally applied to normal distribution which has a small set of values.

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

$\bar{x}_1$  = mean of first set of values

$\bar{x}_2$  = mean of second set of values

$s_1$  = Standard deviation of first set of values.

$s_2$  = Standard deviation of second set of values.

$n_1$  = Total no. of values in first set

$n_2$  = Total no. of values in second set

The t-test, and any statistical test of this sort, consists of three steps:

- i) Define the null and alternate hypothesis.
- ii) calculate the t-statistics for the data.
- iii) compare  $t_{\text{calc}}$  to the tabulated t-value, for the appropriate significance level and degree of freedom.  
If  $t_{\text{calc}} > t_{\text{tab}}$ , we reject the null hypothesis and accept the alternate hypothesis. Otherwise, we accept the null hypothesis and vice versa.

Regular day exam

16 32

3 22

17 23  $\frac{+ 10.1}{15.1}$

3 13  $\frac{15.2}{15.9}$

19 20  $\frac{15.3}{15.9}$

15 29  $\frac{16.0}{15.8}$

24 11  $\frac{15.8}{16.6}$

23 25  $\frac{15.6}{14.9}$

3 13  $\frac{15.8}{15.8}$

12 20  $\frac{15.0}{15.4}$

11 32  $\frac{15.6}{15.8}$

$\underline{15.181}$

$\underline{16.2}$

$\underline{14}$

$\underline{15.17 - 21.63}$

$$\bar{x}_1 = 15.38$$

$$\bar{x}_2 = 15.68$$

$$S_1 = 0.31$$

$$S_2 = 0.40$$

$$t = \frac{15.38 - 15.68}{\sqrt{\frac{(0.34)^2}{n_1} + \frac{(0.4)^2}{n_2}}}$$

$$\begin{array}{r} 3 \\ \underline{2} \\ 2.6 \end{array}$$

$$\underline{2.6}$$

$$6$$

$$\underline{0.31 \times 0.31}$$

$$0.31$$

$$0.93$$

$$0.00$$

$$\underline{0.0961}$$

$$16$$

## ANOVA (Analysis of Variance)

After comparing two groups, the natural next step is comparing multiple groups.

$$F = \frac{\sum_{i=1}^k n_i (\bar{Y}_i - \bar{Y})^2 / (k-1)}{\sum_{i,j} (Y_{ij} - \bar{Y}_i)^2 / (N-k)}$$

$\bar{Y}$  is overall mean  $\bar{Y}_i$  is mean of group.

$Y_{ij}$  observation  $j$  in group  $i$ .

$n_i$  is no. of observations in group  $i$ .

$N$  - total no. of observations.  $k$  - no. of groups.

Grp								
1	4.5	A	10	6.	B	19	6.	C
2	5.1	A	11	8.	B	20	7.	C
3	4.	A	12	4	B	21	6.	C
4	3.	A	13	5	B	22	6.	C
5	2.	A	14	4	B	23	7	C
6	4.	A	15	6	B	24	5	C
7	3.	A	16	5	B	25	6	C
8	4.	A	17	8	B	26	5	C
9	4.	A	18	6	B	27	5	C
3.7	4							
	3.04444444							

$$F = \frac{9(3.7 - 5.22)^2 (3-1)}{7(4-3.7)^2 (27-3)} = \frac{(0.09)(24)}{(0.09)(24)}$$

If the between group variation is significantly greater than the within group variation then it is likely that there is a statistically significant difference between the groups.

→ If the b/w group variation is significantly larger than much larger than the within group variation, the means of different samples will not be equal (i.e. there is a statistically significant difference between the groups). If the b/w and within the group variations are approximately the same size.

⇒ aov() is used to perform ANOVA model.

Manipulating Strings  
(Static strings) —  
Paste: The function

Ex: Paste("Hello", "Jared", "and others")

[1] "Hello Jared and others".

→ sep, that determines what to put in b/w entries

Ex: Paste("Hello", "Jared", sep = "/")

[1] "Hello/Jared".

→ We can also use vectors.

Paste(c("Hello", "Hey", "Howdy"), c("Jared", "Bob", "David"))

[1] "Hello Jared" "Hey Bob" "Howdy David"

→ paste("Hello", c("Jared", "Bob", "David"))  
[1] "Hello Jared" "Hello Bob" "Hello David".  
→ paste("Hello", c("Jared", "Bob", "David"), c("Goodbye", "Seeya"))  
[1] "Hello Jared Goodbye", "Hello Bob Seeya", "Hello David Goodbye".  
→ vectorOfText <- c("Hello", "Everyone")  
paste(vectorOfText, collapse = " ")

(1)

sprintf  
While paste is convenient for static string. for dynamic  
things we use sprintf() function.

