

Date: 11-05-2022

## Poisson Distribution

The discrete prob. + that expresses the prob. of a given number of events occurring in a fixed interval of time or space if those events occur with a known constant average rate and independently of the time since the last event.

or

It is used to show how many times an event is likely to occur over a specified period of time.

Formula:

$$P(X, \lambda) = \frac{e^{-\lambda} \lambda^x}{x!}, x = 0, 1, 2, 3, \dots$$

$\lambda$  = average rate of change  $x!$

Application:

if  $n > 100$ ,  $p < 0.1 \rightarrow$  Poisson Dist.  
when, "n" is large & "p" is close to zero.

The limiting form of Binomial Distribution.

Note:

We use Poisson dist. as a limiting form of Binomial Dist. when  $n$  is large ( $n > 100$ ) and  $p$  is close to zero ( $p < 0.1$ )  
 $\rightarrow$  Poisson dist. is used for rare events.

$$\text{Mean} = \text{Variance} = \lambda$$

Sunny®

Date: .....

Example:

If approximately 2% of the people in a room of 200 people are left handed, find the prob. that exactly 5 people are left handed???

Solution: Given that =

$$p = 0.02$$

$$n = 200$$

$$P(X=5) = ?$$

$\mu = np$  Mean of Binomial Distribution.

$$\lambda = np$$

$$\lambda = (200)(0.02) \Rightarrow \lambda = 4$$

$$P(X=5) = \frac{e^{-4} (4)^5}{5!} = 0.156 \approx 0.16 = 16\%$$

~~Q2~~

If the no. of accidents occurring in an industrial area during a day following a poisson r.v with parameter 3. Find the prob. that on random day

(1) No accident will occur

(2) At most two accidents will occur.

Solutions:

$$\lambda = 3$$

i)  $P(X=0) = \frac{e^{-3} (3)^0}{0!}$

(2)  $P(X \leq 2) = P(X=0) + P(X=1) + P(X=2)$

Sunny®

Date:

$$P(X=1) = \frac{e^{-3} (3)^1}{1!}$$

$$P(X=2) = \frac{e^{-3} (3)^2}{2!}$$

(iii) Prob. for at least 2 accidents

$$P(X \geq 2) = 1 - (P(X < 2)) \\ = 1 - [P(X=0) + P(X=1)].$$

Properties of Poisson Dist.:

→ Events are independent

→ The average Number of successes in given period of time alone can occur.

→ No two events can occur at the same time.

"Continuous probability Distribution"

→ Uniform dist.

→ Exponential dist.

→ Normal dist. "most practical life"

S.D.

$$P(X, \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}$$

"Range for Z"

Mean

$-\infty < x < +\infty$

$-\infty < \mu < +\infty$

$Z = \frac{x-\mu}{\sigma}$  ("Standard Normal r.v")

Z-score

Standard form.

Sunny®

Date:

## "R-language"

### Data Entry with Vector

Case sensitive

$$X = c(1, 2, 3, 4, 5, \dots, 30)$$

→ mean(x)

To Add required package

→ median(x)

→ install a package & load file

→ sd(x)

plot(x) → plotting pie chart 3D

→ Quartile(x, 0.75)

→ pie(x) # construct the pie chart of data.

→ hist(x) # construct the histogram of data.

→ plot(x) # construct the scattered diagram.

### Binomial Dist

→ Random number generation "rbinom(num, n, p)"

→ prob computation dbinom(x, n, p)

→ Cumulative prob computation pbinom(2, 5, 0, 1)

$$P(X \geq 2) = 1 - P(X < 2)$$

### Poisson Distribution

→ rpois(x, λ)

dpois(x, λ) ↓ Mean probability calculation

How many numbers?

### Normal Distribution:-

→ rnorm(x, mean, sd)

### Normal Distribution

$$P(X < 2) = \Phi(2)$$

$$P(2 < X < 3) = \Phi(3) - \Phi(2)$$

$$\therefore P(X > 2) = 1 - \Phi(2)$$

Suppose that plant height is distributed normally with average height 75 cm and variance of Sunny®

Date: .....

25 cm. What is the probability that a plant has

$$\text{Avg. height} = 75 \text{ cm}$$

$$\text{Variance height} = 25 \text{ cm.}$$

$$\text{height} < 70 \text{ cm} ???$$

(2) What is the probability that a plant has more than 85 cm ???

