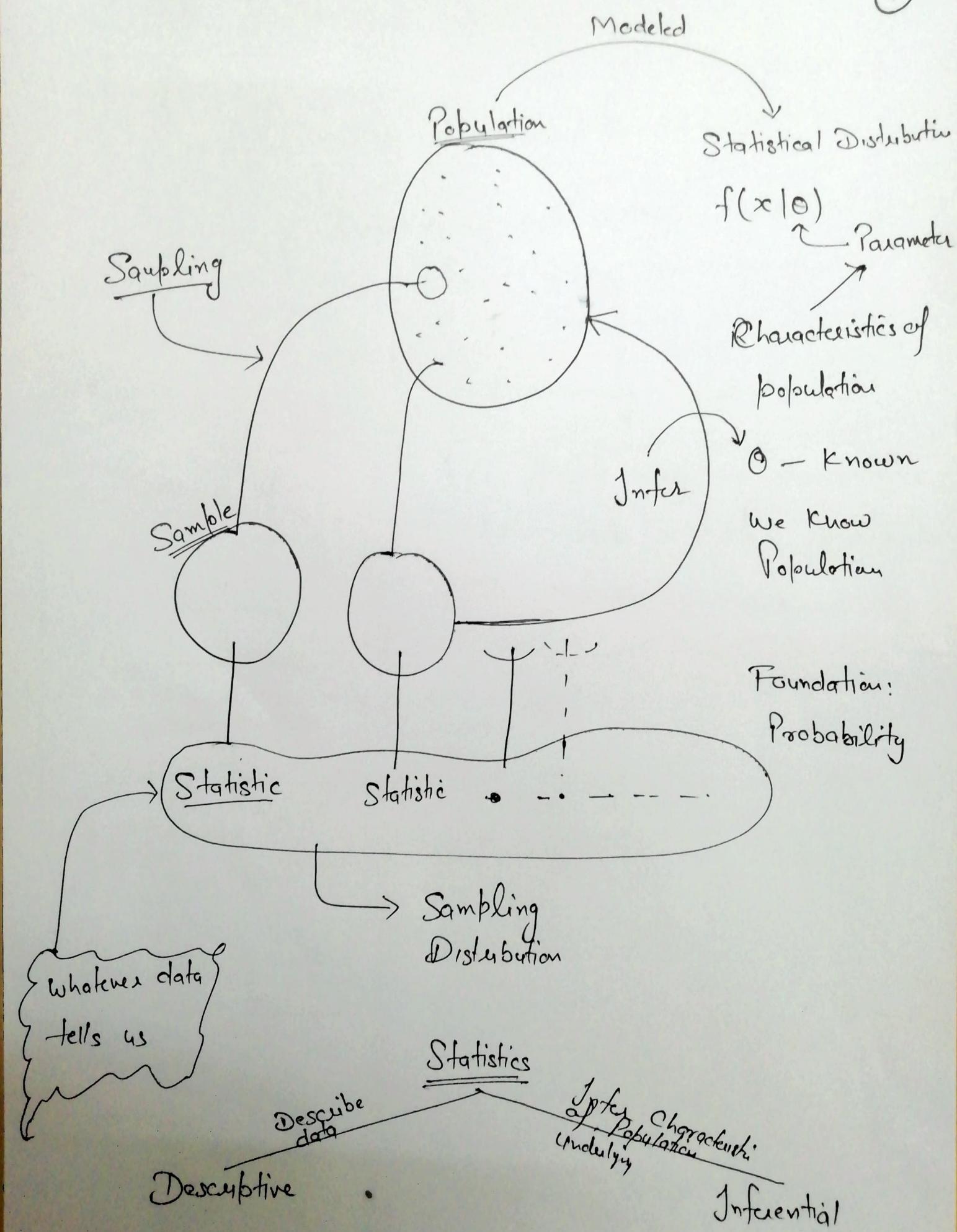


①



# Descriptive Statistics

## Calculating Descriptive Statistics

### → Measures of Central tendency

- Mean
- Median
- Mode
- —

### → Measures of variation

- Variance
- Standard deviation
- Range
- —

## Displaying Descriptive Statistics

Summarizing data in tables, charts and graphs

### Frequency Distribution

A simple way to summarize raw data in tables and make the information more useful

