

Descriptive statistics

Date 11/1/25
Page

Types of Descriptive Statistics

Distribution

Measures frequency of each value in a data set.

Central Tendency

Mean, Median, Mode

Variability

standard deviation, Range, Variance

* Data/statistical Variable -

A collection of actual observations or scores in a survey or an experiment.

Any statistical analysis is performed on data.

* Qualitative Data

1) Deals with descriptions.

2) Data can be observed but not measured.

3) colors, textures, smells, tastes, appearance, beauty etc.

4) length, height, area, volume, weight, speed, time, temp., humidity, sound levels, cost, members, ages etc.

4) Qualitative \rightarrow Quality

4) Quantitative \rightarrow Quantity

* Frequency Distributions

Represents

How many times that observations occurs is called as frequency.

Represents the pattern of how frequently each value of a variable appears in a dataset. It shows the no. of occurrences for each possible value within the dataset.

Ungrouped frequency distribution.

Q. Make the frequency distribution table for the ungrouped data given as follows: 10, 20, 15, 25, 30, 10, 15, 10, 25, 100, 15, 10, 30, 25.

~~mito~~ \Rightarrow Ungrouped frequency distribution table:

Value	anilasM frequency	anilasM
10	4	32
15	3	
20	2	
25	3	16
30	2	10

Dekb as begannig si gepland bestriede. vaf

Grouped frequency distribution.

Q. Make the frequency distribution table for the ungrouped data given as follows:

23, 29, 21, 14, 43, 37, 38, 41, 55, 11, 35, 15, 21, 24, 57, 35,
29, 10, 39, 42, 27, 17, 45, 52, 3, 36, 39, 38, 43,
46, 32, 37, 25

Min. no. of pips and max. no. of pips

10 include & ~~frequency~~ Class Interval

~~to exclude~~

5

~~67.150000~~ ← 20.530 ~~61.110000~~

15

30-40

1

40 - 50

6) comparative superlative

3

Q. Consider a data set of 26 children of ages 1-6 years.

→ Ungrouped frequency distribution

Age	1-2	3-4	4-5	5-6	6-7
frequency	5	3	6	5	4

→ Grouped frequency distribution

Age group	1-2	3-4	5-6
frequency	8	11	6

* Relative frequency distribution

Relative frequency = frequency of event

Total Number of Events

This distribution displays the proportion or % of observations in each interval or class.

It is useful for comparing different data sets or for analysing the distribution of data within the set.

Q. find relative frequency distribution.

Score Range 10-20 21-40 41-60 61-80 81-100

frequency 5 10 20 10 5

Relative frequency $\frac{5}{50} = 0.10$ $\frac{10}{50} = 0.20$ $\frac{20}{50} = 0.40$ $\frac{10}{50} = 0.20$ $\frac{5}{50} = 0.10$

$0.10 + 0.20 + 0.40 + 0.20 + 0.10 = 1.00$

* Cumulative frequency distribution

It is defined as the sum of all the frequencies in the previous values or intervals up to the current one.

The distributions which represent the frequency distributions using cumulative frequencies are called cumulative frequency distributions.

There are two types of cumulative frequency distributions.
Less than type: We sum all the frequencies before the current interval.

More than type: We sum all the frequencies after the current interval.

Q.

The table below gives the values of runs scored by Virat Kohli in the last 25 T-20 matches. Represent the data in the form of less-than-type cumulative frequency distribution.

45

34

50

51

54

56

57

58

59

56

63

65

66

67

68

69

70

71

70

71

72

73

74

75

76

77

78

72

73

74

75

76

77

78

79

80

75

76

77

78

79

80

81

82

83

⇒

Runs

frequency

0-10

2

0-10

10-20

02

03

01

02

01

02

0-10

20-30

01

02

01

02

01

02

30-40

4

40-50

01

02

01

02

01

50-60

01

02

01

02

01

60-70

5

70-80

01

02

01

02

01

80-90

01

02

01

02

01

90-100

1

01

02

01

02

01

Runs scored by Virat Kohli	C.F	Runs scored by Virat Kohli	Cumulative frequency (CF)
less than 10	2	more than 10	25
less than 20	4	more than 20	23
less than 30	5	more than 30	21
less than 40	9	more than 40	20
less than 50	13	more than 50	16
less than 60	18	more than 60	12
less than 70	19	more than 70	6
less than 80	22	more than 80	3
less than 90	24	more than 90	1
less than 100	25		

Measures of Central Tendency

→ Mean = sum of all obs. / no. of obs. (denoted by \bar{x})

→ Mode (or) most frequent observation

→ Median (middle observation)

NOTE: If the data is grouped data then each x_i is the class representative

If the data is grouped data then each x_i is the class representative of that class i.e. $x_i = \frac{x_i^l + x_i^h}{2}$ where, x_i^l is lower bound and x_i^h is upper bound for i^{th} class.