→ sprintf("Hello %s, Your party of %s will be seated in  
%s minutes", person, partySize, waitTime)

Note: %s was replaced with its corresponding variable.

Extracting Text

→ Often text needs to be ripped apart to be made useful,  
and while R has a no.of

→ readHTMLTable() to parse the table.

Ex: load("data/presidents.rdata")

theURL <- "http://

→ require(stringr)

> YearList <- strsplit(string = presidents\$YEAR, pattern = " - ")

> head(YearList)

# combine them into one matrix

> YearMatrix <- data.frame(reduce(rbind, YearList))

> head(YearMatrix)

give the columns good names

names(YearMatrix) <- c("start", "stop")

# bind the new column onto the data-frame

Presidents <- cbind(Presidents, YearMatrix)

# convert the start and stop columns into numeric.

Presidents\$start <- as.numeric(as.character(Presidents\$start))

Presidents\$stop <- as.numeric(as.character(Presidents\$stop))

Options:

strsplit(string, pattern = .n = 8nf, simplify = TRUE)

str-extract Extract matching pattern from a string.

str-count Count the no. of matches in a string.

str-ends Detects the presence or absence of pattern  
at the beginning or end of a string.

→ str\_sub = (presidents \$president, Start = 1, end = 3)  
# get the first 3 characters.

str\_sub(string, Start = 1L, end = -1L)

→ presidents [str\_sub(string = Presidents \$Start, Start = 4,  
end = 4) == 1, c("Year", "president", "Start", "Stop")]

### Regular Expression

# returns TRUE/FALSE if john was found in the name.

← str\_detect(Presidents \$president, "john")  
Pattern =

→ warTime [str\_detect(string = warTime, Pattern = " - ")]

→ which(str\_detect(string = warTime, Pattern = " - "))

→ TheStart ← sapply(theTimes, FUN = function(x) x[i])

→ head(TheStart)

# just return elements where January was detected.

theStart [str\_detect(string = theStart, Pattern = "January")]

# id is a shorthand for "[0-9]"

head(str\_extract(string = theStart, " \d{4}"), 20)

→ # extract 4 digits at the beginning of the text.

head(str\_extract(TheStart, Pattern = " \d{4}"), 20)

extract 4 digits at the beginning & ending of the text.

→ " " if (TheStart, Pattern = " \d{4}"), 20)

replace the first all digits seen with x.

head(str\_replace\_all(string = theStart, Pattern = " \d{4}."),

Write any Aspx for Implement Grid View

## LINEAR MODELS

- In its simplest form regression is used to determine the relationship b/w two variables. That is given one variable.
- That is given one variable, it tells us what we can expect from the other variable.
- This powerful tool, which is frequently taught and can accomplish a great deal of analysis with minimal effort, is called simple linear regression.
- II → The outcome variable (what we are trying to predict) is called the "response" and the input variable (what we are using to predict) is the "predictor".
- predictor to come up with some average value of the response.

$$Y = a + bx + c$$

•  $y$  varies with  $x$

- The aim of linear regression is to find a mathematical equation for a continuous response variable  $y$  as a function of one or more.

$$\text{ssy} = |\beta_1 - [1]| + |\beta_2 - [2]|$$

→ `linearmod <- lm(dist ~ speed, data = cars)`

`print(linearmod)`

`summary(linearmod)`

## partitioning Algorithms

partitioning method: construct a partition of a database  $D$  of  $n$  objects into a set of  $K$  clusters, s.t., min, max min sum of squared distance.

$$\sum_{m=1}^K \sum_{t_{mi} \in K_m} (e_m - t_{mi})^2$$

- Given  $\alpha K$ , find a partition
- k-means clustering Method
- Given  $K$ , the k-means algorithm  $P_0$  implemented in 4 steps:
  - partition objects into  $K$  nonempty subset.
  - compute seed points by the centroid of the clusters of the current partition (the center by the center, i.e., mean point, of the cluster)
  - Assign each object

$$A_1 \quad A_2 \quad A_3 \quad A_4 \quad A_5 \quad A_6 \quad A_7 \quad A_8$$
$$0 \quad \sqrt{25} \quad \sqrt{56} \quad \sqrt{19} \quad \sqrt{ } \quad \sqrt{ } \quad \sqrt{ } \quad \sqrt{ }$$

