

► 1.1 ELEMENTS OF STRUCTURED DATA AND FURTHER READING

- Structured data is the data which conforms to a data model, it has well-defined structure and it follows a consistent order and can be easily accessed and used by a person or a computer program.
- Structured data is generally stored in well-defined schemas such as Databases.
- It is generally tabular with column and rows that define its attributes.
- Data conforms to a data model and has easily identifiable structure.

► 1.1.1 Example : (Characteristics)

- (i) Data is well-organised : Definition, Format and Meaning of Data is explicitly known.
- (ii) Data is kept in fixed fields within a record or file.
- (iii) Similar entities are grouped together to form relations or classes.
- (iv) Entities in the same group have same attributes.
- (v) Easy to access and query.
- (vi) Data elements are addressable, it is efficient to analyse and process.

► 1.1.2 Sources of Structured Data

- (i) SQL Databases
- (ii) Spreadsheets such as Excel
- (iii) OLTP systems
- (iv) Online forms
- (v) Sensors such as GPS or RFID tags.
- (vi) Network and web server logs
- (vii) Medical devices.

► 1.1.3 Advantages of Structured Data

- (i) Structured data have a well-defined structure and it helps in accessing of data.

- (ii) Data can be indexed based on text string as well as attributes. This makes search operation hassle-free.
- (iii) Data mining is easy, i.e. knowledge can be easily extracted from data.
- (iv) Due to its well structured form, operations such as updating and deleting is quite easy.
- (v) Business intelligence operations such as Data warehousing can be easily undertaken.
- (vi) Easily scalable if there is an increment of data.
- (vii) Ensuring security to data is easy.
- (viii) Structured data accounts for only about 25% of data but because of its high degree of organisation and performance make it the foundation of Big Data.

► 1.2 RECTANGULAR DATA

- Rectangular data are the staple of statistical and machine learning models.
- Rectangular data are **multivariate cross-sectional data** (i.e., not time-series or repeated measure).
- Here, each column is a variable (feature), and each row is a case or record.
- Multivariate cross-sectional data (i.e. not time-series or repeated measure) are indicated by rectangular data; Each column of it is a variable (feature) and each row is a record.

► 1.2.1 Procedure of Representation of Rectangular Data

- (i) First procedure is to map rectangle data on to a higher-dimensional point data and use point-based data structure procedures such as the **grid-file**, **PR quadtree**, **point quadtree** and **k-d-tree**.

Procedure mapping of the rectangular data to a four-dimensional point can be represented as : x and y coordinates of one corner and the width and height.

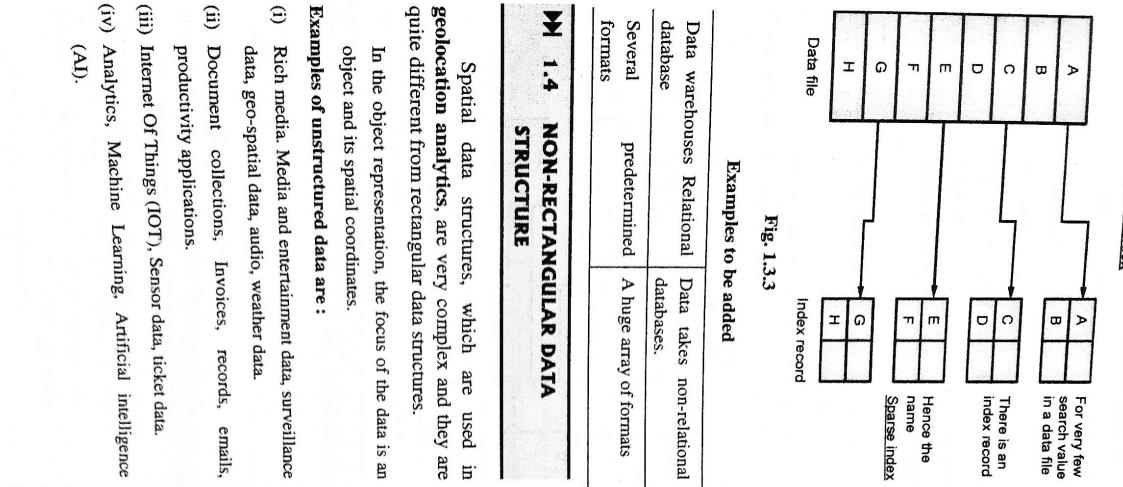


Fig. 1.3.3

(ii) Standard Deviation

Let f_i ($i = 1$ to n) be the frequencies of x_i ($i = 1$ to n).

Let M be the mean. If σ is standard deviation (SD) then,

$$\sigma^2 = \frac{\sum_{i=1}^n [(x_i - M)^2 f_i]}{\sum_{i=1}^n f_i}$$

$$\text{if, } r = 1, \mu'_r = \frac{\sum f_i x_i}{\sum f_i} = \text{Mean } M.$$

Again $\mu'_2 = \frac{\sum f_i x_i^2}{\sum f_i}$

$$\mu'_3 = \frac{\sum f_i x_i^3}{\sum f_i}$$

$$\mu'_4 = \frac{\sum f_i x_i^4}{\sum f_i}$$

Short-cut method to find σ :

$$\text{Let, } x' = \frac{x - x_0}{h}$$

$$\text{Then } \sigma^2 = h^2 \left[\frac{\sum f_i x_i'^2}{\sum f_i} - \left(\frac{\sum f_i x_i'}{\sum f_i} \right)^2 \right]$$

Note that, σ^2 is also called as Variance and denoted by $\text{Var}(x)$.

(iii) r^{th} moments about mean

Let, (x_i, f_i) be the given frequency distribution, then the 1^{st} moment about mean M is given by,

$$\mu_r = \frac{\sum f_i x_i (x_i - M)^{r-1}}{\sum f_i}$$

If $r = 1$, then $\mu_1 = 0$

$$\therefore \mu_1 = \frac{\sum f_i (x_i - M)}{\sum f_i}$$

$$= \frac{\sum f_i x_i}{\sum f_i} - \frac{\sum f_i M}{\sum f_i}$$

$$= \frac{\sum f_i x_i}{\sum f_i} - M$$

$$= \frac{\sum f_i x_i}{\sum f_i} - M \cdot \frac{\sum f_i}{\sum f_i}$$

$$= M - M \cdot 1 = 0$$

Also note that,

$$\text{If } x' = \frac{x - x_0}{h}$$

$$\text{Then, } \mu'_r = h^r \frac{\sum f_i x_i'^{r-1}}{\sum f_i}$$

(iv) Karl Pearson's coefficients of kurtosis

- (i) $\beta_1 = \text{Measure of skewness} = \frac{\mu_3}{\mu_2^2}$
- (ii) $\beta_2 = \text{Measure of flatness of single humped distribution}$