* Mean (Arithmetic Average)

for ungrouped data: $\bar{x} = \frac{\text{sum of all observations}}{\text{no. of observations}}$

Mean = Sum of all observations / no. of observations

Sum of all observations = $\sum x_i$

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

for grouped data: $\bar{x} = \frac{\sum f_i x_i}{\sum f_i}$

where, f_i = frequencies of class i and x_i = midpoint of interval.

case of median, to find the median class if $\frac{N}{2} = CF$ then
 the respective class interval of CF is the median class.
 $\frac{N}{2}$ is not equally equal to CF then find nearly CF
 for the value of $N/2$ and i.e. the Median class.

Median = The middle value of the dataset when arranged
 in ascending order.

for ungrouped data,

firstly, arrange the data in ascending order.

① If n is odd, Median = $\frac{(n+1)^{\text{th}} \text{ observation}}{2}$

② If n is even, Median = $\frac{(n/2)^{\text{th}} + (n/2 + 1)^{\text{th}} \text{ observation}}{2}$

for grouped data,

Step 1: Construct the cumulative frequency distribution.

Step 2: find the median class. Median class is the class
 in which the value of $\frac{N}{2}$ falls in cumulative
 frequency distribution.

Step 3: find the median by using the following formula.

$$\text{Median} = L + \left(\frac{\frac{N}{2} - C.F.}{f} \right) \times h$$

where; N = Total frequency

L = lower limit of median class

f = frequency of median class

C.F. = cumulative frequency of the class before
 the median class.

$h = \text{class width}$

Finding the modal class, if more than one class interval have

fi then the modal class will be the class interval which
 largest lower limit.

10-20 20

20-30 15

30-40 20

30, 30-40 is our modal class

and 30 is the lower limit of the modal class.

→ If there is one mode, the data is unimodal

→ If there are two modes, the data is bimodal

→ If there are more than two modes, the data is multimodal

→ If no value repeats, there is no mode

* Mode:

Value that occurs ^{most} frequently in data set.

for ungrouped data:

Mode = number that occurs the highest number of times.

for grouped data:

Step 1: find the modal class. Modal class is the class with maximum frequency.

$$\text{Mode} = L + \left[\frac{f_m - f_1}{2f_m - f_1 - f_2} \right] \times h$$

where, L = lower limit of the modal class

f_m = frequency of the modal class

f_2 = frequency of the class succeeding the modal class

f_1 = frequency of the class preceding the modal class

* Variance:

let x_1, x_2, \dots, x_n are 'n' observations of a variable x .

Then the variance of this variable,

for ungrouped.

$$\text{Var}(x) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

for grouped data,

$$\text{Var}(x) = \frac{\sum_{i=1}^n f_i (x_i - \bar{x})^2}{\sum_{i=1}^n f_i}$$

standard deviation (σ)

The square root of variance, which gives the spread of data points.

$$\text{S.D} = \sqrt{\frac{\sum f_i (x_i - \bar{x})^2}{n}} \quad (\bar{x} = \bar{X})$$

Coefficient of variation:

A standardized measure of dispersion calculated as the ratio of the standard deviation to the mean.

$$C.V = \frac{\text{Standard deviation}}{\bar{x}} \times 100$$

find Variance, standard deviation and coefficient of variation, class interval frequency, find part = sum

0-10

10-20

20-30

30-40

5

3

2

1

find mean, mode, median, variance, standard deviation, coefficient of variation.

class interval

frequency (f_i)