## Frequency distributions

Showing your data in a table

(3)

Data set: The daily demand for hammers at a hardware store over the last 20 days

Daily demand

2	1	0	2	1
.	.	.	.	.
3	0	2	4	0
3	2	3	4	2
2	2	4	3	0

A frequency distribution is a two-column table. In left column, list each values in the data set from least to greatest.

Daily demand      Frequency      Cumulative frequency      Relative freq.

0	<u>4</u>	4	$4/20 = 0.20$
---	----------	---	---------------

1	2	6	0.10
---	---	---	------

2	7	13	0.35
---	---	----	------

3	4	17	0.20
---	---	----	------

4	3	20	0.15
---	---	----	------

Count the no. of times each value appears and record data in right column

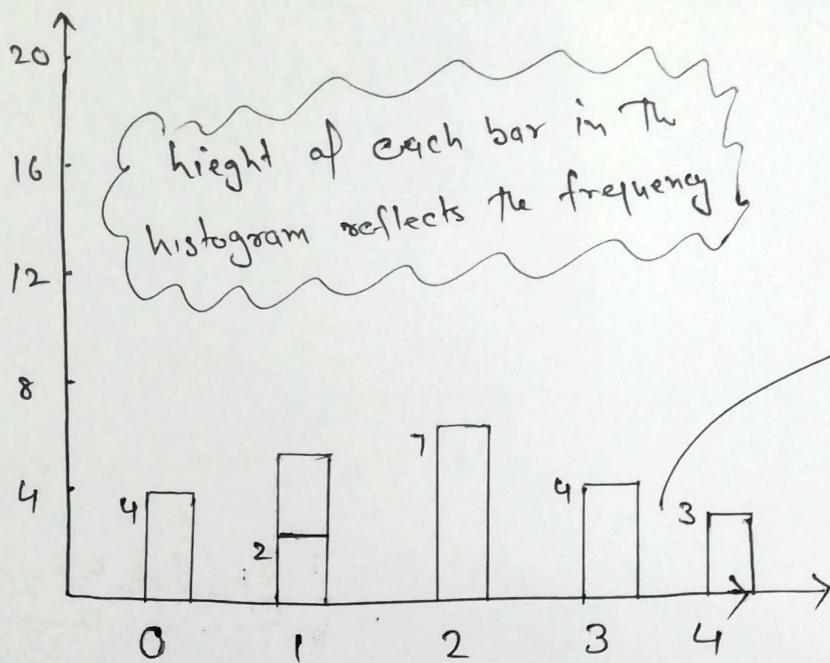
20

1

Always equal to 1

# Histograms — frequency distribution in a chart

(4)



A gap exists between columns because the data is discrete. Discrete data is data that can take on a countable no. of possible values.

Data: The mileage of a specific car with a full tank of gas.

## Miles per tank

302	315	265	296	289	301	308
280	285	318	267	300	309	312
299	316	301	286	281	311	272
295	305	283	309	313	278	284
296	291	310	302	282	287	307
305	314	318	308	280		

(5)

$$W = \frac{\text{Largest value} - \text{Smallest value}}{\text{number of classes}}$$

So we have

$$\text{Largest value} — 318$$

$$\text{Smallest value} — 265$$

$$\text{total data points} — 40$$

Therefore using  $2^k \geq n$  rule

$$\text{Set } k=6 \quad \therefore 2^6 \geq 40$$

Subsequently

$$W = \frac{318 - 265}{6} = 8.833 \equiv 10$$

Set the size of each class equal to 10.

Count the no. of values contained in each class

Miles per tank under	Frequency
260 — 270	2
270 — 280	2
280 — 290	10
290 — 300	5
300 — 310	12
310 — 320	9

Now data has many possible outcomes.  
So group into classes

$$k - \text{no. of classes}$$

$$n - \text{total no. of data points}$$

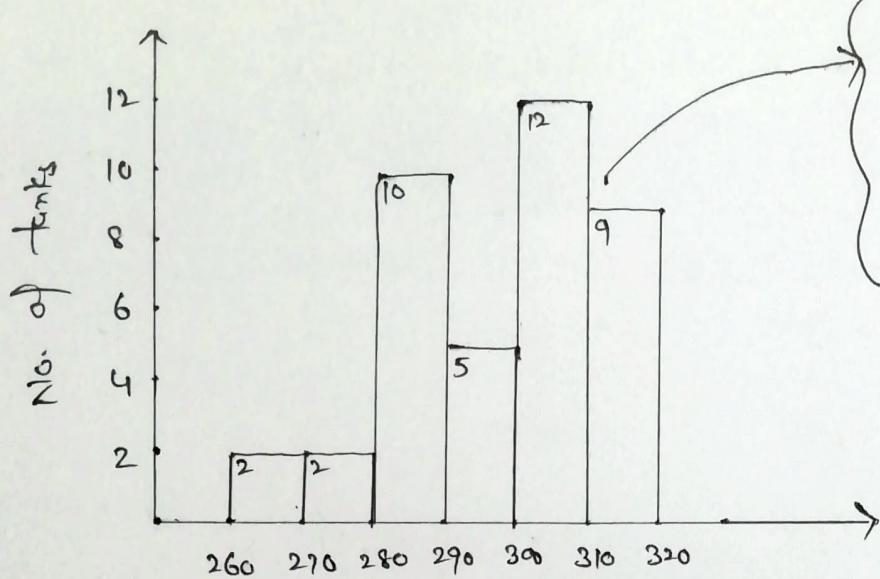
$$2^k \geq n$$

least  $k$  satisfy

$$W = \frac{L-S}{K}$$

Width of class      Largest value      Smallest value

(6)



No gap exists because  
The data is continuous.  
Continuous data can assume  
any value in an interval

Mileage per tank of gas

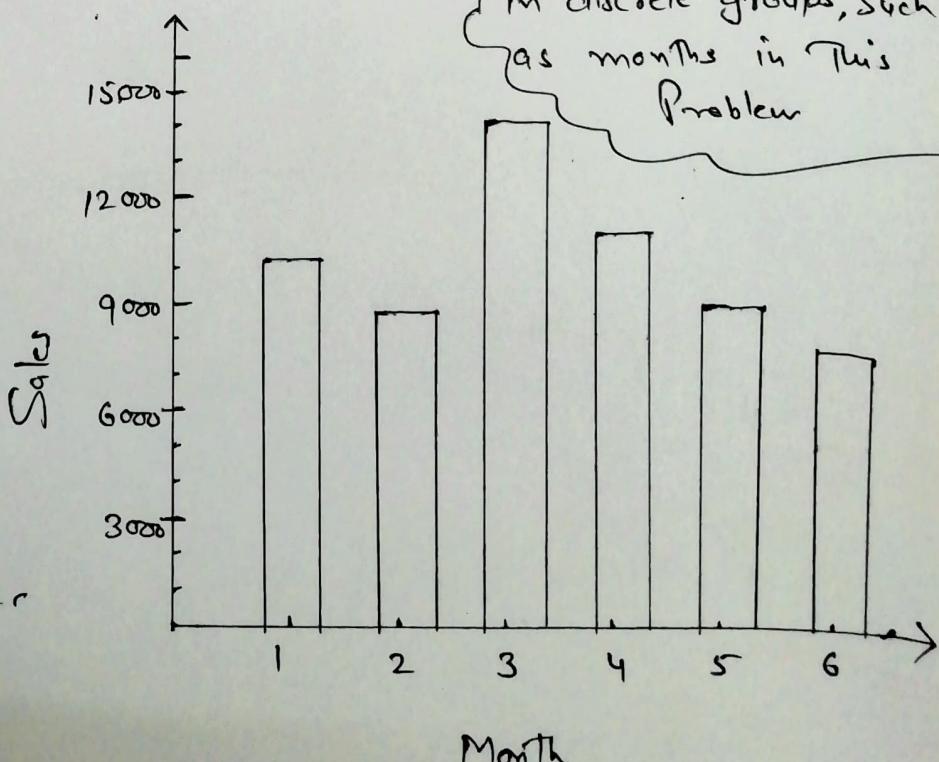
\* for Categorical data

Bar Charts

Setting The bar for visual data

Data: A Company's monthly sales totals

Month	Sales
1	10734
2	8726
3	14387
4	11213
5	9008
6	8430



Categorical data is data that is organized in discrete groups, such as months in this problem

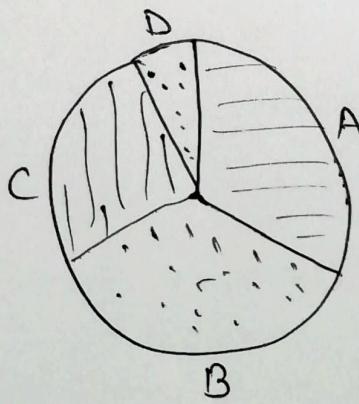
## Pie Charts

Showing your categorical data in a circle

(7)

Data: a grade distribution for a college class

Grade	No. of Students	Relative frequency	Central angle
A	9	$9/30 = 0.30$	$0.30 \times 360^\circ = 108^\circ$
B	12	0.40	$144^\circ$
C	7	0.23	$83^\circ$
D	2	0.07	$25^\circ$
Total	30	1.00	$360^\circ$



Sorry I am very  
poor in drawing  
Plot in your lab

Showy Categorical Data

→ Box chart

→ Pie chart

Histogram

Discrete data

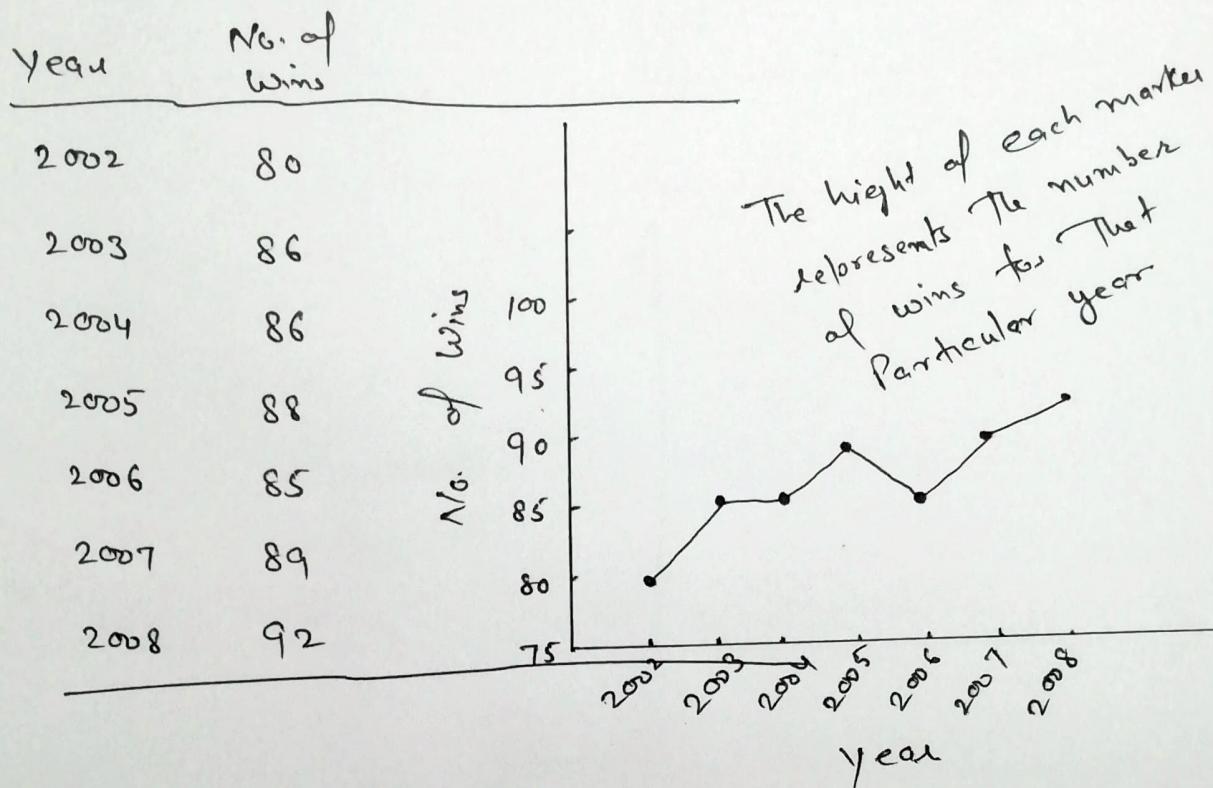
gap

Continuous classes

no gap

## Line Charts : data over time in a chart

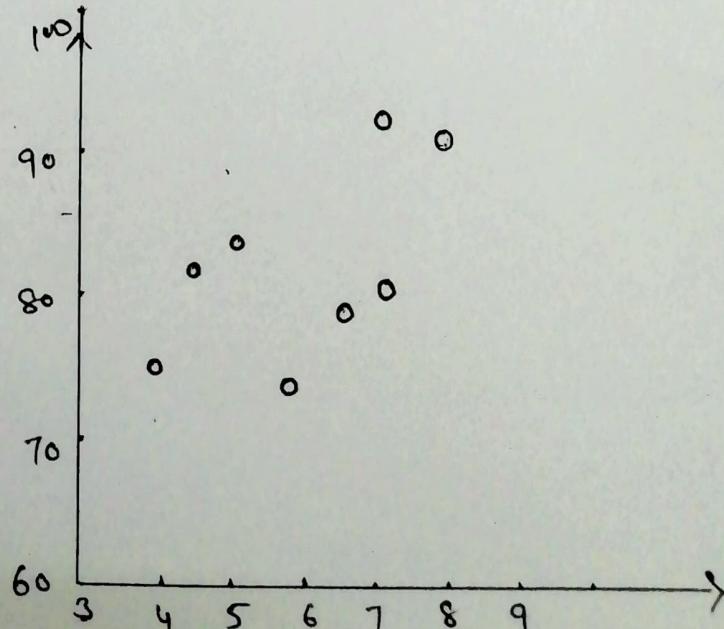
Data: The number of wins recorded by  
The India for seven seasons



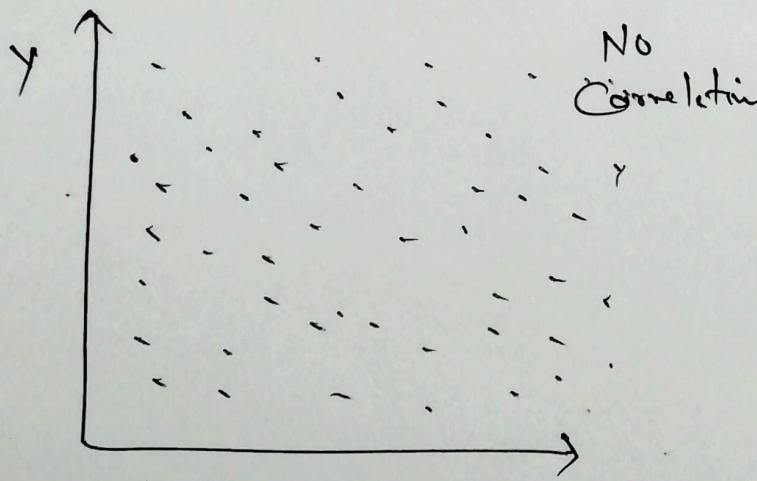
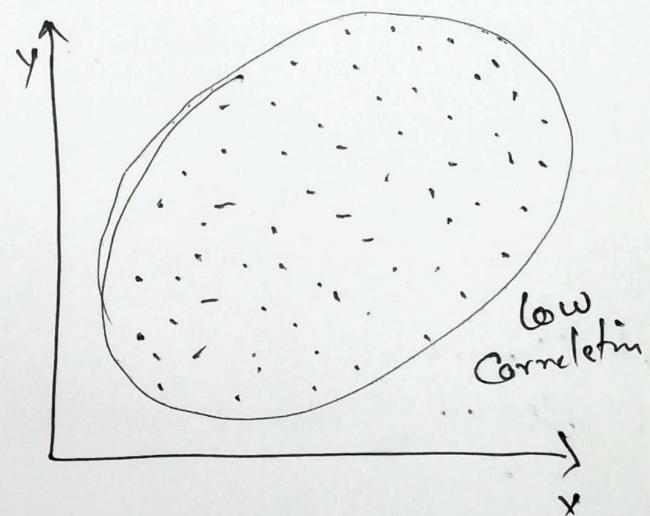
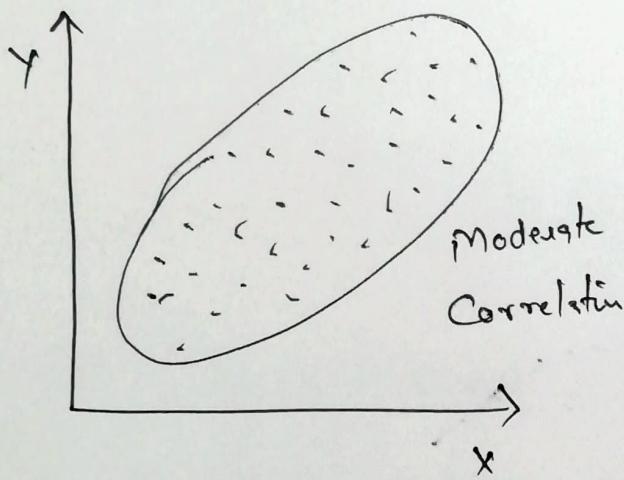
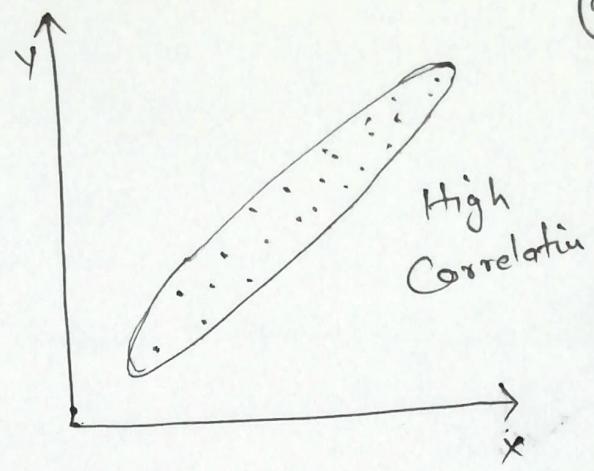
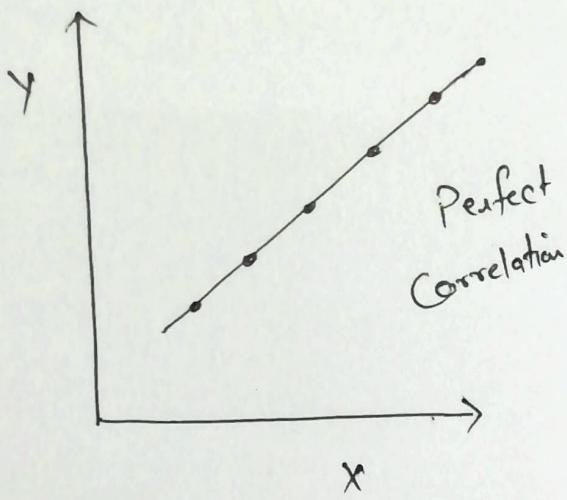
## Scatter Charts illustrate relationship between two variables

Data: The number of hours eight students studied  
for an exam and the scores they earned on the exam.

Study hours	Exam Score
5	84
7	92
4.5	82
7	80
8	90
6.5	78
5.5	74
4	75



9



+ve relation

-ve relation

# Calculating Descriptive Statistics

(10)

↳ Measures of Central Tendency

Center of data

↳ Mean, Median, Mode, Percentile

Mean.

$\bar{x} = (x_1, x_2, \dots, x_n)$  observations

Then

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

Let  $y = x + a$ ,  $a$  is any number

$$y_1 = x_1 + a$$

$$y_2 = x_2 + a$$

.....

$$y_n = x_n + a$$

Then

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^n (x_i + a)$$

$$= \frac{1}{n} \sum_{i=1}^n x_i + \frac{1}{n} \cdot na$$

$$\bar{y} = \bar{x} + a$$

Similarly  $y = x - a$  Then  $\bar{y} = \bar{x} - a$

Let  $y = ax$

(11)

$$y_1 = ax_1$$

$$y_2 = ax_2$$

- - -

$$y_n = ax_n$$

$$\text{Then } \bar{Y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^n ax_i$$

$$= a \frac{1}{n} \sum_{i=1}^n x_i$$

$$= a\bar{x}$$

Similarly  $y = x/a$  Then  $\bar{y} = \frac{1}{a}\bar{x}$

More Precisely

$$\text{Let } y = \frac{x - A}{h}$$

$$\text{Then } \bar{Y} = \frac{1}{h}x - \frac{A}{h}$$

$$\bar{Y} = \frac{1}{h}\bar{x} - \frac{A}{h}$$

$$h\bar{Y} = \bar{x} - A$$

$$\bar{x} = A + h\bar{Y}$$

Let us take data having  $N$  data points / observations

(12)

$x_1, x_2, x_1, x_3, x_4, \dots, x_m$

$x_2, x_1, x_3, x_m, x_1, \dots, x_2$

- - -

$x_3, x_4, x_m, x_3, x_1, \dots, x_m$

True in The data

$x_1$  is occurring, say  $f_1$  times

$x_2$  is occurring, say  $f_2$  times

- - -

$x_m$  is occurring, say  $f_m$  times

frequency

Therefore we can write data using frequency distribution

Data Points	frequency
$x_1$	$f_1$
$x_2$	$f_2$
- - -	- - -
$x_m$	$f_m$

$$N = \sum_{i=1}^m f_i$$

So  $\bar{x} = \frac{\text{Sum of all observations}}{\text{Total no. of Data Points}}$

$N$

$x_1 + x_1 + \dots + f_1$  times

+  $x_2 + x_2 + \dots + f_2$  times

- - -

+  $x_m + x_m + \dots + f_m$  times

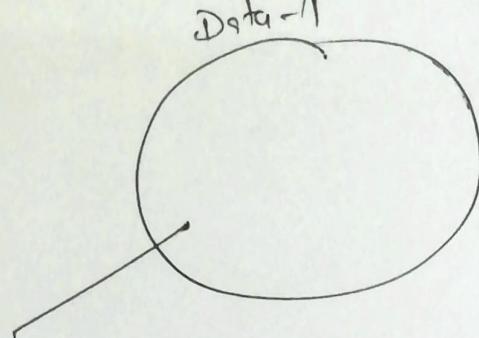
$$= f_1 x_1 + f_2 x_2 + \dots + f_m x_m$$

$$= \sum_{i=1}^m f_i x_i$$

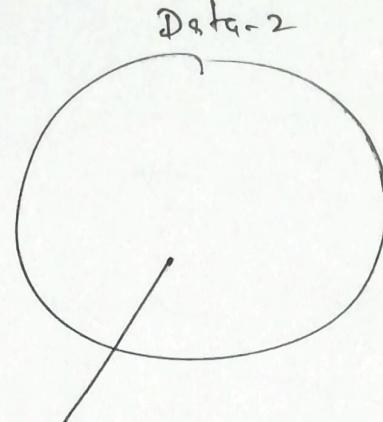
$$\bar{x} = \frac{1}{N} \sum_{i=1}^m f_i x_i$$

$$\text{where } N = \sum_{i=1}^m f_i$$

Composite Series / Combined Mean



$$x_{11}, x_{12}, \dots, x_{1n_1}$$



$$x_{21}, x_{22}, \dots, x_{2n_2}$$

$$\bar{x}_1 = \frac{1}{n_1} \sum_{i=1}^{n_1} x_{1i}$$

$$\bar{x}_2 = \frac{1}{n_2} \sum_{i=1}^{n_2} x_{2i}$$

$$\sum_{i=1}^{n_1} x_{1i} = n_1 \bar{x}_1$$

$$\sum_{i=1}^{n_2} x_{2i} = n_2 \bar{x}_2$$

Mean of data points together of data-1 and data-2

$$\bar{x} = \frac{\text{Sum of all data points}}{\text{Total no. of data points}}$$

$$\bar{x} = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2}{n_1 + n_2}$$

## Data.    Grouped data

(14)

Class/Interval	frequency	Take mid values.
under $a_1 - a_2$	$f_1$	$m_1 = \frac{a_1 + a_2}{2}$
$a_2 - a_3$	$f_2$	$m_2 = \frac{a_2 + a_3}{2}$
$a_3 - a_4$	$f_3$	$m_3 = \dots$
$\dots$		
$\dots - a_m$	$f_m$	$m_m = \dots$

$$\bar{x} = \frac{1}{N} \sum_{i=1}^m f_i m_i$$

where  $N = \sum_{i=1}^m f_i$

—————

## Weighted Mean

Averaging using different weights

Data Points	X	Weights
$x_1$		$w_1$
$x_2$		$w_2$
$x_n$		$w_n$

$$\bar{x} = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$$

Subject | Credit | Score/m

Subject	Credit	Score/m
A	4	70
B	4.5	65
C	5	85
D	6.5	60

CGPA — ?

Median

median value divides the data into two halves, one is having values greater than median value and the other having less than median value.

Sort in increasing / decreasing and then look at the middle value.

$$\begin{array}{c}
 \overline{\overline{n}} \\
 \left\{ \begin{array}{l} \text{odd} \\ \text{even} \end{array} \right. \rightarrow \begin{array}{l} \text{Median} \\ x \left[ \frac{n+1}{2} \right] \text{ The observation} \\ \hline x \left[ \frac{n}{2} \right] + x \left[ \frac{n}{2} + 1 \right] \\ \hline 2 \end{array}
 \end{array}$$

In grouped data.

Median class  $\frac{N}{2}$

look cumulative frequency

$$\text{Median value} = l + \frac{(N/2 - cf)}{f} \times h$$

lower value of median class      frequency of median class      cumulative frequency preceding the median class.

Mode : value that occurs more frequently  
in the given data

(16)

→ may not exist

if exist

→ may or may not be unique

→ always from the data points

→ least affected by extreme values

→ visualized graphically, using Histogram

In grouped data

Mode = L +

$$\frac{(f_1 - f_0) * h}{2f_1 - f_0 - f_2}$$

frequency process

The modal class

difference b/w  
the class

frequency in  
modal class

lower value of  
modal class

frequency next  
to modal class

## Partitioning The data

(17)

Arrange data least to greatest value

2 halves	Median	$\frac{N}{2}$
4 equal parts	Quartiles	$Q_i : i \frac{N}{4}, i = 1, 2, 3$ $Q_2 = \text{Median}$
10 equal parts	Deciles	$D_i : i \frac{N}{10}, i = 1, 2, \dots, 9$ $D_5 = Q_2 = \text{Median}$
100 equal parts	Percentiles	$P_i : i \frac{N}{100}, i = 1, 2, \dots, 99$ $P_{50} = D_5 = Q_2 = \text{Median}$