(3) Prob. for between 70 and 85 ?

Solution: Let  $X$  be the height of plant

$$P(X < 70) = \frac{70 - 75}{\sqrt{25}} = \frac{-5}{5} = -1$$

$$S.D. = \sqrt{\text{Variance}}$$

$$\phi(-a) = 1 - \phi(a)$$

$$= 1 - 0.8413 \Rightarrow 0.1587$$

$$P(X > 85) = 1 - P(X < 85)$$

$$X = 85$$

$$= 1 - P\left(\frac{85 - 75}{5}\right)$$

$$\mu = 75$$

$$\sigma = 5$$

$$= 1 - \phi(2) \Rightarrow 1 - 0.9772 = 0.0228$$

$$= 0.0228$$

$$P(70 < X < 80) = ? \Rightarrow P(80) - P(70)$$

$$= P\left(\frac{80 - 75}{5}\right) - P\left(\frac{70 - 75}{5}\right)$$

$$= \phi(1) - \phi(-1)$$

$$= \phi(1) - [1 - \phi(1)]$$

$$= 0.8413 - [1 - 0.8413]$$

$$= 0.8413 - 0.1587$$

$$= 0.6826 \Rightarrow 68.26\%$$

Sunny®

Date: .....

i)  $P(X < a) = P\left\{ \frac{X - \mu}{\sigma} < \frac{a - \mu}{\sigma} \right\}$   
 $= \Phi\left(\frac{a - \mu}{\sigma}\right)$

ii)  $P(X < a) = \Phi(-a) = 1 - \Phi(a)$

Q: A certain type of storage batteries lasts, on average 3 hrs. with the standard deviation of 0.5 hrs. Assuming that the battery lives are normally distributed. Find the probability that the given battery will last for less than 2.3 hrs ???

$$P(X < 2.3) = ? \quad X = 2.3$$

$$\mu = 3$$

$$= P\left\{ \frac{X - \mu}{\sigma} \right\} \quad \sigma = 0.5$$

$$= P\left\{ \frac{2.3 - 3}{0.5} \right\} \Rightarrow \Phi(-1.4) \Rightarrow 1 - \Phi(1.4)$$

$$= 1 - 0.9192 = 0.0808 \Rightarrow 8.08\%$$

There is 8% chance that a given battery will last for 2.3 hrs.

Q: An electric bulb manufacturing company claims that average life of bulbs is 2 years with the s.d. of 0.25 years. Assuming that electric bulbs manufacturing is normally distributed. Find the prob. that a bulb burns between 2 & 3 years ???

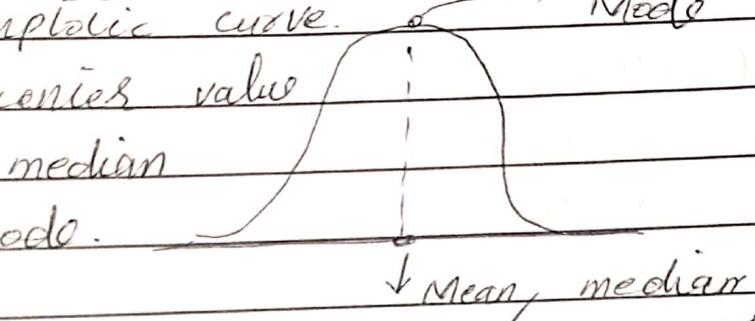
Sunny®

Date: \_\_\_\_\_

$$\begin{aligned} P(2 < X < 3) &=? \Rightarrow P(X < 3) - P(X > 3) \\ &= P\left[\frac{3-2}{0.25}\right] - P\left[\frac{2-2}{0.25}\right] \\ &= \phi(4) - \phi(0) \Rightarrow 0.9998 - 0.5000 \\ &= 0.4998 \Rightarrow 49.98\% \end{aligned}$$

### Characteristics:

- Area under the normal curve is 1
- Asymptotic curve. → Mode
- The center value is mean, median and mode.