10-20

3

20-30

5

30-40

2

40-50

1

\sum

midpoint (x_i) = $\frac{\text{lower limit} + \text{Upper limit}}{2}$

classmate

Date _____

Page _____

class interval	f_i	x_i	CF	$(x_i - \bar{x})^2$	$f_i(x_i - \bar{x})^2$
0-10	4	5	4	272.25	1089
10-20	6	15	10	42.25	253.5
20-30	5	25	15	12.25	61.25
30-40	3	35	18	182.25	546.75
40-50	2	45	20	502.25	1045

$$\bar{x} = \frac{\sum f_i x_i}{\sum f_i} = \frac{20+90+125+105+90}{20+10} = 21.5$$

$$21.5 = \frac{20+90+125+105+90}{20+10} = 21.5$$

$$\text{for Mean, } \bar{x} = \frac{\sum f_i x_i}{\sum f_i} = \frac{20+90+125+105+90}{20+10} = 21.5$$

$$21.5 = \frac{20+90+125+105+90}{20+10} = 21.5$$

for Median,

$$\text{Median} = L + \left(\frac{N}{2} - C.F \right) \times h$$

$$L = 10, N = 20 \Rightarrow N/2 = 10, C.F = 4, h = 10$$

$$f = 6.$$

$$\text{Median} = 10 + \left(\frac{10-4}{6} \right) \times 10 = 20$$

$$\text{for Mode, } \text{Mode} = L + \left(\frac{F_m - F_1}{2F_m - F_1 - F_2} \right) \times h$$

$$L = 10, F_m = 6, F_1 = 4, F_2 = 5, h = 10$$

$$\text{Mode} = 10 + \left(\frac{6-4}{12-4-5} \right) \times 10 = \frac{50}{3} = 16.67$$

$$\text{Var}(x) = \frac{1}{n} \sum_{i=1}^n f_i (x_i - \bar{x})^2 = \frac{30.55}{20} = 152.75$$

$$\sum_{i=1}^n f_i (x_i - \bar{x})^2 = \sum_{i=1}^n f_i (x_i - 21.5)^2 = \sum_{i=1}^n f_i (x_i - 21.5 + 10.5)^2 = \sum_{i=1}^n f_i (x_i - 21.5)^2 + 10.5^2 \sum_{i=1}^n f_i = 30.55 + 10.5^2 \times 20 = 30.55 + 1102.5 = 1133$$

Standard deviation, $S.D = \sqrt{\text{Var}} = 12.35$	1
$\sum f_i^2 = 81$	2
Coefficient of variation $= \frac{S.D}{\bar{X}} \times 100 = \frac{12.35}{21.5} \times 100 = 57.44$	3
Class interval f_i x_i cf $(x_i - \bar{x})^2$ $f_i(x_i - \bar{x})^2$	4
20-30 3 15 3 216.09 648.27	5
\rightarrow 20-30 5 25 8 22.09 110.45	6
\rightarrow 30-40 2 35 15 128.09 196.63	7
40-50 2 45 17 234.09 468.18	8
	17
for mean, $\bar{x} = \frac{\sum f_i x_i}{\sum f_i} = \frac{45 + 125 + 245 + 90}{29.70} = 29.70$	18
for median, median = $L + \left(\frac{N}{2} - CF \right) \frac{h}{f} = 30 + \left(\frac{20}{2} - 17 \right) \frac{10}{5} = 34$	19
$L = 30, N = 20, CF = 17, f = 5, h = 10$	20
median = $30 + \left(\frac{17}{2} - 17 \right) \frac{10}{5} = 31$	21
for mode, mode = $L + \left(\frac{f_m - f_1}{2f_m - f_1 - f_2} \right) h = 30 + \left(\frac{30 - 27}{2 \cdot 30 - 27 - 5} \right) 10 = 31$	22
Mode class = 30-40.	23
$L = 30, f_m = 30, f_1 = 27, f_2 = 5, h = 10$	24
mode = $30 + \left(\frac{30 - 27}{2 \cdot 30 - 27 - 5} \right) 10 = 32.85$	25

$$\text{for variance, } \text{Var}(X) = \frac{\sum f_i (x_i - \bar{x})^2}{\sum f_i} = \frac{1423.53}{17} = \underline{\underline{83.73}}$$

the variance function is $\sigma^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$

$$\text{standard deviation, } s.d = \sqrt{\text{Var}} = \sqrt{83.73} = \underline{\underline{9.150}}$$

$$\text{coeff of variation, c.v} = \frac{\text{std deviation}}{\text{mean}} \times 100 = \frac{9.150}{29.70} \times 100 = \underline{\underline{30.80\%}}$$

lengths to calculate

$$\therefore \text{Mean} = 29.70 \quad \text{Variance} = 83.73$$

Median = 31 shown because $\sigma = 9.150$ is an even number

$$\text{Mode} = 32.85 \quad \text{C.V} = 30.80\%$$

calculate Pearson's Skewness coefficient for a dataset of exam

scores : 85, 88, 92, 94, 96, 98, 100, 100, 100, 100

Mean: mean sum of total marks no. of students

$$\text{Mean} = \frac{85 + 88 + 92 + 94 + 96 + 98 + 100 + 100 + 100}{10} = \underline{\underline{95.3}}$$