- (i) Rich media, Media and entertainment data, surveillance data, geo-spatial data, audio, weather data.
- (ii) Document collections, Invoices, records, emails, productivity applications.
- (iii) Internet Of Things (IOT), Sensor data, ticket data.
- (iv) Analytics, Machine Learning, Artificial intelligence (AI).

$$\text{Let, } x' = \frac{x - x_0}{h}; \text{ where } x_0 \text{ is assumed mean}$$

and h is length of the class interval

$$\text{Then } M = x_0 + hA$$

$$\text{Where, } A = \frac{\sum f_i x_i'}{\sum f_i}$$

and for, $r = 2$

$$\mu_2 = \text{Variance} = \frac{\sum f_i (x_i - M)^2}{\sum f_i}$$

$$= \sigma^2$$

= Square of standard deviation

If $\beta_2 < 3$, the distribution is flat compared to normal curve and is known as plato-kurtic.



Fig. 1.6.1

(vii) The expected value of x : (mean value of X) :

If X is a random variable then the expected value of X is denoted by $E(X)$ and means the value, on average, that X takes.

□ Definition : If $X = x_i$, $i = 1$ to n is a discrete random

variable with frequencies f_i , $i = 1$ to n , then

$$E(X) = \sum_{i=1}^n x_i f_i$$

Note : The expected value of X is also called as the Mean i.e. $M = E(X)$

Properties

(i) The r^{th} moment about origin is also written as,

$$\mu'_r = E(x^r) = \sum_{i=1}^n (x_i^r) f_i$$

Clearly, $E(x) = \mu'_1 = \sum_{i=1}^n x_i f_i$

$$E(x^2) = \mu'_2 = \sum_{i=1}^n x_i^2 f_i$$

$$E(x^3) = \mu'_3 = \sum_{i=1}^n x_i^3 f_i \text{ and}$$

$$E(x^4) = \mu'_4 = \sum_{i=1}^n x_i^4 f_i \text{ and so on.}$$

(ii) Moment about the mean \bar{x} are defined as,

$$\mu_r = E[x_i - \bar{x}]^r = \sum_{i=1}^n (x_i - \bar{x})^r f_i$$

and is called as r^{th} moment about mean \bar{x} .

$$\text{Clearly, } \mu_1 = 0$$

$$\mu_2 = E(x_i - \bar{x})^2 = \mu'_2 - \mu_1^2$$

$$\mu_3 = E(x_i - \bar{x})^3$$

$$\mu_4 = E(x_i - \bar{x})^4$$

$$= \mu'_3 - 3\mu'_2 \mu_1 + 2\mu_1^3$$

$$= 0.0600$$

$$= (0.0600)^3 - 3(0.4546)(0.0375) + 2(0.0375)^2$$

$$= 0.05074 - 3(0.0669)(0.0375)$$

$$= 3.5 + (0.5)(0.075) = -3.538$$

$$= 6(0.0375)^2 (0.4546) - 3(0.0375)^4$$

$$= 3.2385$$

$$= 32.22$$

(iii) Standard deviation : Using the result (ii) in section 10.3, we have

$$\sigma^2 = h^2 \left\{ \frac{\sum f x^2}{\sum f} - \left(\frac{\sum f x'}{\sum f} \right)^2 \right\}$$

$$= (0.5)^2 \left[\frac{582}{320} - \left(\frac{24}{320} \right)^2 \right] = 0.453$$

(iv) By definition of β_1 and β_2 we get,

$$\beta_1 = \frac{\mu_3}{\mu_2} = \frac{(0.0600)^2}{(0.0453)^2}$$

$$= 1578.21$$

Since $\beta_1 > 3$, the distribution is lepto-curtic i.e. it is peaked up sharply than the normal distribution.

UEX. 1.6.2 (Ret - Dec. 96)

From the following frequency distribution compute the standard deviation of 100 students.

Table P.1.6.2

Mass in kg.	No. of students
60-62	5
63-65	18
66-68	42
69-71	27
72-74	8

☒ Soln. : We construct the table,

Let $\bar{x} = 67$ be assumed mean
and $h = 3$
 $=$ class width

$$\text{Let } u = \frac{x-67}{3}$$

$$\mu_2 = \mu'_2 - \mu_1^2 = (0.4546) - (0.0375)^2 = 0.0453$$

Table P. 1.6.2(a)

Class of masses	Midpoint of Classes x	Number of students f	$\frac{x-u}{3\sigma}$	f ₀	f ₀ ²
60-62	61	5	-2	-10	20
63-65	64	18	-1	-18	18
66-68	67	42	0	0	0
69-71	70	27	1	27	27
Total (Σ)		100	0	15	97

We have $\sum f \cdot u = 15$, $\sum f u^2 = 97$

$$h = 3 \text{ and } N = \sum f = 100$$

By definition, $\sigma = h \sqrt{\frac{1}{N} \sum f u^2 - \left(\frac{1}{N} \sum f_u\right)^2}$

$$= 3 \sqrt{\frac{1}{100} (97) - \left(\frac{15}{100}\right)^2}$$

1.7 THE MEDIAN

To find its value, we use the formula,
 $M = l_1 + \frac{l_2 - l_1}{f_1} (n - C)$

M = Median
 l_1 = The lower limit of the class in which median lies
 l_2 = The upper limit of the class in which median lies
 f_1 = The frequency of the class in which the median lies
 m = Middle item and
 C = Cumulative frequency of the group preceding the median group

Let us take an example.

Ex. 1.7.1 : Find the median of the following distribution :

Class interval Rs.	Frequencies
1-3	6
3-5	53
5-7	85
7-9	56
9-11	21
11-13	16
13-15	4
Total	245

1.8 MODE

The word mode comes from the French word 'la mode' (which means the fashion). The mode is the observation which occurs most frequently in a set.

1.8.1 Calculation of Mode

- The median of a set of n observations x_1, x_2, \dots, x_n is the middle value when the observations are arranged in the ascending order of magnitude.
- If n is odd, the middle value which is $\left(\frac{n}{2} + 1\right)^{th}$ the the ascending order of magnitude is unique and is the median.
- If n is even, there are two middle values and the average of three values is the median.

For example, the median of the set 2, 3, 5, 6, 7 is 5 and that of the set -3, -1, 0, 1, 2, 3 is $\frac{0+1}{2} = 0.5$.

- The median is the value which divides the set of observations into two equal halves, such that 50% of the observations lie below the median and 50% above the median.