## Partitioning Algorithms: Basic concept:

Partitioning a database  $D$  of  $n$  objects onto a set of  $k$  clusters, such that the sum of squared distance  $P_S$  minimized (where  $C_P$  is the centroid or mediod of cluster  $C_P$ )

$$F = \sum_{i=1}^k \sum_{P \in C_i} (P - C_i)^2$$

→ Given  $k$ , find a partition of  $k$  clusters that optimizes the chosen partitioning criterion.

K-means

Given  $N$ , the K-means algorithm is implemented in four steps.

- partition objects into  $K$  nonempty subsets.
- compute seed points at the centroids of the clusters of the current partitioning (the centroid is the center i.e mean point, of the cluster)
- Assign each point object to the cluster with the nearest seed point.
- Go back to step 2, stop when the assignment does not change.

dis

- The K-means algorithm is sensitive to outliers.  
Since an object with an extremely large value may substantially distort the distribution of the data.

K-Medoids (PAM)

↓  
Partition around Medoids

- instead of taking the mean value of the object to a cluster as a reference point, medoids can be used, which is most centrally located object in a cluster.

Point	X-coordinates	Y-coordinates
1	4	6
2	2	6
3	3	8
4	8	5
5	7	4
6	4	7
7	6	0
8	4	3
9	6	4
10	3	4

Let us choose  $\underline{(3,4)}$  and  $\underline{(7,4)}$  are the medoids.

$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

$$(3,4) \quad (7,6)$$

$$x_1, y_1 \quad x_2, y_2$$

$$= \sqrt{(7-3)^2 + (6-4)^2}$$

$$= \sqrt{(4)^2 + (2)^2}$$

$$= \sqrt{16+4} = \sqrt{20}$$

$$(7,4), (7,6)$$

$$= \sqrt{(7-7)^2 + (6-4)^2}$$

$$= \sqrt{0+4} = \sqrt{4} = 2$$

$$(3,4) \quad (2,6)$$

$$(7,4), (2,6)$$

$$= \sqrt{(2-3)^2 + (6-4)^2}$$

$$= \sqrt{(2-7)^2 + (6-4)^2}$$

$$= \sqrt{(1)^2 + (2)^2}$$

$$= \sqrt{5+4} = \sqrt{9}$$

$$= \sqrt{5} \checkmark$$

$$= \sqrt{25+4} = \sqrt{29}$$

$$(3,4) \quad (3,8)$$

$$\sqrt{(3-3)^2 + (8-4)^2}$$

$$(7,4), (3,8)$$

$$\sqrt{(3-7)^2 + (8-4)^2}$$

$$\sqrt{0+16} = 4$$

$$(3,4), (8,5)$$

$$\sqrt{(8-3)^2 + (5-4)^2}$$

$$= \sqrt{(5)^2 + (1)^2} = \sqrt{26}$$

$$(3,4), (7,4)$$

$$= \sqrt{(7-3)^2 + (0)^2}$$

$$= \sqrt{(4)^2} = 4$$

$$(3,4), (4,7)$$

$$= \sqrt{(1)^2 + (3)^2}$$

$$= \sqrt{10}$$

$$(3,4), (6,2)$$

$$\sqrt{(8)^2 + (2)^2}$$

$$= \sqrt{9+4} = \sqrt{13}$$

$$(3,4), (7,3)$$

$$\sqrt{(4)^2 + (1)^2}$$

$$= \sqrt{16+1} = \sqrt{17}$$

$$(3,4), (6,4)$$

$$\sqrt{(8)^2 + 0} = \sqrt{9} = 3$$

$$(3,4), (3,4)$$

$$= 0$$

$$\sqrt{16+16} = \sqrt{32} = 4\sqrt{2}$$

$$(7,4), (8,5)$$

$$\sqrt{(8-7)^2 + (5-4)^2}$$

$$= \sqrt{1+1} = \sqrt{2}$$

$$(7,4), (7,4)$$

$$\sqrt{(7-7)^2 + (4-4)^2}$$

$$= 0$$

$$(7,4), (4,7)$$

$$\sqrt{(3)^2 + (3)^2}$$

$$= \sqrt{9+9} = \sqrt{18}$$

$$(7,4), (6,2)$$

$$\sqrt{(1)^2 + (2)^2}$$

$$= \sqrt{5}$$

$$(7,4), (7,3)$$

$$\sqrt{(0)^2 + (1)^2}$$

$$= \sqrt{1} = 1$$

$$(7,4), (6,4)$$

$$\sqrt{(1)^2 + 0} = \sqrt{1} = 1$$

$$(7,4), (3,4)$$

$$\sqrt{(4)^2 + 0} = 4$$

cluster 1 (3, 4)