10 scores sorted in ascending order

so 5th and 6th positions will be median

median: write in ascending order first. 10 students

$$85, 88, 92, 94, 96, 98, 100, 100, 100, 100$$

Total 10 terms \rightarrow that are even = 8

so median is the average of 5th and 6th values when sorted in ascending order:

$$\text{Median} = \frac{96 + 98}{2} = \underline{\underline{97}}$$

$$\text{C.V} = \frac{100 - 95.3}{95.3} \times 100 = \underline{\underline{5.05\%}}$$

$$\text{Standard deviation} = \sigma = \sqrt{\frac{1}{N} \sum (x_i - \bar{x})^2} = \sqrt{\frac{1}{10} (85 - 95.3)^2 + (88 - 95.3)^2 + \dots + (100 - 95.3)^2}$$

$$= (85 - 95.3)^2 + (88 - 95.3)^2 + \dots + (100 - 95.3)^2$$

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{10} [(87 - 94)^2 + (91 - 94)^2 + \dots + (100 - 94)^2] = 26.81.$$

$$S.D = \sigma = \sqrt{26.81} = 5$$

Mode : 100 (100 is most frequently occurring data value in the data) (100 is bimodal)

S.K with respect to mean and median. Median is 94.5.
 $S.K = \frac{\text{mean} - \text{median}}{\text{standard deviation}} = \frac{94.5 - 94}{5} = 0.025$ (slight negative skewness in the distribution of exam scores)

S.K with respect to mean and mode $|S.K| = \frac{\text{mode} - \text{mean}}{\text{standard deviation}} = \frac{100 - 94.5}{5} = 1.1$

This means that the tail of the distribution is slightly longer on the left side, and most of the scores are concentrated on the right side of the mean. $94.5 + 5 + 100 + 104 + 108 + 112 + 116 + 120 + 124 + 128 = 1034$

Bowley's Measure : $B = \frac{Q_3 + Q_1 - 2Q_2}{Q_3 - Q_1}$ (Bowley's measure of skewness)

It is especially useful when dealing with data that may not follow a normal distribution or when a robust measure of skewness is required. Skewness : $B = \frac{Q_1 + Q_3 - 2Q_2}{Q_3 - Q_1}$ (Bowley's measure of skewness)

Quartiles : Q_1 is the first quartile (25th percentile) Definition
 Q_2 is the second quartile (50th percentile) Median
 Q_3 is the third quartile (75th percentile) Upper quartile

Median : $M = \frac{Q_1 + Q_3}{2}$ Median = (Q1 + Q3)/2

Interquartile Range : $IQR = Q_3 - Q_1$ Interquartile range = Q3 - Q1

Outliers : $Q_1 - 1.5 \times IQR$ and $Q_3 + 1.5 \times IQR$ Outliers = data points below Q1 - 1.5 * IQR or above Q3 + 1.5 * IQR

Box Plot : Q_1 , M , Q_3 , $Q_1 - 1.5 \times IQR$, $Q_3 + 1.5 \times IQR$ Box plot = box + whiskers

- Kurtosis** - Measures the degree of peakedness of the distribution.
- # Measure is denoted by $B_2 = \frac{\text{Var}(x^2)}{[\text{Var}(x)]^2}$
 - # Lepto-kurtic - A distribution with heavy tails and a sharp peak ($B_2 > 3$). Curve is peaked.
 - # Platy-kurtic - A distribution with light tails and a flatter peak ($B_2 < 3$). Curve is flat-topped.
 - # Mesokurtic - A normal distribution ($B_2 = 3$).
 - # Curve is Normal.
- Kurtosis** - measures degree of peakedness of the distribution.
-
- (SKEWNESS) Measures the degree of asymmetry of the distribution.
- # Moments: First moment about origin is zero for symmetric distributions.
 - # Moments are statistical measures that give certain characteristics of the distribution.
 - # formulae to calculate moments about the Mean:
 - # First Moment (about the Mean) of distribution is given by $\mu_1 = 0$ (since it is always zero for symmetric distributions).
 - # Second moment (about the Mean) is the Standard deviation σ .
 - # μ_2 (Variance) $\mu_2 = \sum f(x - \bar{x})^2$
 - # Third Moment (about the Mean) $= 4 \cdot 157$
 - # μ_3 (Skewness) $\mu_3 = \sum f(x - \bar{x})^3$

