

SATURDAY

DAY (076-289)

11th Week

17

Appointment

Notes

Work to do

2018  
2018

DATA

ANALYTICS

18 SUNDAY

Data Analyst

S	M	T	W	T	F	S	S	M	T	W	F	S	S	M	T	W	T	F	S	S	M	T
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23

APR  
'18

APR

MAY

JUN

2018

TUESDAY

DAY (079-286)

12th Week

20

Appointment

Notes

Work to do

## -: Statistics :-

Beyond the Data & Navigating the World of Statistics.

- Introduction to Statistics.
- Descriptive Statistics
- Probability.
- Inferential Statistics

APR

MAY

JUN

S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24

APR  
'18

2018

THURSDAY

DAY (081-284)

12th Week

>> 01

22

Notes

Work to do

## Appointment » Introduction of Statistics :-

- What is Statistics ?
- Types of Statistics - Descriptive and Inferential
- Types of Data
- Population and Sample
- Sampling Techniques
- Statistical Data Analysis Steps

## # Descriptive Statistics

- » Measure of central tendency - Mean, Median and Mode
- » Measure of dispersion - Range, Variance, Standard Deviation, Percentiles and Quantiles
- » Graphical Representations - Boxplots, Histogram, Scatterplot.
- » Outliers and understanding their impact
- » Correlation and Covariance

S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	APR
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	'18

23

FRIDAY  
DAY (082-283)  
12th Week

March

- Appointment Notes Work to do
- ### # Probability
- >> Basic Probability concepts : sample space, event etc.
  - >> Types of Events :
    - (i) Disjoint or Non-disjoint event
    - (ii) Independent or Dependent event.
  - >> Complement of Probability
  - >> Conditional Probability
  - >> Bayes Theorem
  - >> Probability Distributions
  - >> Random Variables and its Types
    - (Discrete & Continuous)
  - >> PMF and PDF
  - >> PMF : probability mass Function
  - >> PDF : " density

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
18	1	2	3	4	5	6	7	8	9	10	11	12	13

SATURDAY  
DAY (083-282)  
12th Week

24

2018

- Appointment Notes Work to do
- ### # Discrete Distributions
- Binomial (binomial)
  - Bernoulli (bernoulli)
- ### # Continuous Distributions
- Uniform
  - Normal
- ### # Standard Normal Distribution
- ### # Standardization
- ### # Normalization
- ### # Empirical Rule

25 SUNDAY

S	M	T	W	T	F	S	S	M	T	W	F	S	S	M
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15

APR  
18

26

MONDAY  
DAY (085-280)  
13th Week

March

Appointment	Notes	Work to do
# Inferential Statistics		

- Relationship with Descriptive Statistics
- Point and Interval Estimation
- Confidence Interval (Z/T distribution)
- Hypothesis testing &
- >> Types of Hypothesis : Null and Alternative
- >> Level of Significance and p value
- >> Types of Errors
- >> One tailed or two tailed test.
- Types of test in statistics ( z test, t test, ANOVA, Chi square etc...)

MAR	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	F	S
18	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23

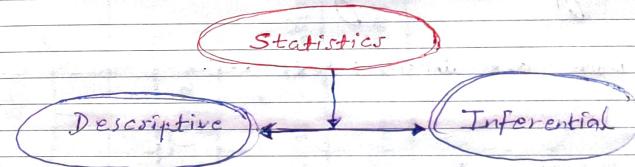
2018

TUESDAY  
DAY (086-279)  
13th Week

27

Appointment	Notes	Work to do
# Statistics		

- # Statistics :- The study and manipulation of data, including ways to gather, review, analyze, and draw conclusion from data.
- \* collecting
  - \* analyzing
  - \* interpreting
  - \* drawing conclusion...



### Type :-

#### i.) Descriptive Statistics :-

- > Collection, Analyzing and Interrelating of data.
- >> Understanding main features of data.
- >> Organizing and Summarizing information from data.
  - graphs
  - tables
  - average
  - variation
  - etc...

MAR	S	M	T	W	F	S	S	M	T	W	F	S	S	M	T	W	F	S	S	M	T	APR	
18	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23

28

WEDNESDAY  
DAY (087-278)  
13th Week

Appointment

Notes

Work to do

ii.) Inferential Statistics :-

- > drawing conclusion from data.  
(representing complete population).
- >> sample ...
- >>> statistical method. (statistical methods) ...
- Confidence Interval
- Estimation
- Hypothesis Testing ...
- >> Both are interrelated ...

# DATA :-

- > collection of facts, observations or measurements used for analysis.
- >> It can be ;
  - Numerical
  - Categorical
  - Combination of both ...

MAR	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M					
18	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31

March

2018

THURSDAY  
DAY (088-277)  
13th Week

29

Appointment

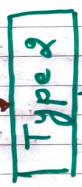
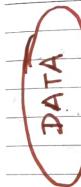
Notes

Work to do



&gt;&gt; Univariate

&gt;&gt; Multivariate



&gt;&gt; Cross-Sectional

&gt;&gt; Time Series



• Structured

• Un-structured

APR

MAY

JUN

APR	S	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M									
18	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	18

30

FRIDAY  
DAY (089-276)  
13th Week

March

## \* Structured Data :-

- > Data having structure.
- >> Data organized in the form of row of columns.

Example :-

- \* Tables
- \* Spreadsheet
- \* relational database etc...

## \* Unstructured Data :-

- >> Data not having a particular structure.

Example :-

- \* multimedia content (image, audio...)
- \* text (emails, blog etc...)
- \* web page

## \* Cross Sectional Data :-

- >> Data collected at a single point of time.

## \* Time Series Data :-

- >> Collecting data over a sequence of time intervals.

Eg :-

- \* daily stock prices
- \* monthly sales data

MAR	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F
18	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16

SATURDAY  
DAY (090-275)  
13th Week

31

2018

## \* Univariate Data :-

- >> Single variable (single data collection).

## \* Multivariate Data :-

- >> Having two or more than two variables in data set.

## # Variables :-

- >> Variables may also called a data item.

Example :-

- \* age
- \* sex
- \* country of birth etc...

## - types of Variables :-

- |                  |                                     |
|------------------|-------------------------------------|
| (i) Nominal      | : gender, colour                    |
| (ii) Ordinal     | : education levels, customer rating |
| (iii) Numerical  | : integer, age, price               |
| (iv) Categorical | : car, product categories           |
| (v) Interval     | : temperature, IQ scores            |
| (vi) Ratio       | : height, weight                    |

S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16

APR

MAY

JUN

18

01

SUNDAY

DAY (091-274)

13th Week

April

Appointment

Notes

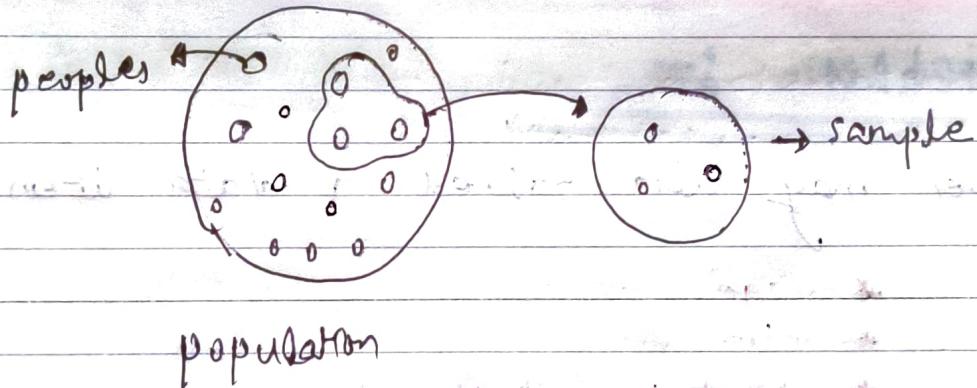
Work to do

## # Population and Sample :-

» population is the entire group of individuals.

example :- people in India ;  
all users on social media

» Sample is a subset of population.



### \* Why Sample :-

- » To reduce the cost of data collection.
- » When a full census data cannot be taken.

APR '18	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25

02

MONDAY  
DAY (092-273)  
14th Week

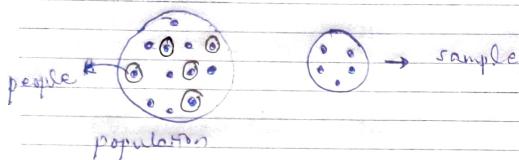
Appointment

Notes

Work to do

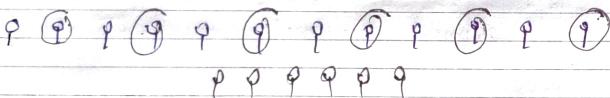
## # Sampling Techniques :-

i) Random :-



ii) Stratified :- based on characteristics.  
(gender).

iii) Systematic :- system follows. ( $k^{\text{th}}$  element)  
(number wise).



iv) Clustered :-

» population divided in cluster.

» randomly select complete cluster.

April

2018

03

TUESDAY  
DAY (093-272)  
14th Week

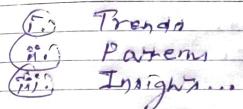
Appointment

Notes

Work to do

## # Statistical Analysis :-

» The process of collecting large volumes of data and then using statistics and other data analysis techniques to identify :-



\* Steps follow while performing :-

- i) Define the problem or research question?
- ii) Data Collection ?
- iii) Data Cleaning ...
- iv) Exploratory Data Analysis.
- v) Data Transformation ?
- vi) Hypothesis Formulation
- vii) Statistical Testing.
- viii) Interpretation of Results ...
- ix) Draw Conclusions ?
- x) Document the Analysis Process / Report Making ...

APR	S	M	T	W	F	S	S	M	T	W	F	S	S	M	T	W
18	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16

T	W	F	S	S	M	T	W	F	S	S	M	T	W	F	S	S	M	T	W
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20

04

WEDNESDAY

DAY (094-271)

14th Week

April

## # Measure of Central Tendency :-

- » The numerical values that are used to represent mid-value or central value of a large collection of numerical data.
- » Part of Descriptive statistics.
- » Describe centre of dataset.

→ Three measures :-

(i) Mean

(ii) Median

(iii) Mode

(i) MEAN : used for the arithmetic mean of data.  
→ also known as "Average" (Avg).

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

mean = sum of all observations ÷ total no. of observations

e.g. If there are 5 observations, which are 27, 11, 17, 19, and 21 then the mean ( $\bar{x}$ ) is given by?

$$\begin{aligned}\Rightarrow \bar{x} &= (27+11+17+19+21) \div 5 \\ &= 95 \div 5 \\ &= 19\end{aligned}$$

APR	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T
78	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17

2018

THURSDAY

DAY (095-270)

14th Week

05

Appointment

Notes

Work to do

## (ii) MEDIAN :

- » middle value of the dataset when it is ordered.
- » Numerical variables.
- » Less influenced by outliers.
- » Even numbers, median is the average of two middle values.

example :  $x = \{2, 3, 1, 7, 5, 9\}$

$$x = \{1, 3, 5, 7, 9\}$$

median

e.g.  $x = \{1, 3, 5, 7, 9\}$

$$\bar{x} = \frac{5+7}{2}$$

$$= \frac{12}{2}$$

2 6 → median

## (iii) MODE :-

- » most frequent value (repeated value).
- » for categorical variables
- » Numerical variables, unique count be less

e.g.  $x = \{2, 4, 3, 4, 2, 6, 7, 2\}$

2 2 → mode

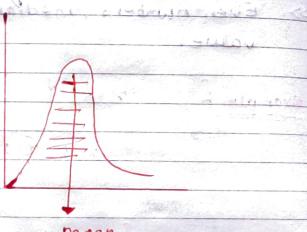
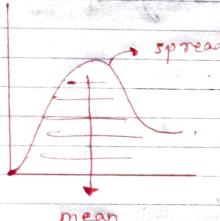
MAY	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T
18	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17

Appointment Notes Work to do

## # Measure of Dispersion :-

&gt; Degree of variation.

&gt; Known as (Spread, Variability).



- > Range
- > Interquartile Range
- > Quartile Deviation
- > Percentiles
- > Variance
- > Standard Deviation

i) Range is difference b/w maximum and minimum value of data.

$$\text{range} = \text{Maximum} - \text{Minimum}$$

• Sensitive to extreme values.

$$x = \{1, 2, 3, 4\}$$

$$= 4 - 1$$

$$\Rightarrow 3 \rightarrow \text{range}$$

APR	S	M	T	W	F	S	M	T	W	F	S	S	M	T	W	F	S	S	M	T	W
18	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21

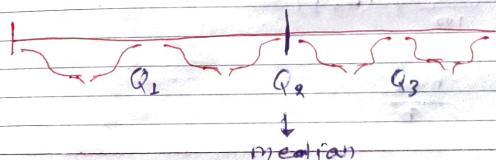
April

2018

Appointment Notes Work to do

(ii.) Quartile Deviation is half of the difference between the third quartile and the first quartile in a given data set.

> divide four equal parts with three quartile in data set.



$$\Rightarrow Q_1 = \frac{n+1}{4} \Rightarrow Q_2 = \frac{n+1}{2} \Rightarrow Q_3 = \frac{3(n+1)}{4}$$

$$\textcircled{1} \quad x = \{2, 4, 6, 7, 8, 10, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30\}$$

$$\Rightarrow Q_1 = \frac{n+1}{4} \Rightarrow \frac{7+1}{4} \Rightarrow \frac{8}{4} \Rightarrow 2$$

$$\Rightarrow Q_1 = 2$$

$$\Rightarrow Q_2 = \frac{n+1}{2} \Rightarrow \frac{7+1}{2} \Rightarrow \frac{8}{2} \Rightarrow 4$$

$$\Rightarrow Q_2 = 4$$

$$\Rightarrow Q_3 = \frac{3(n+1)}{4} \Rightarrow \frac{3(7+1)}{4} \Rightarrow \frac{3 \times 8}{4} \Rightarrow \frac{24}{4} \Rightarrow 6$$

$$\Rightarrow Q_3 = 6$$

09

MONDAY  
DAY (099-266)  
15th Week

April

Appointment

Notes

Work to do

iii) Percentiles :-

» Divide data into 100 equal parts.

$$\begin{aligned} \gg P_1 &= \frac{P}{100} \times (n+1) \\ &= \frac{25}{100} \times (n+1) \\ &= \frac{n+1}{4} \end{aligned}$$

iv.) Interquartile Range :-

» The difference between upper ( $Q_3$ ) and lower ( $Q_1$ ) quartile.

$$\gg Q_3 - Q_1$$

» Less sensitive to extreme values.

v) Variance :- the average of the square deviations from the mean of the given data set.

$$V = \frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N}$$

$x_i$  :- individual data points

$$\bar{x} = \text{mean}$$

APR	S	M	T	W	F	S	S	M	T	W	F	S	S	M	T
'18	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15

TUESDAY  
DAY (108-265)  
15th Week

10

April 18

Appointment

Notes

Work to do

(vi) Standard Deviation :-

» Square root of arithmetic average of the square of the deviations measured from the mean.

» Unit is same to the original dataset.

$$\therefore \text{std} = \sqrt{\text{variance}}$$

# Frequency :-

» Number of times, a value of data occurs.

» Categorical variables.

Frequency Distribution Table :-

Example :-

{1, 3, 3, 2, 4, 1, 2, 0, 1, 2, 3, 5, 4, 1, 2, 1, 3, 4, 1}

⇒ Data Value

Frequency

1	1
2	2
3	3
4	1
5	1

T	W	F	S	S	M	T	W	T	F	S	S	M	T	W
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15

Appointment

Notes

Work to do

\* Relative Frequency :-

- percentage or portion of the data value present in complete dataset.

$$\bullet \text{ R.F.} = \frac{\text{Frequency}}{\text{total no. of observation}} \times 100$$

Data Value	Frequency	Relative Frequency
1	6	$6/20 = 0.3 = 30\%$
2	5	$5/20 = 0.25 = 25\%$
3	4	$4/20 = 0.2 = 20\%$
4	3	$3/20 = 0.15 = 15\%$
5	1	$1/20 = 0.05 = 5\%$
	20	

\* Cumulative Frequency :-

» way to show the running total frequency as moved through the categories.

» statistical method that can be used to organize and examine data.

e.g. Customer complaints at a Company :-

- > Shipping Delay : 15 complaints
- > Billing Issues : 10 complaints
- > Returns : 12 complaints

Apr	S	M	T	W	F	S	S	M	T	W	F	S	S	M	T
18	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15

# Graphical Representations :-

- » A way to representing any data in picturized form.
- » help reader to understand the large set of data very easily as it gives various data patterns in visualized form.

→ Types :-

- Histogram
- Boxplots
- Scatterplots
- Line Graphs
- Bar Graphs
- Line Plot
- Pie Chart



May	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31

13

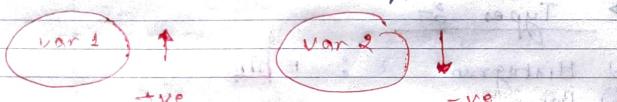
FRIDAY  
DAY (103-262)  
15th Week

## # Covariance :-

» statistical method (measured)

» Describe that how much two variables change together.

» Explain direction (nor the strength of the relationship)...



## # Correlation :-

» describes the strength between two variables.

» Linear relation.

» Direction of the relationship. ( $-1 < p < 1$ )

\*  $0.8 \rightarrow$  strong (+ve)

\*  $0.9 \rightarrow$  (+ve)

\*  $0.5 < p < 0.5 \rightarrow$  no-relation

## ⇒ Limitations :-

- outliers

- it does not measure non-linear relationships.

April

2018

APR	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W												
18	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30

SATURDAY  
DAY (104-261)  
15th Week

14

⇒ X and Y  
⇒ pearson's correlation coefficient.

$$S(x,y) = \frac{\text{cov}(x,y)}{\text{std}(x) \cdot \text{std}(y)} = \frac{\text{cov}(x,y)}{\sqrt{\text{var}(x)} \sqrt{\text{var}(y)}}$$

∴ Correlation

⇒ X and Y

$$\Rightarrow \text{cov}(x,y) = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$$

∴ covariance

## # Causation :-

» Relationship between cause & effect.

» Direct connection in which one variable influences the other.

15 SUNDAY

MAY

APR	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W												
18	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30

16

MONDAY  
DAY (106-259)  
16th Week

Appointment

Notes

Work to do

## # Probability :-

» measure of likelihood of a particular event occurring.

» measure between 0 and 1. ( $0 \leftrightarrow 1$ ).

\* 0 : impossibility

\* 1 : certainty

» 0.5 (50%), 0.6 (60%) : varying degree of likelihood.

### (Basic Concept) :-

- Sample Space
- Random Experiment
- Events
- Probability Function

### (i) Sample Space :-

» set of all possible outcomes of random experiment

example :- Tossing a coin.

$$\rightarrow S = \{\text{Tail, Head}\}$$

April

2018

17

TUESDAY  
DAY (107-258)  
16th Week

Appointment

Notes

Work to do

### (ii) Event :-

» subset of sample space.

e.g. A = Getting an even number.  
 $A = \{2, 4, 6\}$

B = Getting prime numbers.  
 $B = \{2, 3, 5, 7\}$

### (iii) Probability Function :-

» assigns probability to each event.

e.g.  $P(\text{getting head}) = \frac{1}{2} = \frac{\text{No. of favourable outcomes}}{\text{Total no. of outcomes}}$

$$P(\text{getting even no.}) = \frac{3}{6} \Rightarrow \frac{1}{2} \Rightarrow 0.5 \Rightarrow 50\%$$

### (iv) Complement of an event :-

» it consists of all outcomes i.e. not in A.

» taking of probability of an event w.r.t. I should not include probability.

$$\text{e.g. } A = \{2, 4, 6\}$$

$$A' = \{3, 5, 7\}$$

$$\therefore A = 0.5$$

$$A' = 1 - P(A)$$

APR	S	M	T	W	F	S	M	T	W	F	S	M	T	W	F	S	S	M	T	W
18	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20

APR	S	M	T	W	F	S	M	T	W	F	S	M	T	W	F	S	S	M	T	W
18	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20

MAY

18

WEDNESDAY  
DAY (108-257)  
16th Week

April

## Appointment Notes Work to do

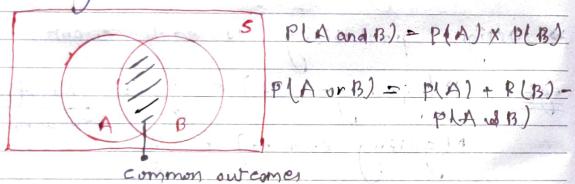
### # Event :-

> subset of sample space.

### # Types :-

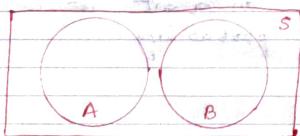
#### (i) Joint Event :-

- non - exclusive events.
- non - mutually exclusive.



#### (ii) Disjoint Event :-

- mutually exclusive.
- do not occur at same time.



APR	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W						
18	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31

2018

THURSDAY  
DAY (109-256)  
16th Week

19

## Appointment Notes Work to do

### (ii.) Dependent Event :-

- these events affect each other.

$$A \rightarrow B$$

eg  $A =$  drawing a red card  
 $B =$  drawing a black card

$$\Rightarrow A = \frac{26}{52}$$

$$\Rightarrow B = \frac{26}{51}$$

$$\therefore P(A \text{ and } B) = P(A) \times P(B|A)$$

$$\therefore P(B|A) = p \text{ of } B \text{ given } A.$$

### (iv.) Independent Events :-

- not affecting each other.

eg  $A =$  tossing a coin  $\{H, T\}$   
 $B =$  "O" and coin  $\{H, T\}$

$$\therefore P(A \text{ and } B) = P(A) \times P(B)$$

APR	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W						
19	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31

MAY  
18

20

FRIDAY  
DAY (110-255)  
16th Week

April

## # Conditional Probability :-

- probability of an event based on the occurrence of previous event.
- previous event has already occurred.

$P(A|B)$  = probability of A given B already occur.

$$\Rightarrow P(A|B) = \frac{P(A \text{ and } B)}{P(B)}$$

## # Bayes Theorem :-

- Reverend Thomas Bayes.
- update the probability of an event based on new evidence/information.

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

- $P(A|B)$  = update probability
- $P(A)$  = initial belief
- $P(B)$  = marginal probability
- $P(B|A)$  = evidence that supports the initial belief.

APR	S	M	T	W	F	S	S	M	T	W	F	S	S	M	T	W	F	S	S	M	T								
18	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30

2018

SATURDAY  
DAY (111-254)  
16th Week

21

Appointment Notes Work to do

## → Use Case :-

- \* Medical diagnosis
- \* Spam Classification (email spam)
- \* Recommendation System
- \* Fraud Detection

example :- Suppose there are 1000 people from which 100 are healthy and 900 are sick and they diagnosed using a test that is 99% effective.

	Diagnosed Healthy	Diagnosed Sick	
Healthy	891	9	error
Sick	1	99	

$$\Rightarrow P(S|D^S)$$

$$\Rightarrow P(S) = \frac{100}{1000} = 0.1$$

$$\Rightarrow P(S') = \frac{900}{1000} = 0.9$$

$$\Rightarrow P(D|S) = 0.99$$

$$\Rightarrow P(D|S') = 0.01$$

$$\Rightarrow P(P) = 0.1 \times 0.99$$

$$\therefore P(S|P) = \frac{P(P|S) \times P(S)}{P(P)}$$

$$= \frac{0.99 \times 0.1}{0.1 \times 0.99 + 0.01 \times 0.9} = \frac{0.99}{0.9967} = 0.9967 / 99.67\%$$

22 SUNDAY

MAY

JUN

23

MONDAY  
DAY (113-252)  
17th Week

Appointment

Notes

Work to do

## # Probability Distribution :-

- » mathematical distribution function that describes the probability of different possible values of a variable.
- » Random Variables :-

↳ Discrete : countable / finite numbers {0, 1, ...}   
 ↳ Continuous : uncountable numbers {0, 1}

## • PMF (Discrete Random Variables) :-

(probability mass function)

## • PDF (Continuous Random Variables) :-

(probability density function)

APR	S	M	T	W	F	S	S	M	T	W	F	S	S	M	T	W	F	S	S	M	T	W
18	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22

April  
2018

24

TUESDAY  
DAY (114-251)  
17th Week

Appointment

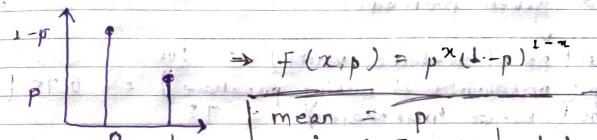
Notes

Work to do

## # Bernoulli Distributions :-

- » Discrete distributions
- » outcomes : {0, 1}

$$\text{PMF} = \begin{cases} p & \text{if } x = 1 \\ q = 1-p & \text{if } x = 0 \end{cases}$$



$$\begin{aligned} \text{mean} &= p \\ \text{variance} &= pq / p(1-p) \end{aligned}$$

- » multiple trials

## → Binomial Distributions :-

- » single trials

$$\Rightarrow (n=1)$$

- » outcomes : {0, 1}

$$\text{PMF} = f(x; n, p) \Rightarrow n(x) p^x (1-p)^{n-x}$$

$$\Rightarrow \frac{n!}{(n-x)! x!} \cdot p^x (1-p)^{n-x}$$

$$\star \text{ mean} = np$$

$$\star \text{ variance} = npq \Rightarrow np(1-p)$$

APR	S	M	T	F	S	S	M	T	W	F	S	S	M	T	W	F	S	S	M	T	W	May	
18	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	18

25

## WEDNESDAY

DAY (115-250)  
17th Week

## Appointment

Example 5: Customer conversion rate

Imagine you're an e-commerce business owner, and you want to analyze the conversion rate of visitors to your website. You're interested in understanding the likelihood of specific no. of purchases (successes) out of fixed no. of website visits during a given period.

$n$  (no. of website visits) = 100

$$p(\text{probability of making purchase}) = 0.75 \quad (\text{5% rate})$$

$x(\text{no. of purchases}) = 75$

$x$  (No. of purchases) = 75

Question 2: What is probability of exactly 75 visitors making a purchase out of website visitors given a 75% conversion rate?

$$\Rightarrow PMF = \frac{100!}{(100-75)! \times 75!} \times 0.75^75 \times (1-0.75)^{25}$$

= 0.091-

= 9.17

2018

**THURSDAY**  
DAY (116-249)  
17th Week

26

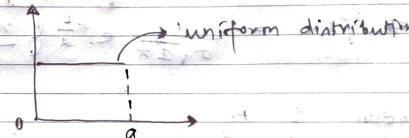
**Appointment:** Dr. John Doe

## → Uniform Distribution :-

» continuous distribution

↳ characterized by a constant probability density (cpd) over our specify interval.

→ Equal chance to occurs.



$$\text{PDF} = f(x) = \frac{1}{b-a}$$

$$\rightarrow a \leq x \leq b$$

$$\text{mean} = \frac{a+b}{2}$$

$$\text{Variance} = \frac{(b-a)^2}{12}$$

27

FRIDAY  
DAY (117-248)  
17th Week

Appointment

Notes

Work to do

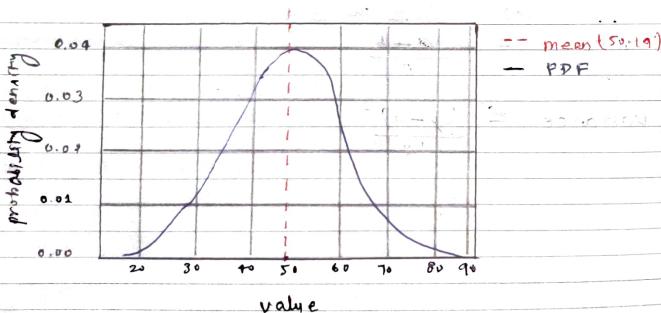
## # Normal Distributions :-

### Gaussian Distribution.

- it appears as a "bell curve" when graphed.
- describes a symmetrical plot of data around its mean value.



$$\text{PDF} \Rightarrow f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

\* mean =  $\mu$ ,\* variance =  $\sigma^2$  (sigma squared)\* std. deviation =  $\sigma$ 

$$\Rightarrow \text{mean} = 50$$

$$\Rightarrow \text{std} = 10$$

Apr	S	M	T	W	F	S	S	M	T	W	F	S	S	M	T
'18	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15

April

SATURDAY  
DAY (118-248)  
17th Week

28

2018

Appointment

Notes

Work to do

## # Standardization :-

- process of converting normal distribution into standard normal distribution.

- mean = 0.
- standard deviation = 1

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

## Normalization :-

- rescales a dataset so that each value falls between 0 and 1.

$$x_{\text{new}} = \frac{x - x_{\text{min}}}{x_{\text{max}} - x_{\text{min}}}$$

example :- Salary : minimum = 3 lakhs  
maximum = 50 "

$$\Rightarrow \frac{4 - 3}{50 - 3} \Rightarrow \frac{1}{47} \Rightarrow 0 \leftrightarrow 1$$

T	W	T	F	S	S	M	T	W	F	S	S	M	T	W
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15

MAY

30

MONDAY

DAY (120-245)

18th Week

April

Appointment

Notes

Work to do

### Normalization

Minimum and maximum value for variable are used for scaling.

It is used when features are of different scales.

Scale value between 0 and 1.

It is affected by outliers.

It is useful when we don't know about the distribution.

Maps the data to a specific range, which may leads to loss of information about the original distribution.

### Standardization

Mean and standard deviation of variables are used for scaling.

It is used when we want to ensure zero mean and unit standard deviation.

It is not bounded to a certain range.

It is less effected by outliers.

It is useful when the feature distribution is normal or gaussian.

Preserves the shape of the original data distribution but shifts and scales it.

APR '18	S	M	T	W	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23

01

**TUESDAY**  
DAY (121-24)  
18th Week

### Appointments

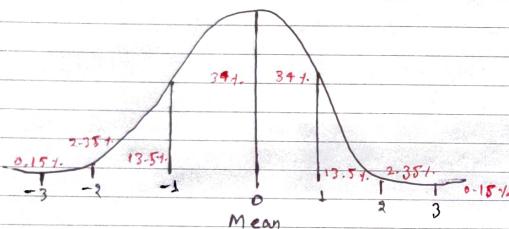
### Note:

#### Work to do

## # Empirical Rule

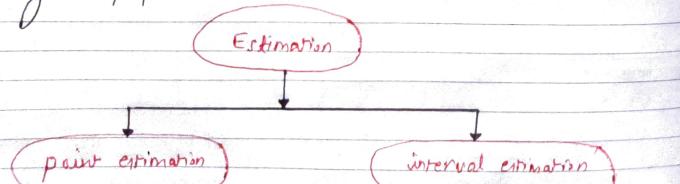
- >> also known as 68-95-99-7

- >> describe how data is distributed in a normal distribution.



## # Inferential Statistics :-

- involves the use of a sample (1) to estimate some characteristics in a large population.
  - (2) to test a research hypothesis about a given population.



May

2018

**WEDNESDAY**  
DAY (122-243)  
18th Week

02

ANSWER

## Notes

### Work to do

- ## • Estimation

E&M a process of drawing conclusion  
for population based on sample data.

- >> population parameter
  - >> Sample statistics

## i.) Point Estimation

- » provide a single value as your best guess for unknown population parameters

## >> Properties :-

- Consistent : larger the sample size  
: more accurate estimation.

## • Unbiased

## Drawbacks

- (i) problem of reliability  
(ii) not enough evidence

<b>MAY '18</b>	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T									
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31

F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S									
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	*	JUN 18

03

THURSDAY  
DAY (123-242)  
18th Week

May

Appointment

Notes

Work to do

### (ii) Interval Estimation :-

- >> provide intervals, in which population parameters fall.
- >> more accurate than point estimation.

### # Confidence Interval :-

- >> an interval of values i.e. computed from sample data i.e. likely to contain the true population parameter.
- >> 95% - 99% (confidence coefficient)

$$\therefore CI = \bar{x} \pm z \frac{s}{\sqrt{n}}$$

 $\bar{x}$  : sample mean $z$  : confidence level value $s$  : sample standard deviation $n$  : sample size

May	T	W	F	S	S	M	T	W	F	S	S	M	T	W	T	F	S	S	M	T	W	F	S
18	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23

FRIDAY  
DAY (124-241)  
18th Week

04

2018

Appointment

Notes

Work to do

### # Student's T Distribution :-

- >> used while assumptions about a mean when we don't know the standard deviation.

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

### # Hypothesis Testing :-

- >> statistical method used to evaluate a claim / hypothesis about population parameter.
- >> proving True / False.
- >> comparing sample data to a hypothesis about a population parameter to determine whether there is enough evidence to reject or correct the hypothesis.
- >> without hypothesis test, you risk drawing wrong conclusion & making bad decisions.

### → Types :-

(i) Null

(ii) Alternate

May	T	W	F	S	S	M	T	W	F	S	S	M	T	W	T	F	S	S	M	T	W	F	S
19	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23

05

SATURDAY  
DAY (125-240)  
18th Week

Appointment

Notes

Work to do

(i) Null Hypothesis :- ( $H_0$ )

- » baseline
- » default assumption
- » " position
- » Accept / Reject

(ii) Alternate Hypothesis :- ( $H_1$ )

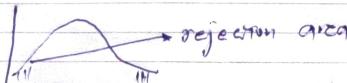
- » Logical opposite of null hypothesis
- » there is significant difference or effect.

Example :- Sanitizer company claims that their product kills 99% germs?

06 SUNDAY  $\rightarrow H_0 : \mu = 99\%$   
 $H_1 : \mu < 99\%$

$\rightarrow$  score / level

$\therefore$  level of significance ( $\alpha$ )  $\Rightarrow$  predetermined threshold



$$\alpha = 0.05$$

$$\alpha = 0.05$$

May

MAY T W T F S M T W T F S M T W T F S M T W T F S S M T W T F S S M T W T F S S M T W T F S S JUN  
18 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31

MONDAY  
DAY (127-239)  
19th Week

07

2018

(ERROR) :-

Appointment

Notes

Work to do

 $\rightarrow$  Error in Hypothesis Testing :-

$\gg$  Error occurs when conclusions are drawn that are incorrect.

\* Types of Errors :-

- (i) Type I error  $\Rightarrow$  (false +ve)
- (ii) Type II error  $\Rightarrow$  (false -ve)

• Type I errors occur when the null hypothesis is rejected when it is true.

• Type II error occur when null hypothesis is accepted when it is false.

	Reject $H_0$	Accept $H_0$
$H_0$ True	Type I error	Correct
$H_0$ False	Correct	Type II error

08

TUESDAY  
DAY (128-237)  
19th Week

Appointment

Notes

Work to do

Example : You are conducting a medical test to detect a rare disease.

$H_0$  : patient is healthy

$H_1$  : patient has disease

⇒ Type I error : patient has a disease when he does not or he is healthy.

⇒ Type II error : patient is healthy when he suffers from that disease.

May  
(2018)May  
2018WEDNESDAY  
DAY (129-236)  
19th Week

09

Appointment

Notes

Work to do

## # TEST :-

» a statistical procedure used to determine if there's enough evidence in a sample data to draw conclusion about a population by comparing the observed data to what is expected under a null hypothesis.

### \* Types of Tests :-

#### (i) One Tailed : (unidirectional)

» critical region for rejecting  $H_0$  is located in only one tail of the distribution.

#### (ii) Two Tailed : (bi-directional)

» critical region for rejecting  $H_0$  is divided between both tails of the distribution.

MAY	T	W	T	F	S	S	M	T	W	F	S	S	M	T	W	T	F
18	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17

F	S	S	M	T	W	F	S	S	M	T	W	F	S	S	M	T	W	F	S	S	JUN
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22

18

Appointment

Notes

Work to do

## # Test in Statistics :-

Z-Test :-

- » compare sample mean with population mean.
- » population standard deviation is known.
- » useful for large sample size ( $n > 30$ )

\* One sample Z-test :-

- when we have single sample to compare with population mean.

$$\left[ \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} \right]$$

\* Two sample Z-test :-

- two independent samples.
- compare their mean

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

MAY	T	W	F	S	S	M	T	W	F	S	S	M	T	W	F	S	S	M	T	W	F	S								
18	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30

2018

Appointment

Notes

Work to do

(iii) T-Test :-» sample size small ( $n < 30$ )

» Compare mean of two groups or mean of a sample to a unknown value.

» when  $\sigma$  (sigma) is unknown in unknown.\* Independent T-test :-

- compare mean of two independent groups to determine if there is any significant difference between them.

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

\* Paired T-test :-

- comparing mean of two related group such as before and after measurement on the same subjects.

$$\Rightarrow \frac{\bar{x}_{\text{diff}}}{\frac{s_{\text{diff}}}{\sqrt{n}}}$$

F	S	S	M	T	W	F	S	S	M	T	W	F	S	S	M	T	W	F	S	S	JUN								
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30

12

SATURDAY  
DAY (132-233)  
19th Week

Anova Test :- (Analysis of Variance)

- » used to determine if there's a statistically significant difference between the means of three or more groups, generalizing the t-test beyond two groups.
- » F-value (F-test)
- » Compare more than two groups.

## \* One-way Anova

- used when you have one independent variable (factor) with multiple levels.

## \* Two-way Anova :

- used when you have two or more independent variables.

13 SUNDAY

May

2018

14

**MONDAY**  
DAY (134-231)  
20th Week

Appointment	Notes	Work to do
(iv) <u>Chi-square Test §</u>	» compare two categorical variables. » Test the association or independence between categories.	

### \* Chi-square Test for independence :-

- used to determine if two categorical variables are related.

## \* Chi-square Test for Goodness-of-fit

- check if a sample distribution matches an expected distribution.

F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S	S	M	T	W	T	F	S								
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31