(2, 6), (3, 8), (4, 7)

(3, 4)

cluster 2 (7, 4)

(7, 6), (8, 5), (7, 4), (6, 2)

(7, 3), (6, 4)

k-Mean

→ based on Mean

k-Medoids

→ based on Median.

k-Medoids

Algorithm.

1. Initially select k random points as the medoids from given n data points of the dataset.
2. Associate each data point to the closest medoid by using any of the most common distance metrics.
3. for each pair of non-selected object h and selected object s, calculate the total swapping cost  $T_{ch}$

If  $T_{ch} < 0$  s is replaced by h.

→ Now calculating the cost which is nothing but the sum of distance of each non-selected point from the selected point.

$$\begin{aligned}
 \text{Total cost} &= \text{cost}((3, 4), (2, 6)) + \text{cost}((3, 4), (3, 8)) + \text{cost}((3, 4), (4, 7)) \\
 &\quad + \text{cost}((7, 3), (7, 6)) + \text{cost}((7, 3), (8, 5)) + \text{cost}((7, 3), (6, 2)) \\
 &\quad + \text{cost}((7, 3), (7, 4)) + \text{cost}((7, 3), (6, 4)) \\
 &= 3 + 4 + 4 + 3 + 3 + 2 + 1 + 2 = \underline{\underline{22}}
 \end{aligned}$$