fourth Moment (about the mean)

$$M_4 (\text{kurtosis}), M_4 = \sum f(x - \bar{x})^4 \text{ (only 21 students)}$$

- ↳ four class marks: Mean = $\frac{A + 2B}{3}$ (middle point of the interval)
- ↳ $f(x) = f$ (constant) $\cdot (x - \bar{x})^4$

The kurtosis coefficient is defined as $\beta_1 = \frac{M_4}{M_2^2}$.

To calculate β_1 (Beta 1) and β_2 (Beta 2) for grouped data using a calculator, we need to first understand what these measures represent.

β_1 used as measure of skewness of distribution.

β_2 or measure kurtosis ("the "tailedness" of the distribution")

formula:

$$\beta_1 = \frac{M_3^2}{M_2^3}, \quad \beta_2 = \frac{M_4}{M_2^2}$$

where, M_2 is the second central moment (variance)

M_3 is the third central moment (used to measure skewness)

M_4 is the fourth central moment (used to measure kurtosis)

Same skills needed to calculate β_1 and β_2 as for calculating M_2 .

Steps to Calculate β_1 and β_2 (using) formula to calculate the mean, standard deviation $\sigma = \mu$

calculate the deviations ($x - \bar{x}$)

M_2, M_3 and M_4 = $\sum f(x - \bar{x})^2, \sum f(x - \bar{x})^3, \sum f(x - \bar{x})^4$

calculate β_1 and β_2 .

NOTE:

① If $\beta_2 = 3$, the data is normal or mesokurtic

② If $\beta_2 > 3$, the data is peaked or leptokurtic

③ If $\beta_2 < 3$, the data is flat-topped or platykurtic

Class Interval Frequency (f) X_i $f_i X_i$

0-10 4 10 40

10-20 6 15 90

20-30 8 20 160

30-40 5 15 75

40-50 3 13 39

45 13.5

Calculate the Mean :

$$\bar{X} = \frac{\sum f_i X_i + 20 + 90 + 100 + 125 + 135}{\sum f_i} = \frac{6210}{26}$$

$$\bar{X} = 23.8$$

$$N = 26$$

$$\sum (X_i - \bar{X})^2$$

$$0-10 4 5 10.8 -18.8 353.44$$

$$10-20 6 15 8.8 -8.8 77.44$$

$$20-30 8 25 -1.2 1.2 1.44$$

$$30-40 5 35 -11.2 11.2 125.44$$

$$40-50 3 45 -8.2 8.2 64.44$$

$$M_2 = \frac{\sum f_i (X_i - \bar{X})^2}{N}$$

$$M_2 = 1413.76 + 465.04 + 11.52 + 627.20 + 1348.32$$

$$M_2 = \frac{3865.84}{26} = 148.68$$

$$(X_i - \bar{X})^3 = (X_i - \bar{X})^4$$

$$-6644.672 \quad 124919.834$$

$$-681.472 \quad 5996.9536$$

$$-1.728 \quad 2.0736$$

$$1404.928 \quad 15735.1936$$

$$9528.128 \quad 201996.314$$

$$\mu_3 = -26578.68 + (-4088.82) + 18.76 + 7024.6 + 28584.36$$

$$= 190.5853$$

$$\mu_4 = 499679.32 + 35981.7 + 16.58 + 78675.95 + 605988.$$

$$= 48936.24$$

$$\mu_4 = 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$= 48936.24$$

$$=$$

sample:

- A subset of the population selected for analysis. It is used to draw conclusions about the population.
- A group of 100 randomly selected B.Tech students.