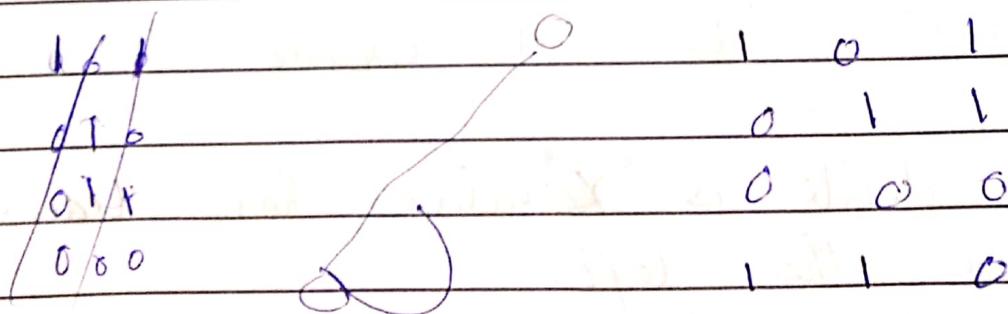
→ "The squeeziness of curve depends on variation, how much it goes far away from center?"

$$\mu \pm 1 = 68.26\%$$

$$\mu \pm 2 = 95.45\%, \mu \pm 3 = 99.73\%$$

↓  
S.D.

visited



Sum

Date: .....

## Normal Distribution Inverse Table

Q:

The average grade for an exam is 74 and the sd. is 7. If 12% ~~percent~~ of the class are given As and the grades are curved to follow a normal distribution, what is the best possible "A" and the highest possible "B"?

$$\mu = 74$$

$$s.d. = 7$$

12% of the students attained grade "A".

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

$$\frac{x - 74}{7}$$

$$x = \frac{x - 74}{7} + 74$$

"Area under the curve logic works here"

$$1 - 0.12 \Rightarrow 0.88$$

Solving eq. of

$$x = (0.88)(7) + 74$$

$$x = 82.18$$

$$x = 82.18$$

Round down to get "B" Grade

Round up to get "A" Grade

$x = 82$  for B Grade

$x = 83$  for A Grade

Q:

Six decile is "X" value that leaves 60% of the area to the left.

$$\text{Six decile} = 0.26$$

Sunny

Date: ..... Interpretation

$X = 70 + \mu$  Almost 76% of  
 $= 74 + 7(0.26)$  the students will get  
 $X = 75.82$  60% marks in this  
 population.

Q:

The IQs of 600 applicants of a certain college are approximately normally distributed, with mean of 115 and sd. 12. If the college requires the IQ of 95%. How many of these students will be rejected on this basis regardless of their other qualifications?

Sol:  $\mu = 115$  Let "x" be the random  $\sigma = 12$  number which indicates the

$$X = (1.65)(12) + 115 \quad \text{IQ of students.}$$

$$X = 134.8 \approx 135 \quad 0.95 \Rightarrow z = 1.65$$

So, almost 136 students will be rejected out of 600 applicants.

### "Regression & Correlation Analysis"

Regression:

The dependence of one variable upon another variable(s)

dependent variable = Y "Response variable"

independent variable = X "Explanatory"

Examples:

Marks of a std. depends on study hours.  
 wheat production depends upon seed Quality.

Speed of laptop depends upon RAM.

Sunny®

Date:-

## Simple Regression (One dependent & One independent Variable)

One dependent variable and more than one independent variables  $\rightarrow$  "Multiple Regression."