Table P. 1.7.1(a)

Class-interval	Frequency	Cumulative frequency
1-3	6	6
3-5	53	59
5-7	85	144
7-9	56	200
9-11	21	221
11-13	16	237
13-15	4	241
15-17	4	245

Now, median = The value of $\frac{N}{2}$ i.e.
 122.5 items, which lies in $5-7$ group,

Applying the formula,
 $M = l_1 + \frac{l_2 - l_1}{f_1} (n - C)$ we get

$$M = 5 + \frac{7-5}{85} (122.5 - 59) = 6.5$$

Table P. 1.8.1

Class-interval (cm)	Frequency
145-146	2
147-148	5
149-150	8
151-152	15
153-154	9
155-156	6
157-158	4
159-160	1
Total	50

Table P. 1.8.1

Δ_1 = Excess of the modal class frequency over the frequency of the class to its left and

Δ_2 = Excess of the modal class frequency over the frequency of the class to its right and
 h = Size of the class-intervals.

Ex. 1.8.1 : Calculate the mode for the given data :

Class-interval	Frequency
1-3	6
3-5	53
5-7	85
7-9	56
9-11	21
11-13	16
13-15	4
Total	245

Soln. : First we find the modal class i.e. the interval which corresponds to the maximum frequency. Here the modal class is 151-152. Class midpoint is 151.5. Its lower boundary is 150.5 and upper boundary is 152.5. Again $h = 2$, $\Delta_1 = 15 - 8 = 7$, $\Delta_2 = 15 - 9 = 6$.

$$\therefore M = 150.5 + \left(\frac{7}{7+6}\right) \times 2 = 150.57$$

Hence $N =$ Total frequency = 245
 $\therefore \frac{N}{2} = 122.5$

We prepare the table with cumulative frequency.

- For grouped data mode is determined as follows :

$$M = L_1 + \left(\frac{\Delta_1}{\Delta_1 + \Delta_2}\right) h$$

Where M is mode,

L_1 = Lower boundary of the modal class.

There are two other averages, the geometric mean and the harmonic mean which are sometimes used.

- The Geometric Mean (G.M.) of a set of observations is such that its logarithm is equal to the Arithmetic mean of the logarithms of the values of the observations. This

Table P. 1.8.1

Table P. 1.8.1

is given by the same formula even if the observations occur with certain frequencies.

Consider the set of values 2, 3, 5, 6 which occur with frequencies 10, 16, 24, 10 respectively. If x is the Geometric Mean, then,

$$\log x = \frac{10 \log 2 + 16 \log 3 + 24 \log 5 + 10 \log 6}{10 + 16 + 24 + 10}$$

$$= 0.5867$$

- (2) The harmonic mean (H.M.) of a set of observations is such that its reciprocal is the arithmetic mean (A.M.) of the reciprocals of the values of the observations.

Consider the above set of values. The harmonic mean of the set of values is given by,

$$\frac{1}{y} = \frac{10 \times \frac{1}{2} + 16 \times \frac{1}{3} + 24 \times \frac{1}{5} + 10 \times \frac{1}{6}}{60}$$

$$= 0.28$$

$$\therefore y = 3.57$$

Note :

1. The Geometric mean can be found only if the values assumed by the observations are positive.
2. It can be shown that $A.M. \geq G.M. \geq H.M.$

1.9.1 Other Measures of Location, Quartiles, Deciles and Percentiles

- We have seen that the median of a set of measurements is the value which divides the set into two equal halves, each containing 50% of the measurements.
- In the same way, some other measures of location can be considered. We define the three quartiles, Q_1 , Q_2 and Q_3 .
- They are such that when the measurements are arranged in increasing order, they divide the set of measurements into four equal parts, the first quartile Q_1 contains the 25% of the measurement, the second quartile Q_2 contains 50% of the measurements and the third quartile Q_3 contains 75% of the measurements.
- Actually the second quartile Q_2 is the median.
- Similarly we define the deciles. The first decile D_1

contains 10% of the measurements, the second decile D_2 contains 20% of the measurements and so on.

D_5 contains 50% of the measurements and so on.

The fifth decile is the median.

In the same manner, we define percentiles. The 99 frequencies P_1, \dots, P_{99} divide the set of measurements into 100 equal parts.

The first percentile P_1 contains 1% of the measurements, the second percentile P_2 contains 2% of the measurements and so on, the 12th percentile contains 12% of the measurements.

The 50th percentile is therefore the median. The method of finding out the quartiles, deciles and percentiles is basically the same as that of finding the median.

The median divides the set of observations into two equal values, each containing 50% of the measurements, the 3rd decile divides the set into two parts, the first part being 30% of the set and the other containing 70% of the observations.

The observation for the first quartile Q_1 corresponds to $\frac{280}{4} = 70^{\text{th}}$ observations, which lies in the interval 600-800, with lower class boundary 600. This interval contains 78 observations and the interval preceding this contains 66 observations. Hence,

Wages (Rs.)	Frequency	Cumulative Frequency
Less than 200	12	12
200-400	16	28
400-600	38	66
600-800	78	144
800-1000	80	224
1000-1200	35	259
1200-1400	14	273
Above 1400	7	280

Step I : To find quartiles

The observation for the first quartile Q_1 corresponds to $\frac{280}{4} = 70^{\text{th}}$ observations, which lies in the interval 600-800, with lower class boundary 600. This interval contains 78 observations and the interval preceding this contains 66 observations. Hence,

$$Q_1 = l_1 + \frac{l_2 - l_1}{f_1} (m - C)$$

where l_1 = Lower limit of the class in which Q_1 lies

l_2 = The upper limit of RM class in which Q_2 lies

f_1 = Positive frequency of the class

$m = \frac{N}{4}$

C = Cumulative frequency of the group preceding the Q_1 class.

$$\therefore Q_1 = 600 + \frac{200}{78} \left(\frac{280}{4} - 66 \right)$$

$$= 610.25 \text{ Rs.}$$

The median which is the second quartile Q_2 , is given by,

$$Q_2 = 600 + \frac{200}{78} \left(\frac{280}{2} - 66 \right)$$

$$= 600 + 189.74 = 789.74 \text{ Rs.}$$

The third quartile Q_3 is given by,

Soln.:

First we prepare cumulative frequency table.