$$(7,5), (7,6)$$

$$\sqrt{0+(5)^2} = 5$$

$$(7,5), (8,6)$$

$$\sqrt{(7)^2 + (8)^2} = \sqrt{65}$$

$$(8,5), (7,6)$$

$$\sqrt{(8)^2 + (7)^2} = \sqrt{113}$$

$$(8,5), (8,6)$$

$$\sqrt{(8)^2 + (1)^2} = \sqrt{65}$$

$$(8,5), (3,8)$$

$$\sqrt{(8)^2 + (3)^2} = \sqrt{73}$$

$$(8,5), (8,5)$$

$$= 0$$

$$(8,5), (7,4)$$

$$\sqrt{(8)^2 + (1)^2} = \sqrt{65}$$

$$(8,5), (4,4)$$

$$\sqrt{(4)^2 + (2)^2} = \sqrt{20}$$

$$(8,5), (6,2)$$

$$\sqrt{(8)^2 + (3)^2} = \sqrt{113}$$

$$(8,5), (7,3)$$

$$\sqrt{(8)^2 + (2)^2} = \sqrt{65}$$

$$(8,5), (0,4)$$

$$\sqrt{(2)^2 + (1)^2} = \sqrt{5}$$

$$(8,5), (3,4)$$

$$\sqrt{(8)^2 + (1)^2} = \sqrt{65}$$

$$\underline{\text{cluster 1}} (3,4)$$

$$(7,6), (2,6), (3,8)$$

$$(4,4), (6,2)$$

$$\underline{\text{cluster 2}}$$

$$(8,5), (7,4), (7,3), (7,6)$$

$$(6,4)$$

$$(3,4) (3,8) \boxed{|x_2 - x_1| + |y_2 - y_1|} (3,8) (8,5)$$

$$|0| + |4| = 4$$

$$|5| + |3| = 8$$

$$(3,4) (7,6)$$

$$(8,5) (7,6)$$

$$|4| + |2| = 6$$

$$|1| + |1| = 2$$

$$(3,4) (2,6)$$

$$(2,6) (8,5)$$

$$|1| + |2| = 3$$

$$|6| + |1| = 7$$

$$(3,4), (8,5)$$

$$(8,5), (8,5)$$

$$|5| + |1| = 6$$

$$= 6$$

$$(3,4) (7,4)$$

$$|4| = 4$$

$$(8,5), (7,4)$$

$$= 1 + 1 = 2$$

$(3,4), (4,7)$

$$|1| + |3| = 4$$

$(3,4), (6,2)$

$$3+2=5$$

$(3,4), (7,3)$

$$3+1=4$$

$(3,4), (6,4)$

$$= 3$$

$\text{II} \quad (3,4) (3,4)$

$$= 0$$

$(8,5) (4,7)$

$$|4| + |2| = 6$$

$(8,5), (6,2)$

$$2+3=5$$

$(8,5), (7,3)$

$$1+2=3$$

$(8,5), (6,4)$

$$2+1=3$$

$(8,5), (3,4)$

$$= 5+1=6$$

$(2,6), (4,7)$

$$\text{Total cost} = 3+4+4+0+5+2+3+3 = 26$$

total cost  $(7,4) < \text{total cost } (8,5)$

$$20 < 26$$

$\therefore (7,4)$  is the medoid.

$(6,4) (3,4)$

$(6,4), (7,6)$

$$= 1+2=3$$

$(2,6) (6,4)$

$$4+2=6$$

$(3,8) (6,4)$

$$3+4=7$$

$(8,5), (6,4)$

$$2+1=3$$

$(4,7), (6,4)$

$$2+3=5$$

$(6,2), (6,4)$

$$= 0+2=2$$

$(7,3) (6,4)$

$$1+1=2$$

$(3,4) (6,4)$

$$= 3$$

cluster (A, B)  
(C, D), (B, E), (D, F)

(A, B)  
(C, D), (B, E), (A, F)  
(C, D) (E, F)

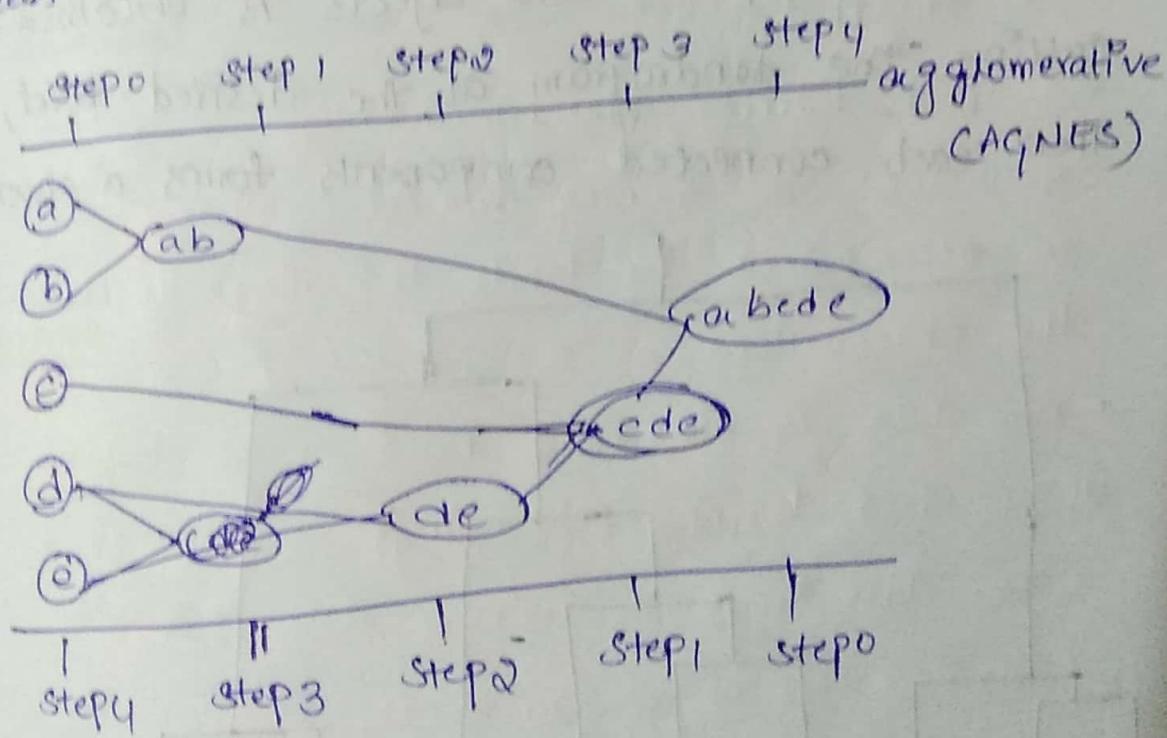
$$\text{Total cost} = 3+4+4+3+3+1+0+0 = 20$$

K-medoids  
~~K~~-cluster :: pam (points, 2)

out	x	y
5	4	4
2	2	6

### Hierarchical clustering

use distance matrix as clustering criteria. This method does not require the no. of clusters K as an input, but needs a termination condition.