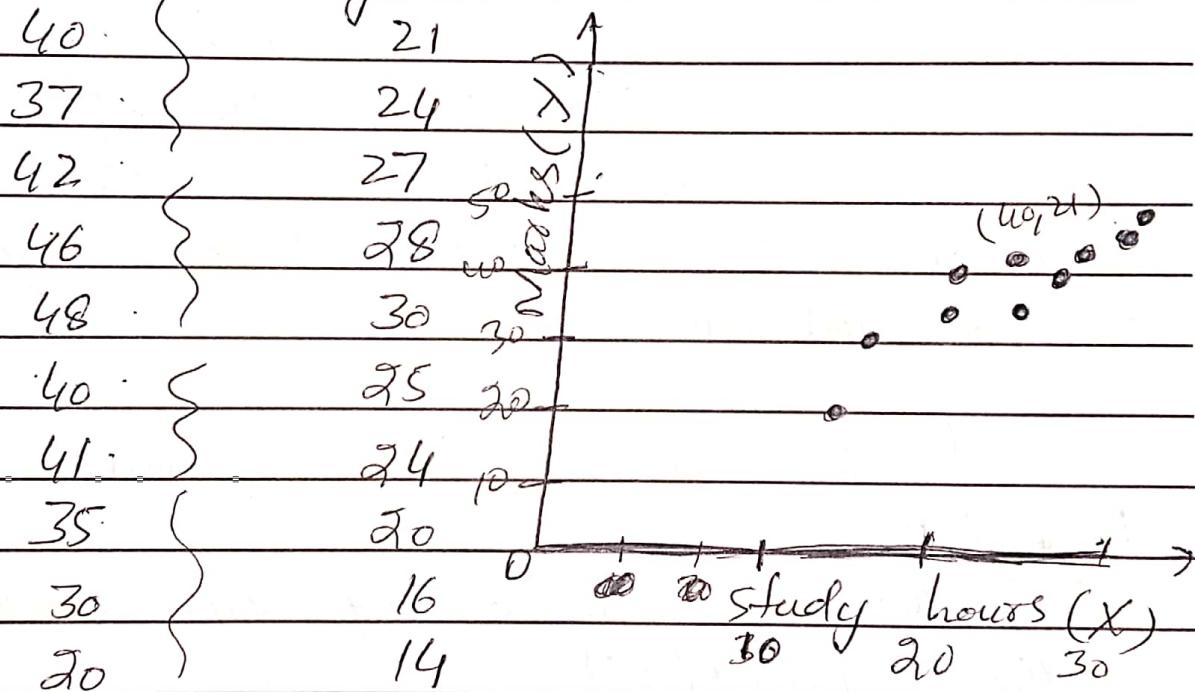
## Multi-Variable Regression

Morse dependency & Independence vars.

# Simple linear Regression

linear      straight line relationship.

marks(y)    Study hours(x)



## Scatter Diagram

## Two dimensional Graphical representation

of data, where we take independent var. on X-axis and dependent var. on Y-axis.

Sunny®

Date: .....

linear = Rate of change is const.

By linear we mean that the rate of change is ~~not~~ const.

### Inference:

As we could see from isend that something is increasing but we can't Quantify that factor.

There is positive and direct relationship between dependent var. and independent var. we are unable to identify how much relationship is present???

Regression Analysis

$$y = \alpha + \beta x + \epsilon$$

PRF  $\epsilon$

Dependent Var

y-intercept

Error term

"population regression function"

Slope / regression co-efficient

Independent Var

$\epsilon$ : The effect of those variables which are not included in the model are studied by the error term.

The survey of whole population is not possible in real or difficult.

Summary

Date: \_\_\_\_\_

estimate of actual  $y$

$$SRF: \hat{y} = a + bx \quad E(\epsilon) = 0$$

$a$  = Estimate of  $\alpha$

$b$  = Estimate of  $\beta$

$$\text{mean of } y \quad a = \bar{y} - b\bar{x} \quad \text{mean of } x$$

$$b = \frac{n \sum XY - \sum X \sum Y}{n \sum x^2 - (\sum X)^2}$$

Ordinary least squares

"Methods of

likelihood estimation

$\hat{\epsilon}^2$

(1) Choose the dependent var. ?

(2) Draw the scattered Diags to check the relationship between variables?

(3) perform regression analysis?

(4) Interpret the regression co-efficient?

(5) Compare the results of part 2 & 4?

(6) Check the model adequacy?

(7) Check the interdependence of variables?

→  $E(3B)$  →

$$a = \bar{y} - b\bar{x}$$

$$b = \frac{n \sum XY - \sum X \sum Y}{n \sum x^2 - (\sum X)^2}$$

Sunny®

Date: \_\_\_\_\_ A \_\_\_\_\_

$\hat{Y}$	$X$	$XY$	$X^2$	$\bar{Y}$
4.92	40	840	1600	35.076
-2.534	37	888	1369	39.534
1.992	42	1134	1764	43.992
0.522	46	1288	2116	these are
-0.45	48	1440	2304	wrong
-1.02	40	1000	1600	41.02 calculations
1.466	41	984	1681	39.534
1.41	35	700	1225	33.59
2.354	30	480	900	-27.646
-4.674	20	280	400	24.674
				378.9/-
	379	9034	14559	5483/-

$$b = \frac{10(9034) - (229)(379)}{10(5483) - (229)^2}$$

$$b = 1.486$$

$$q = \bar{y} - b\bar{x}$$

$$= \frac{\sum y}{n} - b \frac{\sum x}{n}$$

$$= \frac{379}{10} - (1.186) \frac{229}{10}$$

$$a = 3.87$$

Sunny

Date: .....

$$\hat{Y} = 3.87 + 1.486(X)$$

increase  $\beta$  +ve  
decrease  $\beta$  -ve

— (S4B) —

### Interpretation:

On the average the value of "Y" (Marks) will change by a unit increase in "X" (Study hrs)  
OR

For a unit increase in "X" (Study hrs), on the average the "Y" (Marks) will increase by 1.486.

— (S5B) —

The results of scattered diagram and regression co-efficients indicates that there is +ve relationship between marks & study hrs.

### Correlation Analysis

It is the interdependence of variable. It is denoted by  $r$ . Its value lies between -1 & 1

$$r = \frac{n \sum XY - \sum X \sum Y}{\sqrt{[n \sum X^2 - (\sum X)^2] [n \sum Y^2 - (\sum Y)^2]}}$$

$$r = \frac{10(9034) - (379)(229)}{\sqrt{[(10(54.83) - (229)) (10(14959) - (379))]}$$

$$= \frac{3549}{18900.44676}$$

$$\boxed{r = 0.1877}$$

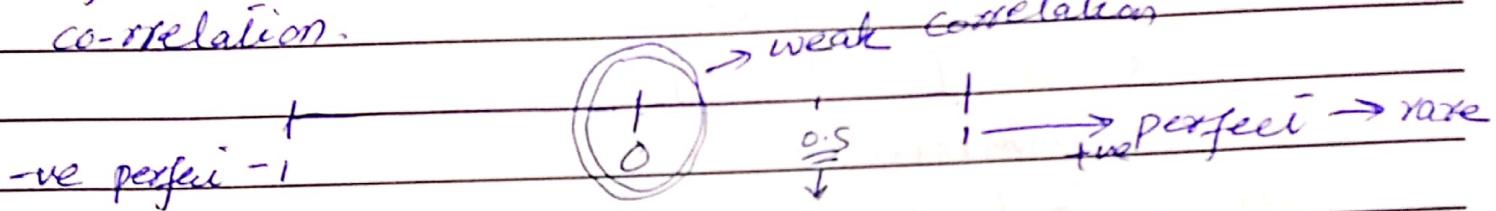
Sunny®

Date: .....

## Interpretation:

Since, the value of  $r$  is very close of +1, it indicates that there is a strong +ve correlation between variables (marks & s.h.)

→ If value of  $r$  is zero, then there is no correlation.



Model adequacy =  $(\text{correlation})^2$

$$r^2 = 0.8862$$

$$r^2 = 88.62\%$$

e.g., if  $r^2 = 0.8862$ , then  $r = \pm 0.937$ .

Model adequacy is computed by taking square of correlation co-efficient " $r$ ".

By model adequacy we mean that how much variation in the model is explained by the included X-variable.

Here, study hours explains 89% of the variation in marks. It means that Study hour is very important variable.

Sunny

Date: .....

# "Sampling Techniques"

Sampling  
Techniques  
from  
a lot.

## Hypothesis Testing

### Constructive Steps:

(1) Null hypothesis:  $H_0$

Alternative hypothesis:  $H_A$  or  $H_1$

(2) Level of significance:  $\alpha$

(3) Test statistic

(4) Critical Region

(5) Computation

(6) Conclusion.

### Example:

#### Case 1

$$\mu = 10.2 \text{ kg } ?$$

$$\sigma = 0.6 \text{ kg}$$

#### Case 2

$$\mu = 12 \text{ kg}$$

$$\sigma = 1 \text{ kg } X$$

#### Case 3

$$\mu = 9 \text{ kg }$$

$$\sigma = 0.3 \text{ kg } X$$

Null hypothesis: "jis cheage ko test kena"

A hypothesis to be tested for possible rejection. Claim

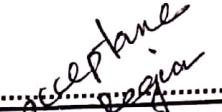
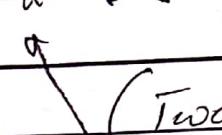
Alternate hypothesis:  $\sim H_0$  "Equal kabhi ni agega"

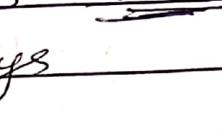
A hypothesis which is accepted when

Null hypothesis is rejected

Sunny®

Date: \_\_\_\_\_

e.g. (1)  $H_0: \mu = 2$  inch   
 $H_1: \mu \neq 2$  inch 

(2)  $H_0: \mu = 3$  days almost   
 $H_1: \mu > 3$  days 

(3)  $H_0: \mu = 250$  ml   
 $H_1: \mu < 250$  ml 

Note:

The statements like "equal", "at least", "almost" will help us to use this statement in  $H_0$ .  
 → "Always equal" will be using in  $H_0$ .  
 → The statements like less than, "greater than", "more than", "shorter", "larger", "superior", "inferior" will be used in  $H_1$ .

level of Significance:  $\alpha$

The probability of committing TYPE-I Errors. Generally we use  $\alpha = 1\%, 5\%, 10\%$

Type-I Errors

Rejection of true  $H_0$

Type-II Errors

Acceptance of false  $H_0$

If not mentioned  $\alpha$  will be 5%.

Sunny®

Date.....

Test Statistics:

$$Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} \quad (\bar{x} - \text{test}) = \frac{\text{Error}}{\text{S.E.}}$$

$\frac{\sigma}{\sqrt{n}} \rightarrow \text{standard error.}$

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \quad (\bar{x} - \text{test})$$

$\sigma$  : population standard deviation

$s$  : sample standard deviation.

Question:

A manufacturing comp. claims that on the average burning hrs of their tube lights is 2000 hrs. with the s.d. of 100 hrs. A random sample of 100 bulbs was taken. On the average they burn for 1950 hrs with the s.d. of 90 hrs.

Assume that the burning hrs follows to the normal distribution. Test whether the functionality time has reduced with  $\alpha = 1\%$ .

Solution:

$$\text{H}_0: \mu = 2000 \text{ hrs.}$$

$$\text{H}_1: \mu < 2000 \text{ hrs.}$$

$\neq$

$$\text{Step 1: } \alpha = 1\%.$$

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

Step-3

Sunny®

$\mu$  = population mean  
 $\bar{x}$  = sample mean.

Date: .....

testing always for population.

Question: A medical company claims that their medicine reduces the recovery time of a disease. A random sample of 50 patients were selected randomly and on the average their recovery time is 7 days with the s.d of 1 days. Test whether the average recovery time  $\mu = 8$  days is true or not. with  $\alpha = 1\%$  ???

Solution:

Step-1

$$H_0: \mu = 8 \text{ days}$$

$$H_1: \mu < 8 \text{ days}$$

Step-2

$$\alpha = 1\%$$

Step-3

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

Continue Q: 1

Step-4

$$H_1: \mu < 2 \text{ days}$$

Case-1 ✓

$$t_c < t_b(0.01, 99)$$

$$t_c < t_\alpha, \text{d.f.}$$

$$t_b = 2$$

$$H_1: \mu \neq 2$$

Case-2

$$|t_c| \geq t_\alpha \frac{2}{\alpha}, \text{dof.}$$

$$H_1: \mu > 2$$

Case-3

$$t_c > t_\alpha, \text{dof.}$$

Degrees of freedom:

$$\text{d.f.} = n-1 \quad \text{for t-test}$$

↓ No. of observations

Sunny

Date: .....

### Step-5

$n = 100 \rightarrow$  "Sample Size"

$\mu = 2000 \rightarrow$  Mean population

$\bar{x} = 1950 \rightarrow$  Mean Sample

$s = 90$

$$t_c = \frac{1950 - 2000}{90/\sqrt{100}} \Rightarrow -5.0$$

$$\boxed{t_c = -5.0} \Rightarrow -5.5556$$

### Step-6

#### Conclusion:

As  $t_c = -5.6$  and  $t_{table} = 2$  which falls in critical region. Hence, we accept  $H_0$ . Hence, we conclude that the burning hours of a bulb has reduced.

$$\bar{x} \pm t_{\alpha/2, d.f.} \times \frac{s}{\sqrt{n}} \quad (\text{Confidence interval})$$

$$= 1950 \pm (2) \times \frac{90}{\sqrt{100}}$$

Hypothesis testing only tells about acceptance and

$$= 1950 \pm 2(9)$$

rejection. C.I

$$= 1950 \pm 18$$

tell the limit

Q: A random sample of 100 recorded deaths in US during past year showed an average life span of 71.8 years. Assuming a population s.d. of 8.9 years does this seem to indicate that the **Sunny**®

Date.....

mean life span today is greater than 70 yrs?  
Use a 0.05 level of significance.

