

## Types of Variables

**Qualitative (categorical)**

⇒ Grouped data

= express a qualitative attribute

such as hairclr, religion,  
marital status

Nominal:

↓  
No ordering/ranking  
is possible

Eg: ⇒ hairclr,  
eyecirc etc

Ordinal:

↓  
Ordering is  
possible

Eg: ⇒ healthy  
can be classified  
as good/bavg/bad

Interval

↓  
ratio of values of variable

do not have any meaning & it

does not have inherently defined  
zero value such as temp.

**Quantitative (numerical data)**

measured in terms of  
no's such as weight/height

No. without fractions

Discrete

↓  
countable & have  
finite no. of possibilities  
such as no. of people

↓  
can't be countable  
& has  $\infty$  no. of slns

Eg: ⇒ height

No. with fractions

Continuous

Ratio

↓  
ratio of values of variable  
have meaning & it have  
an inherently defined zero value  
(Eg: length)

ISM Important topics

## ① Measure of Central Tendency

↳ Mean, median, mode

\* mean:  $\bar{y} = \frac{\sum f y}{N}$

Mean =  $\frac{\text{sum of values}}{\text{tot no. of values}}$

\* Mode: - Most repeated no. of times

One mode  $\Rightarrow$  unimode

Multiple modes  $\Rightarrow$  multimode.

mode =  $3\text{median} - 2\text{mean}$

\* Median:-

↳ Segregate the given no.'s & find middle value

↳ if odd no. of values median = middle value

↳ if even no. of values median = Avg of 2 middle values

## ② Measure of variability

↳ Range, standard deviation, variance

\* RANGE:-

distance covered by scores in a distribution

↳ From smallest to highest value

$$\text{RANGE} = \text{URL} x_{\max} - \text{RL} x_{\min}$$

Eg:  $\Rightarrow [7, 2, 7, 6, 5, 6, 2]$ ; RANGE<sub>high-low</sub> = 7 - 2 = 5

## \* Standard deviation :-

↳ Before going to standard deviation lets first understand what's deviation.

Given data

| DAY  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  | 10 |
|------|----|----|----|----|----|----|----|----|----|----|
| TIME | 39 | 29 | 43 | 52 | 39 | 44 | 40 | 31 | 44 | 35 |

$$\bar{x} = \text{mean} = \frac{\text{sum of val}}{\text{tot val}} = \frac{39+29+43+\dots+35}{10} = \frac{396}{10} = 39.6$$

↳ Now if we subtract each of these values ( $x$ ) with Arthematic mean ( $\bar{x}$ ); that deviation tell about deviation from each of these values from Arthematic mean

$$\textcircled{1} x - \bar{x} \Rightarrow 39 - 39.6 = -0.6$$

$$\textcircled{3} 43 - 39.6 =$$

$$\textcircled{6} 44 - 39.6 =$$

$$\textcircled{9} 44 - 39.6 =$$

$$\textcircled{4} 52 - 39.6 =$$

$$\textcircled{7} 40 - 39.6 =$$

$$\textcircled{10} 35 - 39.6 =$$

$$\textcircled{5} 39 - 39.6 =$$

$$\textcircled{8} 31 - 39.6 =$$

$$\sum x - \bar{x} = 0 \Rightarrow \text{if we do sum of all deviations} = 0$$



It is mandatory that the algebraic sum of the deviations from Arthematic mean shld be always "Zero".

NOW AS SUM OF DEVIATIONS IS ALWAYS ZERO; HOW TO FIND STANDARD DEVIATION??

### STANDARD DEVIATION:-

#### 3 STEPS:-

- ↳ FIND THE DEVIATION OF EACH SCORE & SQUARE EACH DEVIATION
- ↳ SUM OF SQUARED DEVIATIONS.
- ↳ AVG THE SQUARED DEVIATIONS.

$$\text{Variance} = \frac{\sum (x - \bar{x})^2}{n}$$

$$\text{Standard deviation} = \sqrt{\text{Variance}}$$

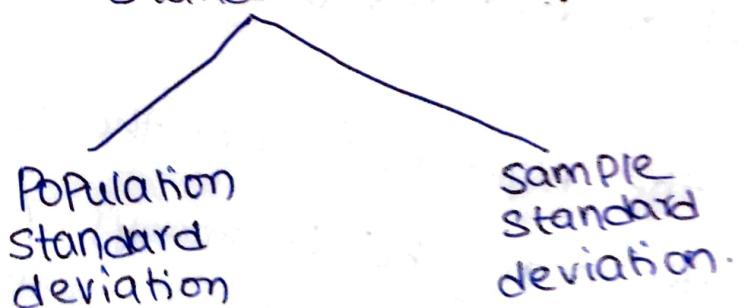
Eg:- {1, 5, 7, 8, 9}

$$\bar{x} = \text{mean} = \frac{1+5+7+8+9}{5} = 6$$

$$\sum (x - \bar{x})^2 = (1-6)^2 + (5-6)^2 + (7-6)^2 + (8-6)^2 + (9-6)^2 = 40$$

$$\text{Variance} = \frac{\sum (x - \bar{x})^2}{n} = \frac{40}{5} = 8$$

$$\text{Standard deviation} = \sqrt{\text{Variance}} = \sqrt{8} \checkmark$$



(i) Population standard deviation: - [entire data]

$$\text{Variance} = \frac{\sum (x - \mu)^2}{n} \quad [\text{if } n \geq 30]$$

(ii) Sample standard deviation:

$$\text{Variance} = \frac{\sum (x - \bar{x})^2}{n-1} \quad [\text{if } n < 30]$$

↳ More concepts :-

$$Q_1 = \left[ \frac{n+1}{4} \right]^{\text{th}} \text{ term value}$$

$Q_1$  = first Quartile  $\Rightarrow 25\%$ ,

$$Q_3 = 3 \left[ \frac{n+1}{4} \right]^{\text{th}} \text{ term value}$$

$Q_3$  = third Quartile  $\Rightarrow 75\%$ .

IQR :- [inter quartile RANGE] :-

↳ measure of variance

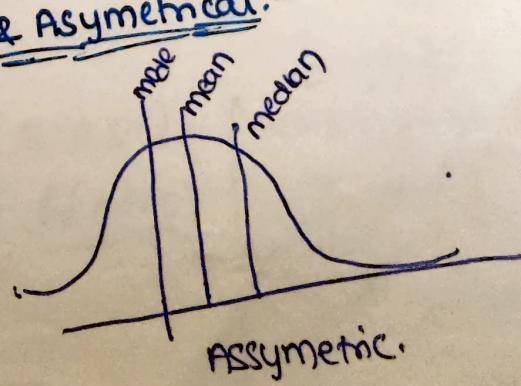
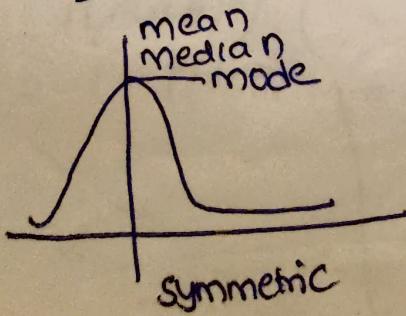
↳ Diff b/w 3rd & 1st quartile

$$\text{IQR} \Rightarrow Q_3 - Q_1$$

⇒ Five Point summary :-

To find five point summary; we need to find  
i) ~~mean~~ minimum ii) median (iii) ~~mode~~ maximum (iv)  $Q_1$  (v)  $Q_3$

⇒ DATA :- Symmetrical & Asymmetrical :-



## DATA distribution :- [skewed]

Negatively (left) : mean < median

Positively (Right) : mean > median

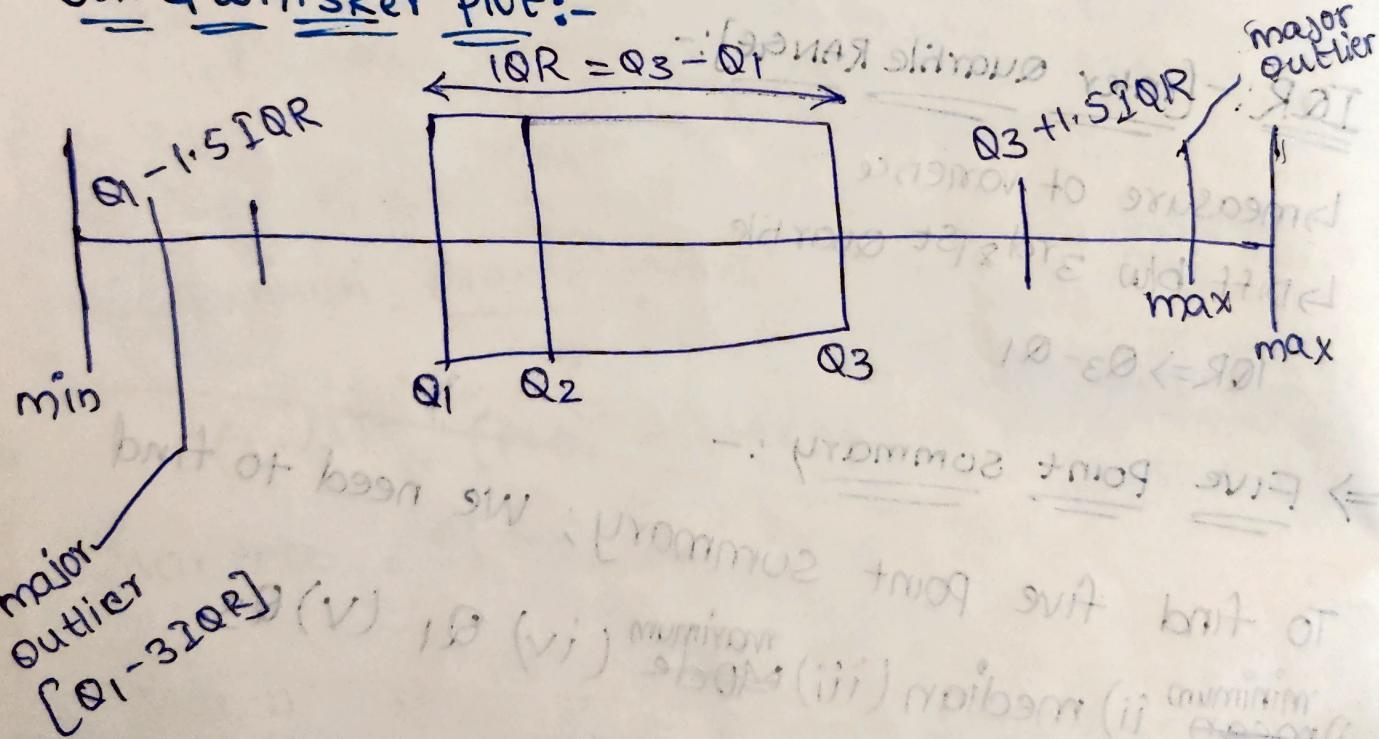
## Potential outliers :-

↳ The lower limit & upper limit of data set are given by

$$\text{lower limit} = Q_1 - 1.5 \times IQR$$

$$\text{upper limit} = Q_3 + 1.5 \times IQR$$

## Box & whisker plot :-



## Basic Probability :-

### ⇒ Random Experiment :-

Random Experiment is used to describe any action whose outcome is not known ~~as~~ in advance.

- Eg:- i) Give a lecture, How many r listening  
 ii) :- Pulling card from deck.

### ⇒ Sample Space :-

The Sample space of a random experiment is set 'S', that includes all possible outcomes of experiment.

### ⇒ Event :-

An Event is a set of outcomes of experiment. This includes the null (empty) set of outcomes & set of all outcomes.

Eg:- dice ⇒ Sample space = {1, 2, 3, 4, 5, 6}

Event = {even}, {odd}, {prime} ---



→ A & B are mutually Exclusive Events.

$$\Rightarrow P(A \cup B) = P(A) + P(B) - P(ANB)$$

### Mutually Exclusive :-

$ANB = \emptyset \Rightarrow$  nullset;  $P(ANB) = 0$

$$P(A) = 1 - P(\bar{A})$$

$$P(A \cup B) = P(A) + P(B)$$

## ⇒ Independent Events:-

we say both A & B are independent if:

$$P(A \cap B) = P(A) \cdot P(B)$$

## ⇒ Conditional Probability

The prob of event happening given that another event has already happened

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(B|A) = \frac{P(B \cap A)}{P(A)}$$

## ⇒ Multiplication Rule:-

let A, B, C are 3 events in Sample Space 'S' then:

$$P(A \cap B \cap C) = P(A) \cdot P(B|A) \cdot P(C|A \cap B)$$

$$P(A_1 \cap A_2 \cap A_3 \cap \dots \cap A_n)$$

$$= P(A_1) \cdot P(A_2|A_1) \cdot \dots \cdot P(A_n|A_1 \cap A_2 \cap \dots \cap A_{n-1})$$

## ⇒ Total Probability :-

The law of tot. Prob.

$$P(B) = \sum_{i=1}^K P(B|A_i) \cdot P(A_i)$$

## Bayer's theorem :-

It will help us to find the prob of something; given that we know something else is happened.

$$P(E_i/A) = \frac{P(E_i) \cdot P(A/E_i)}{\sum_{i=1}^n P(E_i) \cdot P(A/E_i)}$$

(Total prob pca)

## Naive Bayer's theorem :-

$$P(x_1, \dots, x_n | y) = \prod_{i=1}^n P(x_i | y)$$

It works Based on Bayes' theorem & assumes all features are independent to each other.

$$P(x_1, x_2 | y) = P(x_1 | y) \cdot P(x_2 | y)$$

⇒ It is basically used for problems like

$S_1/S_2$ ; Yes/No; Spam/Not spam; Disease/No disease

$$P(c/x_1, x_2, \dots, x_n) \propto P(c) \cdot P(x_1|c) \cdot P(x_2|c) \cdots P(x_n|c)$$

Eg:-

| No | clr   | legs | Hgt   | smelly | Species  |
|----|-------|------|-------|--------|----------|
| 1  | white | 3    | short | Yes    | $S_1, A$ |
| 2  | Green | 2    | tall  | No     | $S_1$    |
| 3  | Green | 3    | short | Yes    | $S_1$    |
| 4  | white | 3    | short | Yes    | $S_1$    |
| 5  | Green | 2    | short | No     | $S_2$    |
| 6  | white | 2    | tall  | No     | $S_2$    |
| 7  | white | 2    | tall  | No     | $S_2$    |
| 8  | white | 2    | short | yes.   | $S_2$    |

- i) find conditional prob of each attribute {clr, leg, hgt, smelly} for species  $\{S_1, S_2\}$
- $$P(S_1) = 4 \quad P(S_2) = 4$$
- |   |   |
|---|---|
| $P(\text{color} = \text{Green}   S_1) = \frac{2}{4}$          | $P(\text{color} = \text{Green}   S_2) = \frac{1}{4}$          |
| $P(\text{white}   S_1) = \frac{2}{4}$                         | $P(\text{white}   S_2) = \frac{3}{4}$                         |
| $P(\text{legs} = 2   S_1) = \frac{1}{4}$                      | $P(\text{legs} = 2   S_2) = \frac{4}{4}$                      |
| $P(\text{legs} = 3   S_1) = \frac{3}{4}$                      | $P(\text{legs} = 3   S_2) = \frac{0}{4}$                      |
| $P(\text{Hgt} = \text{short}   S_1) = \frac{1}{4}$            | $P(\text{Hgt} = \text{short}   S_2) = \frac{2}{4}$            |
| $P(\text{Hgt} = \text{tall}   S_1) = \frac{3}{4}$             | $P(\text{Hgt} = \text{tall}   S_2) = \frac{2}{4}$             |
| $P(\text{smelly} \rightarrow \text{Yes}   S_1) = \frac{3}{4}$ | $P(\text{smelly} = \text{NO}   S_2) = \frac{3}{4}$            |
| $P(\text{smelly} \rightarrow \text{NO}   S_1) = \frac{1}{4}$  | $P(\text{smelly} \rightarrow \text{Yes}   S_2) = \frac{1}{4}$ |

ii) using prob estimate which instance does this belong to?:

$\{C = \text{Green}, \text{leg} = 2, \text{Hgt} = \text{tall}, \text{smelly} = \text{NO}\}$

$$P(S_1/x) = \underbrace{P(x/S_1) \cdot P(S_1)}_{P(x)} ; P(S_2/x) = \underbrace{\frac{P(x/S_2) \cdot P(S_2)}{P(x)}}_{\text{denominator } x = \text{so compare numerators}}$$

$$P(\frac{S_1}{x}) = P(\frac{\text{Green}}{S_1}) \cdot P(\frac{2}{S_1}) \cdot P(\frac{\text{tall}}{S_1}) \cdot P(\frac{\text{NO}}{S_1}) \cdot P(S_1)$$

$$P(\frac{S_2}{x}) = P(\frac{\text{Green}}{S_2}) \cdot P(\frac{2}{S_2}) \cdot P(\frac{\text{tall}}{S_2}) \cdot P(\frac{\text{NO}}{S_2}) \cdot P(S_2)$$

$$P(\frac{S_1}{x}) > P(\frac{S_2}{x})$$

so  $S_1$  is the instance it belongs to.

Naïve Bayes Smoothing - due to diff. conditions but  
 As we saw in naïve Bayes theorem,  $P(A/x)$  is calculated as  
 $P(A_1, A_2, \dots, A_n/x)$  so even if 1 prob is 0 entire thing comes  
 out as zero.

$\Rightarrow$  So to avoid this we add 1 to every count, including unseen event

$$P(w/c) = \frac{\text{count}(w, c) + 1}{\sum \text{count}(w_i, c) + \text{no. of unique words}}$$

Eg:-  $P(\text{A very close Game} / \text{Sports}) = ?$

| text                             | Tag    |
|----------------------------------|--------|
| "A Great Game"                   | Sport  |
| the election was over            | NSport |
| "A Clean but unforgettable Game" | Sport  |
| "it was close election           | NSport |

$$\frac{(2)^9 \cdot (2^8)^4}{(2^9)^{13}} = (x^2)^9 \cdot (x^8)^4 = (x^{18})^9 \cdot (x^{32})^4 = (x^{50})^9$$

$$(2)^9 \cdot \left(\frac{2^8}{2}\right)^4 \cdot \left(\frac{8}{12}\right)^9 \cdot \left(\frac{5}{12}\right)^9 \cdot \left(\frac{10}{12}\right)^4 = (2)^9$$

$$(2)^9 \cdot \left(\frac{2^8}{2}\right)^4 \cdot \left(\frac{4^2}{2}\right)^9 \cdot \left(\frac{5}{2}\right)^9 \cdot \left(\frac{4^2}{2}\right)^4 = (2)^9$$

$$\left(\frac{2}{2}\right)^9 < (2)^9$$

## Random Variable :-

R.V

continuous R.V

Prob density  
fun (PDF)

cumulative distribution  
fun  $f(x)$

Discrete R.V

Prob mass  
fun  $f(x) = P(x)$

cumulative distribution  
fun  $\Rightarrow F(x)$

### Discrete R.V

$$\sum f(x) = 1 ; f(x) \geq 0$$

mean / Expected value of Discrete R.V :-

$$E(x) = \mu = \sum x \cdot f(x)$$

Variance of Discrete R.V is :-

$$\text{Var}(x) = \sigma^2 = E(x - \mu)^2 = \sum (x - \mu)^2 \cdot f(x)$$

(Or)

$$\text{Var}(x) = \sigma^2 = E(x^2) - [E(x)]^2 = E(x^2) - \mu^2$$

$$E(x^2) = \sum x^2 \cdot f(x)$$

$\Rightarrow$  cumulative prob distribution fun: (CDF) of Discrete RV

$$F(x_1) = f_1 = P(X=x_1)$$

$$F(x_2) = f_1 + f_2 = P(X=x_1) + P(X=x_2)$$

$$\vdots$$
  
$$F(x_n) = f_1 + f_2 + \dots + f_n = P(X=x_1) + P(X=x_2) + \dots + P(X=x_n)$$

V.R 2nd year

V.R 2nd year

V.R 2nd year

V.R 2nd year

$\Rightarrow$  Continuous Random Variable :-

i)  $f(x) \geq 0$  ii)  $\int_{-\infty}^{\infty} f(x) dx = 1$  iii)  $P(a \leq X \leq b) = \int_a^b f(x) dx$

↳ Area under  $f(x)$  from  $a \rightarrow b$   
for any  $a \neq b$

$\hookrightarrow$  Basic formulae of  $\int$  & derivatives by Sir:-

Integrations :-

i)  $\int_a^b x^n dx = \frac{1}{n+1} [x^{n+1}]_a^b = \frac{1}{n+1} [b^{n+1} - a^{n+1}]$

ii)  $\int_a^b e^{ax} dx = \left[ \frac{e^{ax}}{a} \right]_a^b$

(iii)  $\int a dx = ax + C$  [const.]

(a)  $f(x) = 3$

## derivatives

$$i) \frac{d}{dx}(x^n) = n \cdot x^{n-1}$$

$$ii) \frac{d}{dx}(e^{ax}) = a \cdot e^{ax}$$

$\Rightarrow$  Mean & variance of continuous R.V :-

$$i) \mu = E(x) = \int_{-\infty}^{\infty} x \cdot f(x) dx \Rightarrow \text{Mean}$$

$$ii) V(x) = \int_{-\infty}^{\infty} x^2 f(x) dx - \mu^2 \Rightarrow \text{variance}$$

$$iii) \text{standard deviation} = \sqrt{\text{variance}}$$

$$e = 2.718$$

## Joint Prob distribution :-

$(x = X, y = Y)$  &  $(x = X - \bar{x}, y = Y - \bar{y})$  be a discrete R.V

let  $x = \{x_1, x_2, \dots, x_n\}$  &  $y = \{y_1, y_2, \dots, y_m\}$

then  $P(x, y) = J_{ij}$  is called joint prop distribution of  $x \& y$

If it satisfies condition

$$i) J_{ij} \geq 0 \quad ii) \sum_{i=1}^m \sum_{j=1}^n J_{ij} = 1$$

|          | $y_1$    | $y_2$    | --- | $y_n$    |            |
|----------|----------|----------|-----|----------|------------|
| $x_1$    | $J_{11}$ | $J_{12}$ | --- | $J_{1n}$ | $f(x_1)$   |
| $x_2$    | $J_{21}$ | $J_{22}$ | --- | $J_{2n}$ | $f(x_2)$   |
| $\vdots$ | $\vdots$ | $\vdots$ | --- | $\vdots$ | $\vdots$   |
| $x_m$    | $J_{m1}$ | $J_{m2}$ | --- | $J_{mn}$ | $f(x_m)$   |
| Sum      | $g(y_1)$ | $g(y_2)$ | --- | $g(y_n)$ | $\sum = 1$ |

Note: Marginal distribution of  $x$  &  $y$  alone are:

\* Discrete  $\Rightarrow g(x) = \sum_y f(x,y)$  &  $h(y) = \sum_x f(x,y)$

\* Continuous  $\Rightarrow g(x) = \int_{-\infty}^{\infty} f(x,y) dy$  &  $h(y) = \int_{-\infty}^{\infty} f(x,y) dx$

conditional Prob of distribution func:-

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad P(B|A) = \frac{P(B \cap A)}{P(A)}$$

if  $x, y$  is Discrete R.V the condition prob dist'n func?

$$P(X=x, Y=y) = \frac{P(X=x, Y=y)}{P(Y=y)}$$

If  $x, y$  is continuous R.V then:

$$f_{X|Y} = \frac{f(x,y)}{f_Y(y)}$$

$f_{x,y}(x,y) \Rightarrow$  Joint PDF of  $x, y$

$f_Y(y) \Rightarrow$  Marginal PDF of  $y$

Bernoulli distribution :-

success (1) with prob  $\Rightarrow P$

success (0) " "  $\Rightarrow (1-P) = q$

mean =  $P$

variance =  $Pq$

$$\text{PMF: } P(X=x) = \begin{cases} P^x q^{1-x} & \text{if } x=1,0 \\ 0 & \text{elsewhere} \end{cases}$$

## Binomial distribution :-

Binomial distribution will be applied from following experimental conditions

- ⇒ No. of trials ( $n$ ) is  $\Rightarrow$  finite.
- ⇒ Trials are Independent to each other
- ⇒ Prob of success  $P$  is constant of each trial
- ⇒ Each trial has 2 mutually exclusive events success/fail

|                             |             |  |
|-----------------------------|-------------|--|
| Mean = $E(x) = np$          | $(1-p) = q$ | ④ PMF<br>$= P(X=x) = \binom{n}{x} p^x (1-p)^{n-x}$<br>$= nC_x p^x q^{n-x}$ |
| Variance = $\sigma^2 = npq$ |             | $x \Rightarrow$ No. of success in $n$ trials<br>$n \Rightarrow$ tot trails |
| $SD = \sigma = \sqrt{npq}$  |             |  |

## Poisson Distribution :-

Poisson distribution is suitable for rare events for which prob of occurrence  $(p)$  is very small & no. of trials  $n$  is very large

$$PMF: \Rightarrow P(x) = \begin{cases} e^{-\lambda} \frac{\lambda^x}{x!} & x = 0, 1, 2, \dots \\ 0, & \text{elsewhere} \end{cases}$$

$$\text{mean} = \lambda$$

$$\mu = \lambda \text{ in Poisson}$$

## Normal Distribution

$$E(X) = \mu = \int_{-\infty}^{\infty} x \frac{1}{\sigma \sqrt{2\pi}} \cdot e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2} dx$$

$$\text{Var}(x) = \sigma^2 = \left( \int_{-\infty}^{\infty} x^2 \frac{1}{\sigma \sqrt{2\pi}} \cdot e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2} dx \right) - \mu^2$$

$$Z = \frac{x-\mu}{\sigma}$$

$$z =$$

mean, SD

→ neither nor  $\Rightarrow (A \cup B)^c$

2) or/and  $\Rightarrow A \cup B \quad \text{and} \Rightarrow A \cap B \quad \text{or} \Rightarrow A \cup B$

$$\left(\frac{n+1}{4}\right)^{th} \text{ term}$$

$$\text{mr } x = \int_{-\infty}^{+\infty} f(x) dy; \quad y = \int_{-\infty}^{+\infty} f(y) dx$$

Naive Bayes  $y/N \Rightarrow \text{sunny}$  kinda Q&A

$$P(Y) = ? \quad P(N) = ?$$

$$P(Y|S) = ? \quad P(YN|S)$$

$$P(Y|S) = \frac{P(S|Y) \cdot P(Y)}{P(S)}$$

in Naive Bayes

$$P(N|S) = \frac{P(S|N) \cdot P(N)}{P(S)}$$

$$P(S|Y) = \frac{P(SNY)}{P(Y)} = \text{find & substitute in above -}$$

# i) Central Limit Theorem (CLT)

$$(i) Z \text{ score} = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$

$n > 30$

Avg mean =  $\mu$   
Prob mean =  $\bar{x}$

$\sigma$  = standard dev

$n$  = tot values

$$(ii) Z \text{ score} = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} \left[ \sqrt{\frac{N-n}{N-1}} \right]$$

→ we use (ii) when  $N$  is given &  $\frac{n}{N} \geq 0.5$

→ if  $\frac{n}{N} < 0.5 \Rightarrow$  we use (i) only       $N \Rightarrow$  not needed in this case.

## 2) Sample distribution of Sample proportion (%) :-

→ In such questions  $x$  will be given (or) directly  
they ask sample proportion

$$i) P = \frac{x}{n}$$

$x \Rightarrow$  sample workers (category)  
 $n \Rightarrow$  tot workers

$$Z = \frac{\hat{P} - P}{\sqrt{\frac{P(1-P)}{n}}}$$

$P \Rightarrow$  Population Prop =  $\frac{x}{n}$ .  
 $\hat{P} \Rightarrow$  sample proportion  
 $q \Rightarrow 1 - P$   
 $n =$  Sample size.

## $\Rightarrow$ Statistical Inference:

- └ Point estimation
- └ Interval estimation
- └ Hypothesis estimation

## $\Rightarrow$ Standard deviation of Sample Prop! -

| finite population  | Infinite Population                          |
|--|--|
| $\sigma_{\bar{P}} = \sqrt{\frac{N-n}{N-1} \cdot \frac{P(1-P)}{n}}$ | $\sigma_{\bar{P}} = \sqrt{\frac{P(1-P)}{n}}$ |

## Confidence Interval

if confidence interval = 95

then  $(1-\alpha) = 95\%$   
 $\alpha = 5\%$ .

This is based on  $N$ .

|   |   |
|---|---|
| $\bar{x} \pm z_{\alpha/2} \left[ \frac{\sigma}{\sqrt{n}} \right]$ | $\bar{x} \pm z_{\alpha/2} \left[ \frac{\sigma}{\sqrt{n}} \right] \left[ \sqrt{\frac{N-n}{N-1}} \right]$ |
|---|---|

∴ confidence interval is

$$\bar{x} - z_{\alpha/2} \left( \frac{\sigma}{\sqrt{n}} \right) < \mu < \bar{x} + z_{\alpha/2} \left( \frac{\sigma}{\sqrt{n}} \right)$$

↳ same way for other too.

⇒ In Qsn if they give conf interval = 95%.

what is  $\alpha \Rightarrow 1 - \alpha = 95\% ; \text{ so } \alpha = 5\%$ .

If  $\alpha = 1\%$ .

$$z_{\alpha/2} = 2.58$$

If  $\alpha = 5\% \Rightarrow 0.5$

$$z_{\alpha/2} \Rightarrow 1.96$$

If  $\alpha = 10\%$ .

$$z_{\alpha/2} = 1.645$$

$$\bar{x} - z_{\alpha/2} \sqrt{\left( \frac{s_1^2}{n_1} + \frac{s_2^2}{n_2} \right)}$$

for 2 means  
& 2 samples

if  $s_p$  is given

↳ Pooled  $s_p$

then .

$$\bar{x} - z_{\alpha/2} [s \in (\bar{x}_1 - \bar{x}_2)]$$

$$* \quad \downarrow \quad s_e(\bar{x}_1 - \bar{x}_2) = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \quad (8) \quad \sqrt{s_p \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}$$

$$- \quad (\bar{x}_1 - \bar{x}_2) \pm z_{\alpha/2} [s_e(\bar{x}_1 - \bar{x}_2)].$$

# HYPOTHESIS

$\Rightarrow$

NULL

$\Rightarrow$  if  $\mu_1 = \mu_2$

$\Rightarrow$  NO diff

$\Rightarrow$  Unbiased/Neutral

Alternative

$\Rightarrow$  if  $\mu_1 \neq \mu_2$   $\mu_1 > \mu_2$

$\Rightarrow$  diff  $\rightarrow$  Yes

$\Rightarrow$  Biased.

## level of significance ( $\alpha$ )

|                                | 1% $\alpha$                 | 5% $\alpha$                 | 10% $\alpha$                 |
|--------------------------------|-----------------------------|-----------------------------|------------------------------|
| two tailed test ( $H_a \neq$ ) | $ z_{\alpha/2}  = \pm 2.58$ | $ z_{\alpha/2}  = \pm 1.96$ | $ z_{\alpha/2}  = \pm 1.645$ |
| Right-tailed test ( $H_a >$ )  | $z_{\alpha} = 2.33$         | $z_{\alpha} = 1.645$        | $z_{\alpha} = 1.28$          |
| left-tailed test ( $H_a <$ )   | $z_{\alpha} = -2.33$        | $z_{\alpha} = -1.65$        | $z_{\alpha} = -1.28$         |

$\Rightarrow$  IN Right tailed

If  $z > z_{\alpha}$   
L reject

in left tailed

If  $z < z_{\alpha}$   
L reject

2-tailed

If  $-z_{\alpha/2} \leq z \leq +z_{\alpha/2} \Rightarrow$  Accept.  
L else reject

## Test significance for proportion :-

$$Z_P = \frac{P - P_0}{\sigma_P}$$

$$\sigma_P = \sqrt{\frac{PQ}{n}} \quad \left( \frac{N-n}{N-1} \right)$$

$Z$  score  
 $n > 30$

$$\text{If } N \text{ is unknown} \Rightarrow \sigma_P = \sqrt{\frac{PQ}{n}}$$

$$\text{If } P \text{ is unknown} \Rightarrow \sigma_P = \sqrt{\frac{Pq}{n}}$$

$P \Rightarrow$  Prop of survival  $\Rightarrow \frac{x}{n}$

$$q = (1-p)$$

$P_0 \Rightarrow$  Given prob proportion

$Q \Rightarrow 1-P \quad n \Rightarrow$  sample size

Difference of proportion  $\Rightarrow$  Test significance

$$Z = \frac{(P_1 - P_2) - (P_1 - P_2)}{\sqrt{\frac{P_1 q_1}{n_1} + \frac{P_2 q_2}{n_2}}}$$

$$Z = \frac{P_1 - P_2}{\hat{P} (1 - \hat{P}) \left[ \frac{1}{n_1} + \frac{1}{n_2} \right]}$$

$$\hat{P} = \frac{x_1 + x_2}{n_1 + n_2}$$

T-test ( $n < 30$ )

$$t = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$

$n < 30$   
T score

## T-test

diff b/w mean of 2 population :-

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n} + \frac{s_2^2}{n}}} \quad (31)$$

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{s_p \left( \frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

$s_p \Rightarrow$  Pooled SD

$$s_p = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1 + n_2 - 2}$$

$n < 30$

$\mu$  isn't

If  $\mu$  isn't given then

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

Steps

- 1) State  $H_0$  &  $H_a$
- 2)  $\alpha$
- 3) t-distribution
- 4) compute t-score
- 5) reject/accept  $H_0$

difference b/w mean of Paired observations  $\bar{d}$  - before  $\bar{d}_1$  after  $\bar{d}_2$

$$t = \frac{\bar{d} - \bar{d}_0}{s_d / \sqrt{n}}$$

$$\bar{d} = \frac{\sum_{i=1}^n d_i}{n}$$

$$s_d = \sqrt{\frac{\sum_{i=1}^n (d_i - \bar{d})^2}{n-1}}$$

### chi-square-test

$$\chi^2 = \sum_{i=1}^m \sum_{j=1}^n \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

$$E_{ij} = \frac{g_i c_j}{n}$$

$$df = (r-1)(c-1)$$

ii) Chi-square  $\Rightarrow$  If calculated  $\chi^2 >$  critical  $\chi^2$

$\hookrightarrow$  Reject null hypothesis

df &  $\alpha$  — See table find critical  $\chi^2$ .

If  $\alpha$  isn't given  $\Rightarrow \alpha = 0.05$

### ANOVA :-

i) one-way ANOVA

q

ii) two-way ANOVA

## One way ANOVA

### Step 1

Sum of all values  $\Rightarrow G$

$$CF = \frac{G^2}{n}$$

### Step 2

Sum of sq of each element  $\Rightarrow \sum \sum x_{ij}^2$

$$\underline{\text{Step 3}} \quad \sum \sum x_i^2 - CF = SST$$

$$\underline{\text{Step 4}} : SSTR = \frac{\sum C_i^2}{n} \Rightarrow \text{each column sum tot}$$

$$\underline{\text{Step 5}} \quad SSTR = SST - SSW$$

(Parametric test)

$\alpha = 0.05$   
L if not given

Eg:-

|   | Grp <sub>1</sub> | Grp <sub>2</sub> | Grp <sub>3</sub> | Grp <sub>4</sub> |
|---|------------------|------------------|------------------|------------------|
| 1 | a                | b                | c                | d                |
| 2 | e                | f                | g                | h                |
| 3 | i                | j                | k                | l                |
| 4 | m                | n                | o                | p                |

### Step 1

$$a+b+c+d + \dots + p = G$$

$$CF = \frac{G^2}{n}$$

### Step 2

$$\sum x_{ij}^2 = a^2 + b^2 + c^2 + d^2 + \dots + p^2$$

$$\underline{\text{Step 3}} \quad SST = \sum x_{ij}^2 - CF$$

$$\underline{\text{Step 4}} \quad \frac{\sum C_i^2}{n} = \frac{(a^2 + e^2 + i^2 + m^2)}{4} + \dots + \frac{(d^2 + h^2 + l^2 + p^2)}{4}$$

Step 5 -

$$SSTR = \frac{\sum c_i^2}{n} - CF.$$

$$K-1 \Rightarrow df_1$$

$$n-K \Rightarrow df_2$$

in  
Anova  
table

$$\& SSE = SST - SSTR$$

Anova table [Step 6]

| source of variance | df  | sum of squares | mean square               | F ratio                |
|--------------------|-----|----------------|---------------------------|------------------------|
| B/w groups         | K-1 | SSTR           | $\frac{SSTR}{K-1} = MSTR$ |                        |
| within groups      | n-K | SSE            | $\frac{SSE}{n-K} = MSE$   | $\frac{MSTR}{MSE} = F$ |
| tot.               | n-1 | (SSTR+SSE)     |                           | $F(\alpha, K-1, n-K)$  |

n  $\Rightarrow$  total values from a to P

k  $\Rightarrow$  No. of groups .

conclusion if calculated  $>$  table value  
(F)  $\rightarrow H_0:$  Reject

## 2-Way ANOVA -

|   | Grp1 | Grp2 | G3 | G4 |
|---|------|------|----|----|
| 1 | a    | b    | c  | d  |
| 2 | e    | f    | g  | h  |
| 3 | i    | J    | K  | L  |

Step 1

$$G = a + b + c + d + \dots + L$$

$$\frac{G^2}{n} = CF$$

Step 2

$$\sum \sum x_{ij} = a^2 + b^2 + c^2 + d^2 + \dots + L^2$$

$$SST = \sum \sum x_{ij}^2 - CF$$

Step 3 :-

$$\frac{\sum C_i^2}{n_i} = \frac{(a+e+i)^2}{n} + \frac{(b+f+J)^2}{n} + \dots + \frac{(d+h+l)^2}{n}$$

Here  $n=3$

$$SSTR = \frac{\sum C_i^2}{n_i} - CF$$

Step 4 :-

SSBL for Row

$$\Rightarrow \frac{\sum R_i^2}{n} - CF = SBR$$

$$\Rightarrow \frac{(a+b+c+d)^2}{n} + \frac{(e+f+g+h)^2}{n} + \frac{(i+j+k+l)^2}{n} = \sum R_i^2$$

$n=4$  here.

$$SSBL = \sum R_i^2 - CF$$

Step 5:-

$$SSE = SST - SSTR - SSBL$$

Step 6: table

| Source of variance | df           | sum of squares | mean square                    | F ratio  |
|--------------------|--------------|----------------|--------------------------------|--|
| b/w treatments     | $k-1$        | SSTR           | $MSTR = \frac{SSTR}{(k-1)}$    | for treatment<br>$F = \frac{MSTR}{MSE}$<br>$F(\alpha, \text{cal value})$ |
| b/w blocks         | $m-1$        | SSBL           | $MSBL = \frac{SSBL}{m-1}$      | for Block<br>$F = \frac{MSBL}{MSE}$<br>$F(\alpha, \text{cal value})$     |
| within Groups      | $(k-1)(m-1)$ | SSE            | $MSE = \frac{SSE}{(k-1)(m-1)}$ |  |

$k=4 \Rightarrow 4$  treatments  $\Rightarrow G_1, G_2, G_3, G_4 \}$

$m=3 \Rightarrow 3$  Blocks  $\Rightarrow 1, 2, 3$

in this example.

If  $F_{\text{cal}} < F_{\text{table}}$

Accept  $H_0$   
Reject  $H_1$

$F_{\text{treatment}} \Rightarrow F(\alpha, (k-1), (k-1)(m-1))$

$F_{\text{Block}} \Rightarrow F(\alpha, (m-1), (k-1)(m-1))$

## Maximum likelihood Estimation

Normal

$$f(x|\mu) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$l(\mu) = -\frac{1}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2$$

$$\frac{\partial l}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^N (x_i - \mu) = 0$$

$\hat{\sigma}^2 = \frac{\sum (x_i - \bar{x})^2}{n}$

Binomial

$$L_p = p^k (1-p)^{n-k}$$

log

$$\ell_p = k \log p + n-k \log(1-p)$$

$$\frac{\partial \ell_p}{\partial p} = \frac{k}{p} - \frac{n-k}{1-p} = 0$$

Binomial  $\theta = k/n$

Poisson

$$\frac{\partial L}{\partial \mu} = -n + \frac{n-x}{\mu}$$

$$\hat{\mu} = \bar{x}$$

$$\bar{x} = \frac{\sum x_i}{n}$$

Least square

$$J(b_0, b_1) = \sum_{i=1}^n \left[ y_i - (b_0 + b_1 x_i) \right]^2$$

$\partial J / \partial b_1 = 0$   
no gradient so  $x^2$

# Line of Regression / Linear regression

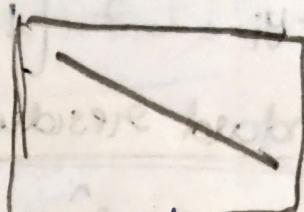
| x | 4 | $x^2$        | $y^2$        | xy          |
|---|---|--------------|--------------|-------------|
| 1 | 5 | 1            | 25           | 5           |
| 2 | 6 | 4            | 36           | 12          |
| 3 | 7 | 9            | 49           | 21          |
| 4 | 8 | 16           | 64           | 32          |
|   |   | $\Sigma x^2$ | $\Sigma y^2$ | $\Sigma xy$ |

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

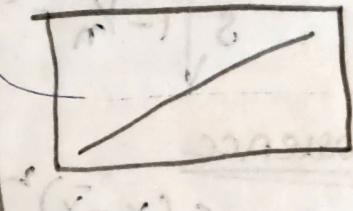
$$\hat{\beta}_1 = \frac{\hat{S}_{xy}}{\hat{S}_{xx}} = \frac{\sum x_i y_i - \frac{\sum x_i \cdot \sum y_i}{n}}{\sum x_i^2 - \frac{(\sum x_i)^2}{n}}$$

$$\hat{\beta}_0 = \frac{\sum y_i}{n} - \hat{\beta}_1 \cdot \frac{\sum x_i}{n}$$

If  $\beta_1 < 0$  diag:



If  $\beta_1 > 0$  diag:



## Sum of sq. of error / Residuals

$$SSE = \sum y_i^2 - \hat{\beta}_0 \sum y_i - \hat{\beta}_1 \sum x_i y_i$$

$$\sigma^2 = S^2 = \frac{SSE}{n-2}$$

## Total sum of squares

$$SST = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum y_i^2 - \frac{(\sum y_i)^2}{n}$$

$$SST = \sum y_i^2 - \frac{(\sum y_i)^2}{n}$$

## Co-efficient of determination

★ ★ ★

$$r^2 = 1 - \frac{SSE}{SST} \quad (\text{blw } 0 \& 1)$$

(or)

$$r^2 = \frac{SSRg}{SST}$$

$$\hat{y}_i - \bar{y}_m = \sigma^2 \cdot \left[ i - \bar{y}_m - \frac{(x_i - \bar{x})^2}{Sxx} \right]$$

at least ≥

★ ★

## Standard residual:-

$$e_i^o = \frac{\hat{y}_i - \bar{y}_m}{\sqrt{i - \bar{y}_m - \frac{(x_i - \bar{x})^2}{Sxx}}}$$

### Covariance

$$\text{var}_x = \frac{\sum (x - \bar{x})^2}{N-1} = \frac{\sum (x - \bar{x}) \cdot (x - \bar{x})}{N-1}$$

$$\text{cov}_{xy} = \frac{\sum (x - \bar{x})(y - \bar{y})}{N-1}$$

## formulae for finding correlation

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

## Fitting a trend to time series - (fit a linear trend)

for odd num:

$$x = \frac{t - \text{middle value of } t}{\text{Interval}(h)}$$

for even no's

$$x = \frac{t - (\text{mean of middle 2 values of } t)}{\frac{1}{2} \text{Interval}(h)}$$

$$\hat{b} = \frac{\sum xy_t}{\sum x^2}$$

$$\hat{a} = \frac{\sum y_t}{n}$$

$$\hat{y}_t = \hat{a} + \hat{b}x$$

| Year<br>t | Production<br>$y_t$ | $x = t - \frac{\text{value}}{\text{Interval}}$ | $x^2$    | $xy_t$                 |
|-----------|---------------------|--|----------|------------------------|
| -         | -                   | -  | -        | -                      |
| -         | -                   | -  | -        | -                      |
| -         | -                   | -  | -        | -                      |
| -         | -                   | -  | -        | -                      |
|           |                     | $\sum y_t$                                     | $\sum x$ | $\sum x^2$ $\sum xy_t$ |

## Exponential smoothing technique :-

L compute moving avg)

(simple exponential smoothing with  $\alpha$ )

$$F_1 = F_2 = y$$

$$F_n = \alpha \cdot y_{n-1} + (1-\alpha) F_{n-1}$$

## Performance measures

$$MSE = \frac{\sum (\Delta t)^2}{n}$$

$$MAD = \frac{\sum |\Delta t|}{n}$$

$$MAPE = \frac{\sum \Delta t}{yt} * 100$$

$$LAD = \max |\Delta t|$$

$$|\Delta t| = |y_t - P_t|$$

### Gmm

|           | RB | NT |
|-----------|----|----|
| literally | 20 | 5  |
| UD 1      | 20 | 5  |
| UD 2      | 30 | 3  |
| UD 3      | 50 | 7  |

cluster 1:  $(C_{20,1}) (10,1)$   
 cluster 2:  $((50,7)(10,1), 0, \rho)$   
 module of  $\mu$

Ans

### cluster 1

$$\mu_1 = 20$$

$$\mu_2 = 1$$

$$\sigma_1^2 = 10$$

$$\sigma_2^2 = 1$$

$$\pi_1 = 20/5$$

$$= 4/5$$

### cluster 2

$$\mu_1 = 50$$

$$\sigma_1^2 = 10$$

$$\pi_2 = 20/5$$

$$\mu_2 = 7 \quad \sigma_2^2 = 1$$

$$N(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[ -\frac{(x-\mu)^2}{2\sigma^2} \right]$$

cluster 1:  $(20, 5)$

for  $D_1 \Rightarrow (20, 5)$  for cluster 1

feature 1 ( $N(\mu_1, \sigma_1)$ )

$$= N(20, 20, 10)$$

$$g_1^2 = \frac{1}{2\pi(10)} \exp \left[ -\frac{(20-20)^2}{2(10)} \right] = 1.4367198 \times 10^{-1}$$

feature 2 for 201 (20,5)

$$t \sim N(x_2, \mu_2, \sigma_2)$$

$$z \sim N(0, 1, 1)$$

$$\rho_2 = \frac{1}{2\pi(1)} \cdot \exp \left[ -\frac{(x_2 - \mu_2)^2}{2\sigma_2^2} \right] = 1.33830226 \times 10^{-4}$$

Product  $\rightarrow$  Gaussian 1  $\times$  Gaussian 2

$$= 1.4687 \times 10^{-7} \times 1.33830226 \times 10^{-4} = 1.9296807 \times 10^{-11}$$

$$\text{Weight} = \pi \cdot P$$

$$= 0.5 (\text{product})$$

$$= 0.5 (1.9296807 \times 10^{-11}) = 9.647512 \times 10^{-12}$$

$$V_{11} = \frac{\text{Weight}_1}{\text{Weight}_1 + \text{Weight}_2}$$

$$V_{12} = \frac{\text{Weight}_2}{\text{Weight}_1 + \text{Weight}_2}$$

Same way

cluster 2 for feature 1 & feature 2

then next point [30/3]

do cluster 1  $\rightarrow$  feature 1  
 $\rightarrow$  feature 2

do cluster 2  $\rightarrow$  feature 1  
 $\rightarrow$  feature 2

Spearman

(2.06)

rank of correlation coeff

$$\rho R = 1 - \frac{6 \sum d_i^2}{n(n^2-1)}$$

Spearman