Table P. 1.9.1(a)

$$Q_3 = 800 + \frac{200}{80} \left[280 \times \frac{3}{4} - 144 \right] = 965 \text{ Rs.}$$

Step II : The observation for 4th decile corresponds to the $\frac{280 \times 4}{10} = 112^{\text{th}}$ observation, which lies in the interval 600-800. Hence,

$$D_4 = 600 + \frac{200}{78} (112 - 66) = 717.95 \text{ Rs.}$$

Step III : The observation for 66th percentile corresponds to $280 \times \frac{66}{100} = 184^{\text{th}}$ observation which lies in the interval 800-1000. Thus the 66th percentile P_{66} is given by,

$$P_{66} = 800 + \frac{200}{80} (184.8 - 144) = 902 \text{ Rs.}$$

In the same way,

$$P_{10} = 400 + 0 = 400$$

$$\text{and } P_{90} = 1000 + \frac{200}{35} \left[\frac{280 \times 90}{100} - 220 \right] = 1182.96 \text{ Rs.}$$

1.10 THE RANGE

The range of set of numbers is the difference between the largest and the smallest items of the set. The range is a very crude measure. It does not tell us about the distribution of the values of the set relative to the average.

1.10.1 The Semi-Interquartile Range

This is a more refined form of range. It is defined by,

$$Q = \frac{Q_3 - Q_1}{2}$$

And is called as semi-interquartile range.

Q_1 = First quartile

Q_3 = Third quartile

M 1.11 THE MEAN DEVIATION

Ex. 1.11.2 : Calculate the mean deviation from the mean of the following distribution.

$$\text{M.D.} = \frac{\sum_{i=1}^n |x_i - \bar{x}|}{n}$$

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

is the arithmetic mean and $|x_i - \bar{x}|$ is the absolute value of the deviation.

Ex. 1.11.1 : Find the mean deviation of the set of measurement 1, 3, 8.

Soln.: Here the arithmetic mean,

$$\bar{x} = \frac{1+3+8}{3} = 4$$

$$\therefore \text{M.D.} = \frac{|1-4| + |3-4| + |8-4|}{3} = 2.67$$

1.11.1 Mean Deviation for Grouped Data

Let x_1, x_2, \dots, x_n occur with the corresponding frequencies f_1, f_2, \dots, f_n , then

$$\text{M.D.} = \frac{\sum_{i=1}^n f_i |x_i - \bar{x}|}{\sum_{i=1}^n f_i}$$

$$\text{where, } \bar{x} = \frac{\sum f_i x_i}{\sum f_i}$$

Note that the above formula is also applicable in the case of a frequency distribution whose class intervals have mid-points x_1, x_2, \dots, x_n and the classes have frequencies f_1, f_2, \dots, f_n .

Ex. 1.11.3 : The mean annual salary paid to all employees of a company was Rs. 5000. The mean annual salaries paid to male and female employees were Rs. 5200 and Rs. 4200 respectively. Determine the percentage of males and females employed by the company.

Ex. 1.11.4 : The first of the two samples has 100 items with mean 15 and standard deviation 3. If the whole group has 250 items with mean 15.6 and standard deviation $\sqrt{13.44}$, find the standard deviation of the second group.

Soln.: We have,

$$\begin{aligned} n_1 &= 100, \quad \bar{x}_1 = 15, \quad \sigma_1 = 3 \\ n &= n_1 + n_2 = 250 \\ \bar{x} &= 15.6 \\ \sigma &= \sqrt{13.44} \\ \therefore n_2 &= 250 - 100 = 150 \end{aligned}$$

$$\begin{aligned} \therefore \text{Percentages are } 80 \text{ and } 20. \\ \text{From Equations (1) and (2),} \\ n_1 = 80, n_2 = 20, \\ \therefore 26n_1 + 21n_2 = 2500 \\ \therefore 5000 = \frac{100}{100(9) + 150\sigma_2^2} \quad \dots(2) \\ \therefore 250 \times 13.44 = 900 + 36 + 40 + 150\sigma_2^2 \\ \therefore 13.44 = \frac{100(-0.0)^2 + 250(0.4)^2}{250} \\ \therefore \sigma_2 = 4 \end{aligned}$$

Soln.: We first calculate mean then find mean deviation.

Table P.1.11.2(a)

Marks	Number of students	Marks	Number of students
0-10	5	0-10	5
10-20	8	10-20	8
20-30	15	20-30	15
30-40	16	30-40	16
40-50	6	40-50	6
Total	50	Total	50

Here, Mean = $25 + \frac{10}{50} \times 10 = 27$ marks and

$$\text{Mean deviation} = \frac{\sum f_i |x_i - \bar{x}|}{\sum f_i}$$

$$= \frac{472}{50} = 9.44 \text{ marks}$$

(MU-New Syllabus w.e.f academic year 22-23) (M5-12B)

Tech-Neo Publications..A SACHIN SHAH Venture

Ex. 1.11.5 : Calculate the first four moments of the following distribution about mean and hence find β_1 and β_2 .

Soln.: We first calculate the first four moments about the point $x = 4$.

Table P.1.11.5

x	f	d = x - 4	fd	fd ²	fd ³	fd ⁴
1	8	-3	-24	54	-216	648
2	28	-2	-56	112	-224	448
3	56	-1	-56	56	-56	56
4	70	0	0	0	0	0
5	56	1	56	56	56	56
6	28	2	56	112	224	448
7	8	3	24	72	216	648
8	1	4	4	16	64	256
Σ	N = 256	-	0	512	0	2816

We prepare the table.

(Note the formula)

$$\begin{aligned} \text{Where } d_1 &= \bar{x}_1 - \bar{x} = 15 - 15.6 = -0.6 \\ \text{and } d_2 &= \bar{x}_2 - \bar{x} = 16 - 15.6 = 0.4 \\ \therefore \sigma &= \sqrt{13.44} \\ &= \sqrt{\frac{100(-0.0)^2 + 250(0.4)^2}{250}} \\ &= \sqrt{100(0) + 150\sigma_2^2} \\ &= \sqrt{150\sigma_2^2} \\ &= \sqrt{150} \sigma_2 \\ &= \sqrt{150} \times 4 \\ &= 12\sqrt{5} \end{aligned}$$

(MU-New Syllabus w.e.f academic year 22-23) (M5-12B)

Tech-Neo Publications..A SACHIN SHAH Venture

$$\mu'_r = \frac{1}{N} \sum f_i d^r$$

(Get acquainted with the notation)

where $d = x - \bar{x}$

$$\mu'_1 = \frac{1}{N} \sum f_i d$$

Now, $\mu'_1 = \frac{1}{N} \sum f_i d = 0$

$$\mu'_2 = \frac{1}{N} \sum f_i d^2 = \frac{512}{256} = 2$$

$$\mu'_3 = \frac{1}{N} \sum f_i d^3 = 0$$

$$\mu'_4 = \frac{1}{N} \sum f_i d^4 = \frac{2816}{256} = 11$$

Moments about mean are $\mu_1 = 0$

$$\mu_2 = \mu'_2 - \mu_1'^2 = 2 - 0 = 2$$

$$\mu_3 = \mu'_3 - 3\mu'_2 \mu_1' + 2\mu_1'^3$$

$$= 0 - 0 + 0 = 0$$

$$\mu_4 = \mu'_4 - 4\mu'_3 \mu_1' + 6\mu'_2 \mu_1'^2 - 3\mu_1'^4$$

$$= 11 - 0 + 0 - 0 = 11$$

$$\mu_2^2 = \frac{\mu_2}{\mu_2} = \frac{0}{0} = 0$$

$$\text{Now, } \beta_1 = \frac{\mu_2}{\mu_2} = \frac{11}{4} = 2.75$$

And $\beta_2 = \frac{\mu_2}{\mu_2} = \frac{11}{4} = 2.75$

► 1.12 ROBUST ESTIMATES

- Robust estimates** have been studied for the following problems.
- estimating location parameters
 - estimating scale parameters
 - estimating regression coefficients
 - estimation of model-states in models expressed in state-space form.

Example :

- (1) The mean is not a robust measure of central tendency.

If the given data set is the values {1, 3, 5, 8, 9}, then if we add another data point with the value - 2000 or + 2000 to the data, the resulting mean will be different to the mean of the original data.

Similarly, if we replace one of the values with a data point of -ve - 2000 or + 2000 then the resulting mean will be very different to the mean of the original data.

Robust statistical methods are used for many common problems, such as estimating location, scale and regression parameters.

One motivation is to produce statistical methods that are not affected by outliers.

Another motivation is to provide methods with good performance when there are small departures from a parametric distribution.

For example, robust methods gives good result for mixture of two normal distributions with different standard deviations.

(3) The median absolute deviation and interquartile range are robust measures of statistical dispersion, but the standard deviation and range are not.

(4) Trimmed estimators and Winsorised estimators are general methods to make statistics more robust.

L-estimators are a general class of simple statistics, but it is robust. And M-estimators are a general class of robust statistics. In this case solution is preferred even though it involves lot of calculations.

(i) by designing estimators so that a pre-selected behaviour of the influence function is achieved.

(ii) by replacing estimators that are optimal under the assumption of a normal distribution with estimators that are optimal for other distributions.

For example, using t-distribution with low degrees of freedom (i.e. high kurtosis, degrees of freedom between 4 and 6) or with a mixture of two or more distribution.

► 1.12.3 Measures of Robustness

The main tools used to describe robustness are : (i) The breakdown point, (ii) The influence function and (iii) The sensitivity curve.

► (i) The breakdown point

- The breakdown point of an estimator is the proportion of incorrect observations (e.g. arbitrarily large observations).
- An estimator can handle before giving an incorrect (e.g. arbitrary large) result. Usually the asymptotic (infinite sample) limit is quoted as the breakdown point. But the finite-sample breakdown point may be more useful.

For example, consider n independent random variables (X_1, X_2, \dots, X_n) and the corresponding realisations x_1, x_2, \dots, x_n , then we have $\bar{X}_n = \frac{x_1 + x_2 + \dots + x_n}{n}$, to estimate mean. Such as estimator has a breakdown point of 0 (or finite sample breakdown point of $\frac{1}{n}$). It is because we can make \bar{X} arbitrarily large by changing any of x_1, x_2, \dots, x_n .

(2) The median is a robust measure of central tendency. Considering the same dataset {1, 3, 5, 8, 9}, if we add another datapoint of value - 2000 or + 2000, then the median will change slightly, but it will be similar to the median of the original data.

• A breakdown point cannot exceed 50% because if more than half the of the observations are contaminated, it is not possible to distinguish between the underlying distribution. Therefore, the maximum breakdown point of 0.5. There are estimators which achieve such a breakdown point.

• Instead of relying on data, we can use the distribution of the random variables.

• Statistics with high breakdown points are called as resistant statistics.

► (ii) Influence Function

The influence function is also termed as the empirical influence function. It is a measure of the dependence of the estimator on the value of any one of the points in the sample.

It relies on calculating the estimator again with a different sample. Hence it is called as a model-free measure.

In mathematical terms, on influence function is defined as a vector in the space of the estimator, which is in turn defined for a sample which is a subset of the population.

► (iii) EIF

1. (Ω, A, p) is a probability space.

2. (X, Σ) is a measurable space (state space)

3. θ is a parameter space of dimension, $P \in N^*$.

4. $(f, S) = (R, B)$

The empirical influence function is defined as follows :

Let $n \in N^*$ and

$X_1, \dots, X_n : (\Omega, A) \rightarrow (X, \Sigma)$ are i.i.d.

And (x_1, \dots, x_n) is a sample from these variables.

$T_n : (X^n, \Sigma^n) \rightarrow (f, S)$ is an estimator. Let $i \in \{1, \dots, n\}$. The empirical influence function ($E I F_i$) at observation i is defined by

$E I F_i : x \in X \rightarrow n \cdot (T_n(x_1, \dots, x_{i-1}, x_{i+1}))$

Here, we are replacing the i -th value in the sample by an arbitrary value and looking at the output of the estimator.

Thus EIF is defined as the effect scaled by $(n + 1)$ instead of n , on the estimator of adding the point x to the sample.

► (iii) Sensitivity Curve

Also if we replace one of the values with a datapoint of value - 2000 or + 2000, then the resulting median will still be similar to the median of the original data.

- Here we are going to study, what happens to an estimator when we change the distribution of the data slightly :- it assumes a distribution, and measures sensitivity to change in this distribution.
- Here, the empirical influence assumes a sample set, and measures sensitivity to change in the samples.

1.13 ESTIMATES BASED ON PERCENTILES

- In a given dataset, the P^{th} percentile is a value such that at least P percent of the values or less than or at least $(100 - P)$ percent of the values are to be Considered.

For example, to find 70^{th} percentile, we sort the given data. We begin with the smallest data and then proceed 70 percent of the way to the largest value.

We note that the **median** is the same thing as the 50^{th} percentile.

- A common measure of variability is the difference between 25^{th} percentile and 75^{th} percentile. It is termed as Interquartile Range (or IQR).

Consider an example { 5, 3, 9, 7, 2, 1, 3, 6}. We sort the elements and get {1, 2, 3, 3, 5, 6, 7, 9} Now, the 25^{th} percentile is 2.5 and 75^{th} percentile is 6.5

\therefore interquartile range = $6.5 - 2.5 = 4$.

1.13.1 Percentile for Even Data

- If we have data containing even number of elements (even).

Here, we take any value between the order statistics x_i and x_{i+1} , where is satisfies:

$$100 \cdot \frac{i}{n} \leq p \leq 100 \left(\frac{i+1}{n} \right)$$

Also, the percentile for the weighted average :

Percentile (P) = $(1 - w)x_i + wx_{i+1}$ For some weight w between 0 and 1.

1.13.2 Exploring the Data Distribution

The estimates that we have covered so far indicate that the data can be expressed in a single number to describe its location or variability. It also tells us the nature of its distribution.

We mention below the main terms used for exploring the distribution.

- Boxplot** : Tukey has introduce Boxplot. It visualises the distribution of data quickly.

- Frequency Table** : It is the count of numeric data values in a set of intervals (bins).

- Histogram** : A plot of the frequency table with the given intervals on x-axis and the corresponding count on y-axis.

- Density Plot** : It is a smoothed version of the histogram.

1.13.3 Percentiles and Boxplots

- Another way of describing data is with a box plot.
- To construct the boxplot, we proceed as follows:

- Draw a rectangular box whose bottom is the lower quartile (25^{th} percentile) and whose top is the upper quartile (75^{th} percentile)

- Draw a horizontal line segment inside the box to represent the median.

- Extend horizontal line segments from each end of the box out to the most extreme observations.

Box plots can be either vertically drawn or horizontally drawn.

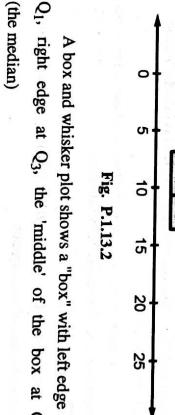


Fig. 1.13.2

A box and whisker plot shows a "box" with left edge at Q_1 , right edge at Q_3 , the 'middle' of the box at Q_2 (the median). The maximum and minimum as 'whiskers'.

The maximum and minimum as 'whiskers'. We observe that the plot divides the data into 4 equal parts.

If you score in the 75^{th} percentile, then 75% of population score lower than you.

Ex. 1.13.1 : Suppose the test score where : 22, 34, 68, 79, 79, 81, 83, 84, 87, 90, 92, 96 and 99. If your score was 75 , in what percentile did you score?

1.13.4 Outliers

Soln. : There were 14 scores reported and there were 4 scores at or below yours.

Now, $\frac{4}{14} \times 100\% = 29$

\therefore You scored in the 29^{th} percentile.

Ex. 1.13.2 : Find Q_1 , Q_2 and Q_3 for the following data set, and draw a box-and-whisker plot {2, 6, 7, 8, 8, 11, 12, 13, 14, 15, 22, 23}.

Soln. : The standard definition for an outlier is a number which is less than Q_1 or greater than Q_3 by more than 1.5 times the interquartile range (i.e., $1 \text{ Q.R.} = Q_3 - Q_1$).

Step I : There are 12 data points arranged: The middle two are 11 and 12 so the median i.e. Q_2 is 11.5.

The 'lower half' of the data set is the set {2, 6, 7, 8, 8, 11} The median of this lower half is 7.5; i.e. $Q_1 = 7.5$.

The 'upper half' of the data set is the set {12, 13, 14, 15, 22, 23}. The median here is 14.5; which is Q_3 . $\therefore Q_3 = 14.5$.

Step II : A box and whisker plot displays the values Q_1 , Q_2 and Q_3 along with the extreme values of the data set {2 and 23 in this problem}.

Step III : There are totally 15 values of arranged in increasing order.

Hence, Q_2 is the 8th data point.

Q_1 is 4th data point, $\therefore Q_1 = 4$ and Q_3 is 12th data point, $Q_3 = 23$.

Now, the interquartile range $I \text{ Q.R.} = Q_3 - Q_1$

$= 23 - 4$

$= 19$

Step IV : To find, if any, the values less than $Q_1 - (1.5 \times I \text{ Q.R.})$, or greater than $Q_3 + (1.5 \times I \text{ Q.R.})$.

Now, $Q_1 - (1.5 \times I \text{ Q.R.})$

$= 4 - (1.5 \times 10)$

$= 46 - 15$

$= 31$

If a data value is very far away from the quartiles (i.e. either much less than Q_1 or much greater than Q_3), it is called as an outlier.

Instead of being shown using the whiskers of the box-and whisker plot, outliers are generally shown as separately plotted points.

and $Q_3 + (1.5 \times IQR)$

$$= 56 + 15 = 71$$

Since 5 is less than 31 and 75 and 102 are greater than 71, hence there are 3 outliers.

we exhibit the box and whisker plot

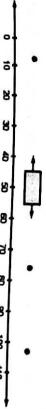


Fig. P. 1.13.3

1.13.5 Density Plots and Estimates

A density plot is a representation of the distribution of a numeric variable. It uses a kernel density estimate to show the probability density function of the variable.

It is the smoothed version of the histogram and is used in the same concept.

Interpretation of Density Curve

- (1) If a density curve is a left skewed, then the mean is less than the median
- (2) If a density curve is a right skewed, then the mean is greater than the median.
- (3) If a density curve has no skew, then the mean is equal to the median.
- (4) The peaks of density plot help display where values are concentrated over the interval.

An advantages of density plot over Histograms is that they are better at determining the distribution shape because they are not affected by the number of bins.

(5) Difference between density plot and histogram is as follows : Histogram provides the visual representation of data distribution. By using the histogram we can represent a large amount of data, and its frequency.

Density plot is the continuous and smoothened version of the histogram estimated from the data. It is estimated through Kernel Density Estimation.

(6) Data density in terms of sampling is how many items of particular information set are examined. Data density affect accuracy. Higher densities mean more accuracy.

Qualitative variates can also be divided into two types :

- (i) They may be continuous, if they can take any value we specify in some range or
- (ii) Discrete if their values change by steps or jumps.

Definition (5)

A continuous variate is a variate which may take all values within a given range.

Definition (6)

A discrete variate is a variate whose values change by steps. Births, 1990. The aim was to improve the survival rate and the care of the babies soon after birth by collecting the data on new born babies. The aim is to gain information about a population, namely 'all Indian births'.

Definition (7)

The frequency distribution of a (discrete) variate is the set of possible values of the variate, together with the associated frequencies.

Definition (8)

The frequency distribution of a (continuous) variate is the set of class-intervals for the variate, together with the associated class-frequencies.

If we classify the whole population according to birth-weights, then instead of looking at the frequency of each variate, we first group the values into intervals, which is the sub-division of the total range of possible values of the variate.

In this example, the variate may be classified as 1-500, 500-1000, 450-5000, 500-1-500 grams.

Definition (9)

A class-interval is a sub-division of the total range of values which a (continuous) variate may take.

Definition (10)

The class-frequency is the number of observations of the variate which fall in a given interval.

Definition (11)

Cumulative frequency distribution is the sum of all observations which are less than the upper boundary of given class interval : or this number is the sum of the frequencies upto and including that class to which the upper class boundary corresponds.

Definition (12)

Points to note while constructing the Tables.

1. Make the table self-explanatory provide a title, a brief description of a source of the data, State in what units the figures are expressed, label rows and columns where appropriate.
2. Keep the table as simple as possible.
3. Distinguish between zero values and missing observations.
4. Make alternations clearly.
5. Give the calculations a logical pattern on the sheet.

Table 1.14.1 : Cumulative Frequency (more than) Table

Class (cm) Interval	Frequency	Cumulative frequency more than
143-146	2	2
149-150	8	15
151-152	15	30
147-148	5	7
153-154	9	39
155-156	6	45
157-158	4	49
159-160	1	50
Total	50	

Table 1.14.2 : Cumulative frequency (less than) Table

Class (cm) Interval	Frequency	Cumulative frequency less than
145-146	2	50
147-148	5	48
149-150	8	43
151-152	15	35
153-154	9	20
155-156	6	11
157-158	4	5
159-160	1	1
Total	50	

1.15 HISTOGRAM

1.15.2 Example

- Histogram is commonly used device for charting continuous frequency distribution.

- It consists in erecting a series of adjacent vertical rectangles on the section of the horizontal axis (X - axis), with bases (sections) equal to the width of the corresponding class intervals and heights are so taken that the areas of rectangles are equals to the frequencies of the corresponding classes.

1.15.1 Construction of Histogram

- The variate values are taken along X-axis and the frequencies along Y-axis.

Case (I) Histogram with equal classes

- If classes are of equal magnitude, each class interval is drawn on X-axis by a section which is equal to the magnitude of the class interval.

- On each class interval erect a rectangle with the height proportional to the corresponding frequency of the class.

- The series of adjacent rectangles (one for each class) so formed gives the histogram of the frequency distribution and its area represents the total frequency of the distribution.

Case (II) : Histogram with unequal classes

- If the classes are not uniform, then the different classes are represented on X-axis by sections which are equal to the magnitude of the corresponding classes and the heights of the corresponding rectangles are to be adjusted so that the area of the rectangle is equal to the frequency of the corresponding class.
- This can be done by taking the height of each rectangle equal to the corresponding frequency density of each class, where,

$$\text{Frequency density of a class} = \frac{\text{Frequency of the class}}{\text{Magnitude of the class}}$$

- Ex. 1.15.1 : Represent the adjoining distribution of marks of 100 students in the examination by a histogram.

- Ex. 1.15.2 : Represent the following data by means of a histogram:

Wages (Rs.)	10-15	15-20	20-25	25-30	30-40	40-50	50-60
Number of workers	7	19	27	15	12	12	8

- Ex. 1.15.3 Exercise

- Describe briefly the construction of histogram of a frequency distribution.

Prepare a histogram from the following data :

- Here the class-intervals are of unequal magnitude, the corresponding frequencies have to be adjusted to obtain the 'frequency-density', so that area of the rectangle is equal to the class frequency.
- We note that the first four classes are of magnitude 5, the class 30-40 is of magnitude 10 and the last two classes 40-60 and 60-80 are of magnitude 20.
- Since 5 is the minimum class interval, the frequency of the class 30-40 is divided by 2 and the frequencies of classes 40-60 and 60-80 are to be divided by 4 as shown :

Marks	Number of students
0-10	4
10-20	6 - 4 = 2
20-30	24 - 6 = 18
30-40	46 - 24 = 22
40-50	67 - 46 = 21
50-60	86 - 67 = 19
60-70	96 - 86 = 10
70-80	99 - 96 = 3
80-90	100 - 99 = 1

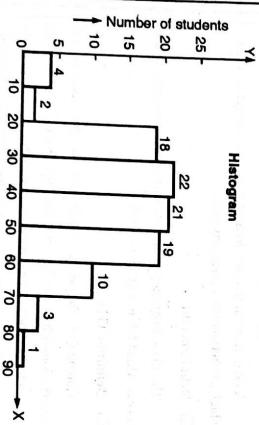


Fig. P. 1.15.1

- Ex. 1.16.1 Steps for construction of Pie-Diagram

Just as rectangles are used to represent the total magnitude and its various components, the circle may be divided into various sections or segments representing certain proportion or percentage of the various component parts to the total. Such a sub-divided circle diagram is known as an angular or pie diagram.

- Ex. 1.16.2 Example

Draw a pie-diagram to represent the following data of proposed expenditure by a state-govt. for the year 1947-98.

Weekly wages (Rs.)	Number of workers	Magnitude of class	Height of rectangle
10-15	7	5	7
15-20	19	5	19
20-25	27	5	27
25-30	15	5	15
30-40	12	10	$\frac{12}{10} = 6$
40-60	12	20	$\frac{12}{20} = 3$
60-40	8	20	$\frac{8}{20} = 4$

- (i) Express each of the component values as a percentage of the respective total.
- (ii) Since the angle at the centre of the circle is 360° , the total magnitude of the various components is taken to be equal to 360° .

The degrees represented by the various component parts of a given magnitude can be obtained as follows :

$$\text{Degree of any component part} = \frac{\text{Component value}}{\text{Total value}} \times 360^\circ$$

- (iii) Pie-diagram is also known as circular diagram

- Ex. 1.16.2 Example

Draw a pie-diagram to represent the following data of proposed expenditure by a state-govt. for the year 1947-98.

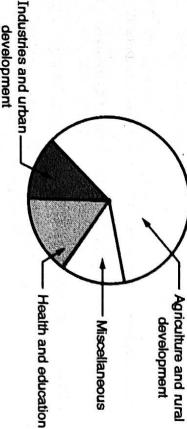
- order is not the same as categorical data, e.g. 'one', 'two', 'three'.

- But the sorting of these variables uses logical order for example, gender is a categorical variable and has categories male and female. And there is no ordering to assign categories.

Soln. :

E^{SP} Calculation of pie-chart

Items	Agriculture and rural development	Industries and urban development	Health and education	Miscellaneous
Proposed expenditure (in million Rs.)	4200	1,500	1000	500
	(1)	(2)	(3) $\frac{(2)}{7200} \times 360^\circ$	
Agriculture and Rural development	4,200	$4200 / 7200 \times 360^\circ = 210^\circ$		
Industries and urban development	1500	$1500 / 7200 \times 360^\circ = 75^\circ$		
Health and education	1000	$1000 / 7200 \times 360^\circ = 50^\circ$		
Miscellaneous	500	$500 / 7200 \times 360^\circ = 25^\circ$		
Total	7200	360°		



E^{SP} Pie-diagram representing proposed expenditure

- To study multivariate analysis, we come across following terms:
- (1) Contingency table :** In statistics, a contingency table (also known as a cross tabulation or crosstab) is a type of table in a matrix format that displays the (multivariate) frequency distribution of the variables. They are used in survey research, business intelligence, engineering and scientific research. A crucial problem of multivariate statistics is finding the direct-dependence structure of the variables contained in the contingency table. If some of the conditional independence are revealed, then the storage of the data can be done in a better way. To achieve this, the relative frequencies from the contingency table are used.
- Ex. 1.16.1 : Considered two variables, sex (male or female) and handedness (right or left handed).**
- Soln. :**
 - We select 100 individuals randomly from a large population to study of sex differences in handedness.
 - From the table (contingency table).

(MU-New Syllabus w.e.f academic year 22-23) (M5-128)

- A contour plot is a curve along which the function of two variables, has a constant value.
- It is a plan section of the three dimensional graph of the function $f(x, y)$ parallel to the x, y plane.
- A Contour line joints points of equal elevation (height) above a given level.
- The contour interval of a contour map is the difference in elevation between successive Contour lines.
- The most common form is the rectangular contour plot, which is shaped like a rectangle.

1.16.5 Hexagonal Binning

Hexagonal binning is a plot of two numeric variables with the records binned into hexagons.

Hexagonal binning is another way to manage the problem of having too many points that start to overlap.

Hexagonal binning plot density, rather than points.

Points are binning into gridded hexagons with distribution (the number of points per hexagon) is displayed using either colour or the area of the hexagons.

Ex. 6.16.2 : Convert binary number 1101010 into hexadecimal number.

Soln. :

- Converting the given number into decimal number.

$$(1101010)_2 = 1 \times 2^6 + 1 \times 2^5 + 0 \times 2^4 + 1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 0 \times 2^0$$

$$= 64 + 32 + 0 + 8 + 0 + 2 + 0$$

$$= 106$$

- Hexadecimal number :

$$(106)_{10} = 6 \times 16^1 + 10 \times 16^0$$

$$= (6A)_{16}$$

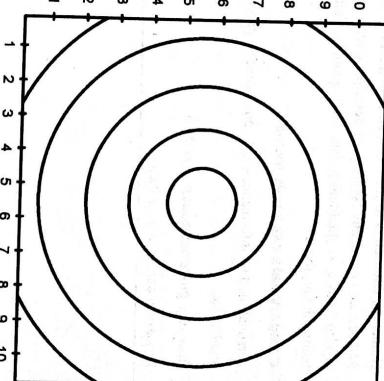


Fig. Ex. 6.16.2

Rectangular contour plot showing symmetrical surface with a Central peak.

1.16.7 Two Category Variables

- If two category variables are independent, then the value of the variable does not change the probability distribution of the other.
- If two category variables are related, then the distribution of one depends on the other.

(MU-New Syllabus w.e.f academic year 22-23) (M5-128)

1.16.8 Chi-Square Test of Independence

- This test is used to determine if category variables are independent or if they are in fact related to one another.
- This test measures the difference in the observed conditional distribution of one variable across levels of the other and compares it to the distribution of that variable.

1.16.9 Categorical and numerical data

- Category data refers to a data type that can be stored and identified based on the names or labels given to them.
- Numerical data refers to the data that is in the form of numbers and not in any language or descriptive form. It is also known as qualitative data as it qualifies data before classifying it.

Category data examples include

- (i) Personal biodata information- full name, gender, phone number etc.

Race, age group and educational level

Numerical data examples include

Module 2

Data and Sampling Distributions

CHAPTER 2

Syllabus :

- 2.1 Random Sampling and Sample Bias, Bias, Random Selection, Size Versus Quality, Sample Mean Versus Population Mean, Selection Bias, Regression to the Mean, Sampling Distribution of a Statistic, Central Limit Theorem, Standard Error, The Bootstrap, Resampling Versus Bootstrapping.
 2.2 Confidence Intervals, Normal Distribution, Standard Normal and QQ-PLOTS, Long-Tailed Distributions, Student's t-Distribution, Binomial Distribution, Chi-Square Distribution, F-Distribution, Poisson and Related Distributions, Poisson Distributions, Exponential Distribution, Estimating the Failure Rate, Weibull Distribution.

2.1 Sampling.....2-3

2.1.1 Definitions.....2-3

2.1.2 The Purpose of Sampling

2-3

2.1.3 Random Sampling.....2-3

2.1.4 Stratified sampling.....2-4

2.1.5 Opportunity sampling.....2-4

2.1.6 Systematic Sampling.....2-4

2.1.7 Cluster sampling.....2-5

2.1.8 Clusters of Different sizes.....2-5

2.1.9 Applications of Cluster Sampling.....2-5

2.1.10 Advantages of Cluster Analysis.....2-5

2.1.11 Disadvantages of cluster analysis.....2-6

2.1.12 Multi-stage Cluster Sampling.....2-6

2.1.13 Advantages of multi-stage cluster sampling.....2-6

2.1.14 Disadvantages of Multi-Stage Cluster Sampling.....2-6

2.2 Size Vs Quality.....2-7

2.3 Sample Mean Vs Population Mean.....2-7

2.4 Selection bias.....2-8

2.4.1 Types of Selection Bias.....2-8

2.5 Regression to the mean sampling distribution of a statistic.....2-9

2.12 Introduction TO Standard Discrete Distributions...2-29

2.6 Central Limit theorem.....2-10

UO. State central limit theorem and explain.

[MU - 2011]

UO. Explain the central limit theorem.

[MU - O 1(b); Dec. 14, O 3(b); May 15, O 1(a); Dec. 15, O 1(c); Dec. 16, O 1(b); Dec. 17, O 1(a); May 19, 2015 Marks]

2.6.1 Examples on Central Limit Theorem.....2-13

2.7 Confidence Intervals

2.8 Normal Distribution (o Gaussian Distribution).....2-17

UO. Discuss normal distribution and its characteristics. Or

State importance of normal distribution.

[MU-O 4(c); Dec. 17, O 6(b); Dec. 15, Dec. 16, O 6(c); May 19, 05 Marks]

2.8.1 Characteristics of Normal Distribution.....2-18

2.8.2 Properties of Normal Distribution (N.D.).....2-18

UO. What is the importance of standard normal variate ? [MU - Dec. 2020]

2.8.3 Mean, Variance and Mode of N.D.....2-19

2.8.4 Moment Generating Function.....2-21

2.8.5 Mean Deviation about Mean.....2-23

2.9 Solved examples Normal Distribution.....2-23

2.10 Standard Normal and QQ-plots.....2-27

2.11 Long tailed distributions.....2-28

2.12 Introduction TO Standard Discrete Distributions...2-29