Solution:

Step-1  $H_0: \mu = 70$

Step-2  $H_1: \mu > 70$

Step-3  $\alpha = 0.05$

Step-4  $t_c > t_{\alpha, d.f}$

$$t_c > t(0.05, 99) = 1.67$$

Step-5

$$t_c = \frac{\bar{x} - \mu}{S/\sqrt{n}}$$

$$= \frac{71.8 - 70}{8.9/\sqrt{100}} = 2.02$$

Step-6

$$t_c > t(0.05, 99)$$

$$2.02 > 1.67$$

So,  $H_1$  stays true.

Q:

Step-1  $H_0: \mu = 46 \text{ kw}$

Step-2  $H_1: \mu > 46 \text{ kw}$

Step-3  $\alpha = 0.05$

$$t_c > t(\alpha, d.f)$$

Sum

Date: \_\_\_\_\_

## Regression Models & Correlation Models using R-language

>  $x_1 = c(7, 9, 11, 13, 15)$

>  $x_2 = c(15, 18, 19, 20, 21, 22)$

>  $y = c(40, 30, 25, 20, 30, 50)$

Formulae:

>  $lm(y \sim x_1)$  [for simple regression]

>  $lm(y \sim x_1 + x_2)$  [for multiple regression]

>  $\text{cor}(x_1, x_2)$  [for correlation]

## Comparison Of Covariance & Correlation.

$$\gamma = n \sum xy - \sum x \sum y$$

$$\sqrt{[(n \sum x^2 - (\sum x)^2)(n \sum y^2 - (\sum y)^2)]}^{1/2}$$

$$\gamma = \frac{\text{Cov}(x, y)}{\sqrt{V(x)V(y)}}$$

$$\text{Cov}(x, y) = E\{(x - \bar{x})(y - \bar{y})\}$$

$$\text{Cov}(x, x) = E\{(x - \bar{x})(x - \bar{x})\}$$

Covariance:

The dependence of two variables. It tells us that how both variable move together

Sunny®

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

$$E(x) = \sum x p(x)$$

Date: \_\_\_\_\_

Covariance  $(-\infty, \infty)$  ]  $\rightarrow$  Range  
Correlation  $(-1, 1)$  ]

Correlation:

The ratio of covariance of variables to the variance of variables

Degrees of freedom:  $n-1$

The difference of no. of parameter and numbers of observations.

$$d_f = n - \text{number of parameters}$$

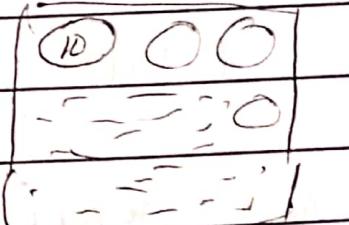
Confidence Interval:

$$\bar{x} \pm t\left(\frac{\alpha}{2}, v\right) \frac{s}{\sqrt{n}}$$

Estimate  $\pm$  Confidence  $\times$  S.E  
population mean limit  
estimation.

95% C.I for population Mean

$$\alpha = 5\%$$



Example:

$$H_0: \mu = 10$$

Since,  $\mu = 10$  falls

$$H_1: \mu > 10$$
 in this interval, so

$$t_c > t_{\alpha/2}$$

we accept  $H_0$ .

$$t = 2.1$$

$$t_b = 2.6$$

$$(9, 10) \rightarrow \text{interval}$$

Sunny®

Date: .....

## Interpretation of Confidence interval:

The given interval indicates that the recovery time is from 9 to 12 days and the company is 95% sure in their statement.

If we take other samples from lot have taken then 95 out of 100 will show the same results.

Sample	population	
$n$	$N$	"Size"
$\bar{x}$	$\mu$	"Mean"
$s$	$\sigma$	"S.d."
$r$		correlation
$b$		slope co-efficient

Standard Normal Range  $-3 \leq z \leq +3$

P-value:

If P-value is very close to zero then reject  $H_0$ . otherwise accept  $H_0$ .

$$> pt(\alpha, r)$$

P-value  $\leq \alpha$ . value



Reject  $H_0$

\*\*\* Significance.

Sunny