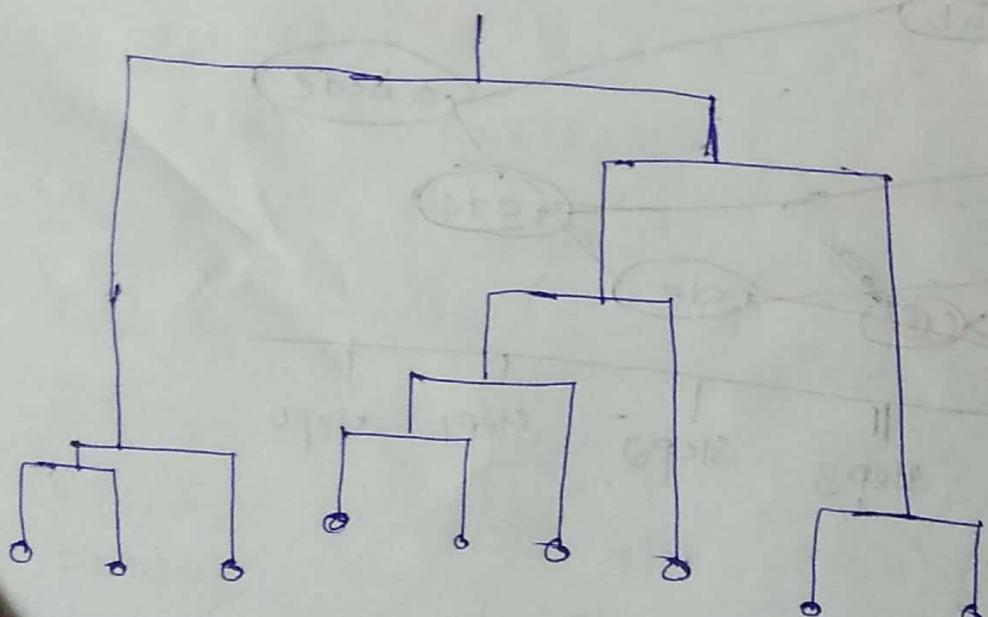


## AGNES

- introduced in Kaufman and Rousseeuw
- implemented in statistical packages e.g. SPSS
- use the single-link method and the dissimilarity matrix
- Merge nodes that have the least dissimilarity
- go in non-descending fashion

1) Dendrogram shows how clusters are merge  
Decompose data objects into several levels of nested partitioning (tree of clusters), called a dendrogram.

- A clustering of the data objects is obtained by cutting the dendrogram at the desired level, then each connected components forms a cluster



### Algorithm

- put each data point in its own cluster.
- identify the closest closest

### Codes

```
clusters <- hclust(dist(iris[1:3]))
```

↓  
method for hierarchical clustering.  
used to create distance matrix.

### To cut the dendrogram

```
clusterCut <- cutree(clusters, 3)
```

```
> table(clusterCut, Pros $ species)
```

out:

clusterCut	Setosa	Versicolor	Virginica
1	56	0	0
2	0	21	50
3	0	29	0

### DIANA

## Classification:

### Classification:

- predicts categorical class labels (discrete or nominal)
- classifies data (constructs a model) based on the training set and the values (class labels) in a classifying attribute and uses it to classify new data.

### Prediction:

- models continuous-valued functions, i.e. predicts unknown or missing values.

### Typical applications

- credit approval
- Target marketing
- Medical diagnosis
- fraud detection.

### Supervised learning (classification)

Supervision: The training data (observations, measurements) are accompanied by labels indicating the class of the observations.

- New data is classified based on the training set.

### Unsupervised learning (clustering)

- The class labels of training data is unknown.
- Given a set of measurements, observations etc with the aim of establishing the existence of classes or clusters in the data.

Issues: (Classification)

• Data cleaning

Preprocess data in order to reduce noise and handle missing values.

• Relevance analysis (feature selection)

Remove the irrelevant or redundant attributes.

• Data transformation

Decision Tree

Algorithm:

Generate, decision tree.

Generate a decision tree from the training tuples of data Partition D.

Inputs: