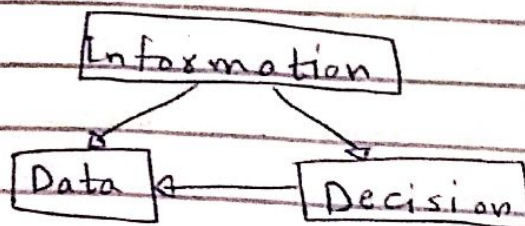# Probability & Statistics - MT2005

## Lec #01

### Statistics:-

- Collection of data.
- Presentation of data.
- Analysis of data.
- Decision about data.

```
         ┌─────────────┐
         │ Information │
         └─────────────┘
          ↙           ↘
   ┌──────┐         ┌──────────┐
   │ Data │ ◄────── │ Decision │
   └──────┘         └──────────┘
```

---

## Lec # 02

### Data:-

Raw facts and figures are called "data".

### Population:-

Population is the entire set of items from which you draw data for statistical study. It can be a group of individuals, set of items.

Noun as census, enumeration,

Cost ↑
Time ↑

| Common person who did not knows about statistics | A person who knows about statistics |
|---|---|
| Collection of data not well. | Collection. |
| Presentation. | Presentation. |
| Analysis of data not well. | Analysis. |
| Decision. | Decision. |

# Lec#03

## Sample:-

Sample is the part/chunk/sub-part of population. Which contains the whole characteristics of the population.

• Sample based studies are called sample-surveys, e.g., In CS how much students have glosses, in class of 39,

No. of students with closses = 19
Total students in class = 39

$$Result = \frac{19}{39} = 49\%.$$

## Data :-

Raw facts & figures are called data.

## Types of data :-

### 1) Qualitative :-

Data about quality or characteristic of variable. e.g., name, color, etc.

### 2) Quantitative :-

Data about numeric variable or can be represented in numeric form. e.g. height, weight, CGPA, etc.
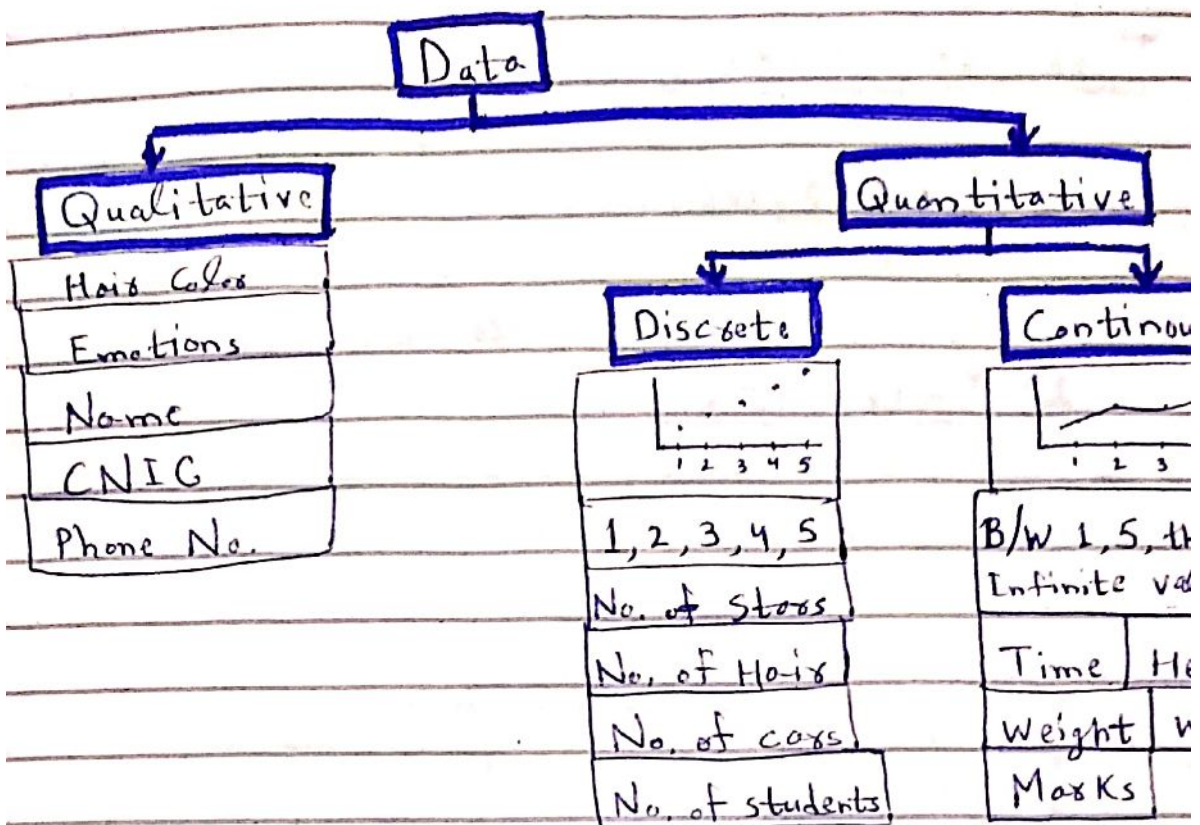
#### 2.1) Discrete :-

A variable or numeric variable that is obtained by counting. e.g., no. of present students.

#### 2.2) Continous :-

A variable that is obtained by measuring. e.g., speed, temperature, energy, e

```
            Data Variable
                  |
        ----------------------
        |                    |
  Qualitative          Quantitative
                            |
                    ----------------
                    |              |
                Discrete      Continous
```

# Lec#04

```
                    ┌──────┐
                    │ Data │
                    └──┬───┘
         ┌─────────────┴──────────────┐
         ▼                            ▼
┌──────────────┐            ┌──────────────┐
│ Qualitative  │            │ Quantitative │
└──────────────┘            └──────┬───────┘
                              ┌─────┴──────┐
                              ▼            ▼
                        ┌──────────┐  ┌──────────┐
                        │ Discrete │  │ Continuou│
                        └──────────┘  └──────────┘
```

| Qualitative |
|---|
| Hair Color |
| Emotions |
| Name |
| CNIC |
| Phone No. |

Discrete:

| 1, 2, 3, 4, 5 |
|---|
| No. of Stors |
| No. of Hair |
| No. of cars |
| No. of students |

Continuou:

| B/W 1,5, th |
|---|
| Infinite va |
| Time | He |
| Weight | v |
| Marks | |

- We cannot perform statistics method to qualitative data. So, we assign code to them. e.g.,

| Gender | Gender. |
|--------|---------|
| M      | 1       |
| F      | 2       |
| M      | 1       |

Assign,

$M = 1$

$F = 2$

- We can convert quantitative data to qualitative like, marks to grade, IQ Level
- We cannot convert qualitative data to quantitative data. e.g., grade to marks is not possible.

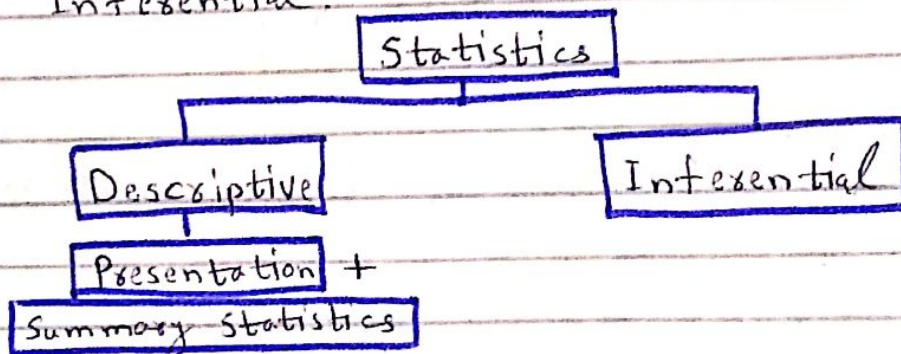<u>Lec#05</u>

## <u>Types of Statistics:-</u>

1) <u>Descriptive:-</u>
   How to describe/Present data/method
   - Presentation + ▓▓ Summary Statistics.

2) <u>Inferential:-</u>
   Generalization of sample statistic
   towards population parameter is called
   Inferential.

```
          ┌──────────────┐
          │  Statistics  │
          └──────┬───────┘
       ┌─────────┴─────────────┐
 ┌─────────────┐       ┌──────────────┐
 │ Descriptive │       │ Inferential  │
 └──────┬──────┘       └──────────────┘
 ┌──────────────┐
 │ Presentation │ +
 │Summary Statistics│
 └──────────────┘
```

## <u>Parameters:-</u>
The characteristics of anything
is called parameter (from population).

## <u>Statistic:-</u>
The characteristics of sample data.

## <u>Types of Collected Data:-</u>

1) <u>Primary Data:-</u>
   - The data that is collected
   first time, or first handed.
   - Raw form of data.
   - More time consuming, but with
   your own choice.
   - Collect data by yourself, don't

- Usually, questionnaire used to collect data.

## 2) Secondary data :-
- When any statistical tool is apply on primary data then it becomes secondary data.
  - Second-handed data.
  - Less time consuming, but we don't have choice.
  - Ready-to-use data, and we have to compromise on it.
  - Use websites to collect organize data.

```
            ┌─────────────────┐
            │  Collected Data │
            └────────┬────────┘
         ┌───────────┴───────────┐
┌──────────────┐         ┌────────────────┐
│ Primary Data │         │ Secondary Data │
└──────────────┘         └────────────────┘
```

# Lec#06

## Presentation of Data:-
Presentation of data has a very big role in our daily life.

## Types of Presentation of data:-

1) Tabular.
2) Graphical.

## 1) Tabular :-

### • Frequency Distribution :-

### • Classes :-

1. All of data should be there in the range of classes.

• Two Important things about classes are :-

1) No. of Classes : 6 or 7
2) Size of Class : Equal Gap

• Equal Gap should be there in the size of classes to make algorithm, etc because different gap is subjective thing & can't make algorithm, etc.

• To check / to know the equal gap for given data :-

$$\frac{max - min}{No. \; of \; Classes} \quad e.g., \quad \frac{max - min}{6 \; or \; 7}$$

**Wrong Approach**

| Classes | Classes |
|---------|---------|
| 0 – 50  | 0 – 2   |
| 50 – 100| 2 – 4   |
|         | 4 – 6   |
|         | 6 – 8   |

**Wrong Approach**

| Classes |
|---------|
| 0 – 10  |
| 11 – 20 |
| 21 – 30 |

| Classes |
|---------|
| 0 – 10  |
| 10 – 20 |
| 20 – 30 |

→ Less than $<$
→ Greater than Equal to $>=$   OR   Vice Versa

### Example :-

Age          Counts/No.

| Classes | Frequency (f) | Commulative Frequency (c.f) | | Relative Frequency (r.f) |
|---------|---------------|------|--------|--------------------------|
|         |               | Normal | Inverse | |
| 0 – 10  | 1             | 1    | 30     | $1/30 = 0.033 = 3.3\%$ |
| 10 – 20 | 3             | 4    | 29     | $3/30 = 0.1 = 10\%$ |
| 20 – 30 | 5             | 9    | 26     | $5/30 = 0.166 = 16.6\%$ |
| 30 – 40 | 10            | 19   | 21     | $10/30 = 0.333 = 33.3\%$ |
| 40 – 50 | 6             | 25   | 11     | $6/30 = 0.2 = 20\%$ |
| 50 – 60 | 4             | 29   | 5      | $4/30 = 0.13 = 13\%$ |
| 60 above| 1             | 30   | 1      | $1/30 = 0.033 = 3.3\%$ |

| Genders | f |
|---------|---|
| M       | 10 |
| F       | 3 |

• Qualitative data don't have Classes. We cannot assign codes to qualitative but we can do all others. Grouping of cities into divisions or provinces.

## 2) Graphical :-

Considered the type of data while collecting.

### Methods :-

1) Bar Chart :-
   1. Simple bar chart.
   2. Multiple bar chart.
   3. Component bar chart.
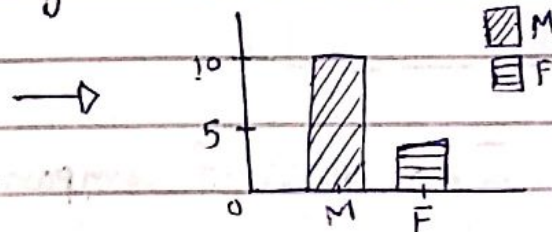2) Pie Chart.
3) Histogram.

### 1) Bar Chart :-

• Bar Chart are just for qualitative data.
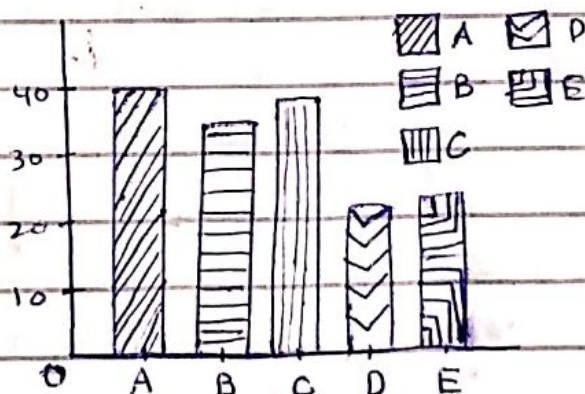
#### 1. Simple bar Chart :-
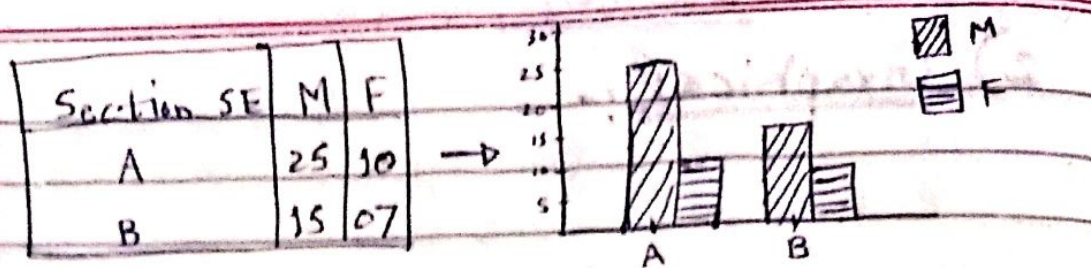
• Only for one variable of data.

1)

| Gender | f |
|--------|---|
| M | 10 |
| F | 3 |

→

2)

| Section | f |
|---------|----|
| A | 40 |
| B | 35 |
| C | 38 |
| D | 21 |
| E | 23 |

#### 2. Multiple bar Chart :-

• Used for two variables of qualitative data. e.g, section with gende

→

| Section SE | M | F |
|---|---|---|
| A | 25 | 10 |
| B | 15 | 07 |



## 3. Component bar Chart:-

· Initially same as multiple bar chart. It the total is meaningfull th.. moves to component bar chart. e.g...

| Section | M | F | Total |
|---|---|---|---|
| A | 25 | 10 | 35 |
| B | 15 | 07 | 22 |



Wrong Approach

| | PaK | India | Ban |
|---|---|---|---|
| GDP | — | — | — |
| NI | — | — | — |
| AP | — | — | — |
| Total | — | — | — |

Meaning less Total

## Examples of Components:-

① Hord disk in PC

C: 

Filled

② Battery



Filled

③ Video Player



Played

## 2) Pie Chart:-

Pie chart is only for one variable of qualitative data, e.g., Cities.

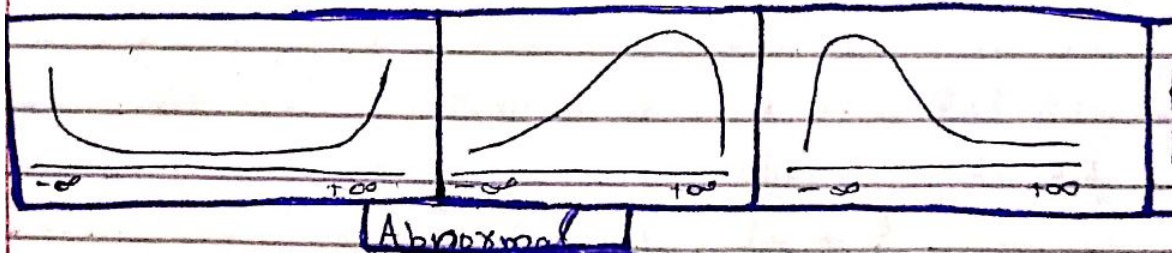| Cities | $f$ | Angle → $\frac{City}{Sum(f)} \times 360°$ |
|--------|-----|--------|
| FSD | 15 | 216° |
| CHT | 7 | 100.8 |
| LHR | 2 | 28.8 |
| ISB | 1 | 14.4 |
| Total | 25 | |

## 3) Histogram:-

Histogram is for ~~all~~ quantitative data.

• **Normal:-**

Bell shape curve.
Symentic curve.
Normal curve.

$-\infty$                    $+\infty$

Abnormal

Example of Histogram. e.g.,

| Classes | $f$ |
|---------|-----|
| 0 – 10 | 1 |
| 10 – 20 | 3 |
| 20 – 30 | 5 |
| 30 – 40 | 10 |
| 40 – 50 | 6 |
| 50 – 60 | 4 |
| 60 above | 2 |

Normal curve

- Also, Histogram tells the shape of data. e.g..

Abnormal $\rightarrow$ p

# Summary Statistics :-

- How to summarize the quantitative data. A single value that represents the whole data. Average/Mean.

- **Average/Mean** :-
  - Majority of data.
  - Centered Value.
  - Balancing Point.

- **Measure of central Location** :-
  1) **Mean** :-
     - Balancing Point of Data.

$$Mean = \overline{X} = \frac{Sum \ of \ all \ Observations}{No. \ of \ Observations} = \frac{X_1 + X_2 + .... + X_n}{n}$$

$$\overline{X} = \frac{\sum_{i=s}^{s} X_i}{n}$$

- But, there is effects on data/Mean by extreme observations. e.g., 1, 2, 3, 4, 500

$$\overline{X} = 510/5 = 102$$

- Can 102 be a representator of 1, 2, 3, 4, 500. Not/No. Now, moves to median.

  2) **Median** :-
     - ~~~~ Middle most observation.
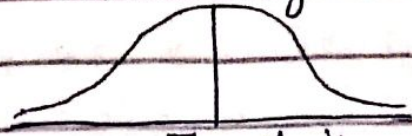     - Even Count of Data,
         1, 2, 3, 4, 5, 6
       Median = 3+4/2 = 3.5
     - Odd Count of Data,
         1, 2, 3, 4, 5
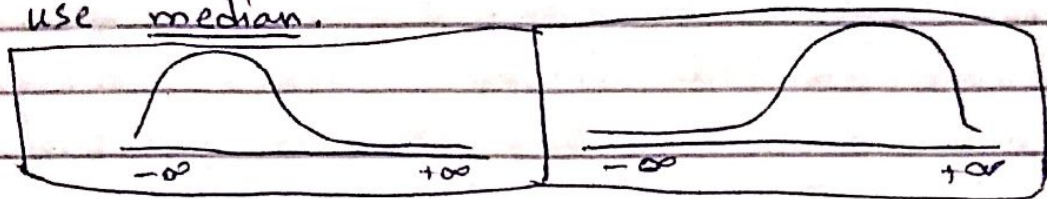       Median = 3

• Median is less sensitive method that cannot effect by outlyers.



Still in Normal → $\bar{X}$ = Median

• In case of no outlyer, prefer to use __mean__, because it has mathemati formula. • But, in given below case use __median__.



1, 2, 3, 4, 500, then use median.

---

<u>Lec#09</u>

## 3) Mode :-

• The most common value in qualitative data. __OR__
• Qualitative data with highest frequency. e.g., 1, 2, 3, 2, 2, 1, 1, 2, 2, 3, 2
$$Mode = 2$$

• Mean, Median, is only 50% of data
• To observe the qualitative or quantitativ data in chunks. e.g., __Quantile__ is used.



, e.g.,



Median = $Q_2$ (50%)

$Q_1$ (50% of data)

$Q_3$ (50% of data)

• Average is not enough to analyze the data. To know the Difference ẞ of data from average. We used dispersion or variation.

## Measure of Dispersion:-

### 1) Range:-

$R = Max - Min$

$R = 5 - 1 = 4$ , e.g., 1, 2, 3, 4, 5

this result tells the max distance or difference of data th. can occures.

$R = (x - \bar{x})$

$$R = \frac{\Sigma(x - \bar{x})}{n} = \frac{0}{n} = 0$$

| Data | $x - \bar{x}$ |
|------|---------------|
| 1 | $1 - 3 = -2$ |
| 2 | $2 - 3 = -1$ |
| 3 | $3 - 3 = 0$ |
| 4 | $4 - 3 = 1$ |
| 5 | $5 - 3 = 2$ |

So, Sum of the deviation of each observation from their mean = 0

$$\frac{\Sigma(x - \bar{x})}{n} = 0$$

To overcome this problem, we have:-

### 1) Variance:-

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

### 2) Standard Deviation:-

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

• Weight of some product mentioned as

$10g \pm 1g \quad (9 - 11)$

$\bar{x} \pm S.D$

| | Example |
|---|---|
| | 3 Sigma | $10g + 1g$ |

| | 3 Sigma | Example $10g + 1g$ |
|---|---|---|
| $\overline{X} \pm S.D * 1$ | 68.5 % | 9 – 11 |
| $\overline{X} \pm S.D * 2$ | 95.7% | 8 – 12 |
| $\overline{X} \pm S.D * 3$ | 99.5% | 7 – 13 |

Another Example: $(57 \pm 12)$ Marks

| | |
|---|---|
| $57 \pm 12 \ (12*1)$ | 45 – 69 |
| $57 \pm 24 \ (12*2)$ | 33 – 81 |
| $57 \pm 36 \ (12*3)$ | 21 – 93 |

### 3) Inter Quantile Range (IQR):

Upper Quantile → Lower Quantile

$$IQR = Q_3 - Q_1$$

1) Good Method for quantitative data with no outlier:-

$$\underline{\overline{X} \pm S.D}$$

2) Good Method for quantitative data with outliers:-

$$\underline{Median \pm IQR}$$

3) Good Method for qualitative data:-

$$\underline{Mode \pm (Max - Min)}$$

Qualitative wise Quantitative Analysis

| Gender | Age |
|---|---|
| M | — |
| F | — |
| M | — |
| M | — |
| F | — |
| F | — |

→

| Gender | Age | S.D | — |
|---|---|---|---|
| M | — | - | - |
| F | - | - | - |

Now, we can make graph of that.