NOTE: size of samples, denoted by 'n' or sample size 'm'

size of population, denoted by 'N':

size of population = size of sample + size of remaining population

VARIABLE

Numerical (quantitative) variables: Categorical (qualitative)

discrete & continuous

Nominal < ordinal

le numerical value is can take any value in ordered, categories which are mutually exclusive.

of visits to hospital

height in cm (eg, male/female mutually exclusive).

discrete = eg, minimal/serve

Shape of Data

Skewness:

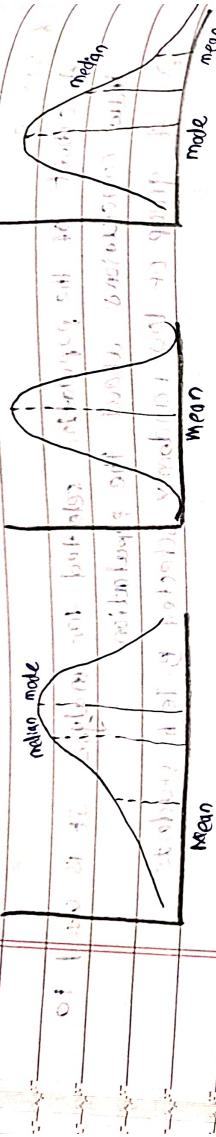
Measured asymmetry of the distribution of values in a data set. Skewness helps understand the underlying distribution of data, which is crucial for decision-making, risk assessment, and predicting future trends.

Relative to the mean: If the plot is skewed to the left, it is negatively skewed.

If the plot is skewed to the right, it is positively skewed.

If the plot is roughly symmetric, it has no skewness.

If the plot is roughly skewed to the right, it is right-skewed.



mean < median < mode \rightarrow both symmetrical data \rightarrow mode < median < mean

mean > median > mode \rightarrow mode < median < mean

mean = median = mode \rightarrow mode < median < mean

\rightarrow Skewness coefficient (Pearson's first coefficient of skewness):

This is a numerical measure of skewness, which determines the skewness when mean and mode are not equal.

If mean > mode, the skewness will consist positive value.

If mean < mode, the skewness will be a negative value.

If mean = mode, the skewness will be zero.

coefficient of skewness as Pearson's Measure of Skewness

Wrt Mean and Median:

$$\text{Coefficient of skewness} = \frac{3(\text{mean} - \text{median})}{\sigma}$$

σ

wrt mean and mode:

True Skew = $3(\text{mean} - \text{mode}) / (\text{standard deviation})^3$

if $S_k > 0$, it indicates a positively skewed distribution where the data is evenly balanced on both sides of the mean

\rightarrow If $S_k > 0$, it suggests a positively skewed distribution where the tail on the right side is longer or fatter, and most points are concentrated on the left side of the mean!

\rightarrow If $S_k < 0$, it indicates a negatively skewed distribution where the tail on the left side is longer or fatter, and most data points are concentrated on the right side of the mean.

Bowley's Measure:

Bowley's skewness coefficient, is a statistical measure used to assess the skewness or asymmetry in a probability distribution.

Coefficient of Bowley's Measure.

If $B=0$, the distribution is perfectly symmetric about the mean

If $B < 0$, the distribution is negatively skewed (left-skewed),

If $B > 0$, the distribution is positively skewed (right-skewed), meaning the tail on the left side of the distribution is longer or heavier.

If $B \neq 0$, the distribution is positively skewed (right-skewed)

indicating that the tail on the right side of the distribution is longer or heavier.

Quartiles, Deciles, Percentiles.

Partition values are statistical measures that divide a dataset into equal parts.

Quartiles - divide data into equal four parts.

(Ans: 1st Quartile = 25%) 2nd Quartile =

1st Q \rightarrow 2nd Q \rightarrow 3rd Q \rightarrow 4th Q

Start \rightarrow 2nd Q \rightarrow 3rd Q \rightarrow End

1st Q \rightarrow 2nd Q \rightarrow 3rd Q \rightarrow 4th Q

Q_2 = median.

$Q_1 = 1^{\text{st}}$ quartile (N : No. of observations)

Quartiles of Ans: Median \rightarrow 2nd Q \rightarrow 3rd Q \rightarrow 4th Q

$Q_1 = \text{first quartile}$ or $Q_1 = \frac{(N+1)^{\text{th}} \text{ term}}{4}$

$Q_2 = \text{second quartile}$ or $Q_2 = \frac{(2(N+1))^{\text{th}} \text{ term}}{4}$

= median.

$Q_3 = \frac{(3(N+1))^{\text{th}} \text{ term}}{4}$

$Q_3 = \text{Third quartile}$ or $Q_3 = \frac{(4(N+1))^{\text{th}} \text{ term}}{4}$.

$$\text{for deciles, } D_k = L + \left(\frac{\frac{kN}{10} - CF}{f} \right) X_h$$

f

for centiles, $P_k = L + \left(\frac{\frac{kN}{10} - CF}{f} \right) X_h$

Note: ~~Median formula is same as for deciles~~

Above formulas are for grouped as well as for ungrouped data.

Calculate the lower and upper quartiles of the following weights

In the family : 25, 19, 32, 11, 40, 36, 13, 51 and 46.

first, organize the given data in ascending order. So it is 11, 19, 25, 32, 36, 40, 46, 51 and 51

for

$$Q_1 = \text{the } (N+1)^{\text{th}} \text{ term} = \frac{1}{4} \text{ of } 10^{\text{th}} \text{ term} = 2.5^{\text{th}} \text{ term.}$$

Given data is 4, 19, 25, 32, 36, 40, 46, 51 and 51

$$= \frac{2^{\text{nd}} \text{ term} + 3^{\text{rd}} \text{ term}}{2} = \frac{11 + 13}{2} = 12 \text{ (second Q1)}$$

$$Q_2 = 12. \quad \text{Median is the } 5^{\text{th}} \text{ term.}$$

Activity: Invert the following steps again to calculate Q3

Note: Suppose $Q_1 = 2^{\text{nd}}$ term \Rightarrow 11, 19, 25, 32, 36, 40, 46, 51 and 51

Note: Suppose $Q_1 = 2.75$ which should be 11, 19, 25, 32, 36, 40, 46, 51 and 51

$$= 2^{\text{nd}} + 0.75 (3^{\text{rd}} \text{ term} - 2^{\text{nd}} \text{ term})$$

$$= 11 + 0.75 (13 - 11) = 12.5$$

$$= 12.5 + 1.5 = 14$$

$$= 14 + 1.5 = 15.5$$

$$= 15.5 + 1.5 = 17$$

$$= 17 + 1.5 = 18.5$$

$$= 18.5 + 1.5 = 19$$

$$= 19 + 1.5 = 20.5$$

$$= 20.5 + 1.5 = 22$$

$$= 22 + 1.5 = 23.5$$

$$= 23.5 + 1.5 = 25$$

$$= 25 + 1.5 = 26.5$$

$$= 26.5 + 1.5 = 28$$

$$= 28 + 1.5 = 29.5$$

In other words, you take $Q_1 = 13$.

as it is near to 3rd term, or it is greater than 2nd

and term of 3rd, and will be 13 within the word are

for $Q_3 = \text{median} = 25$

Activity: Invert the following steps again to calculate Q3

for Q_3 , $Q_3 = \left[\frac{3(N+1)}{4} \right]^{\text{th}} \text{ term} = \left[\frac{3(9+1)}{4} \right]^{\text{th}} \text{ term}$

$$= \left(\frac{3 \times 10}{4} \right)^{\text{th}} = (7.5)^{\text{th}} \text{ term}$$

If given data is grouped and in intervals then use formula $Q_k = L + \left(\frac{kN}{f} - CF \right) \times h$ for quartiles



$$Q_3 = (7.5)^{\text{th}} \text{ term}$$

$$= \frac{8^{\text{th}} \text{ term} + 9^{\text{th}} \text{ term}}{2} = \frac{35 + 40}{2} = 37.5$$

$$Q_3 = 37.5 \text{ and } h = 5$$

9th
quartile

(ii) calculate Q_1 and Q_3 for the data related to the age in years of 99 members in a housing society.

Age (in years) of 99 no. of members | cumulative freq.

Age	No. of members	Cumulative freq.
10	10	10
15	18	28
20	22	50
25	15	65
30	10	75
35	12	87
40	8	95
45	4	99

→ first find Cumulative frequency.

$$N (\text{all odd the frequency}) = 99.$$

$$Q_1 = L + \frac{(N+1)}{4} - CF = 10 + \frac{100}{4} - 28 = 18.$$

Now, the 25th term falls under CF of 25 and the age

against this CF value is 18.

If we have to find exact value of Q_1 for grouped data.

Then formula $N/2$ replaced by $N/4$.

$$\text{Median} = L + \frac{(N/2 - CF)}{F} \times h$$

$$Q_1 = L + \frac{(N/4 - CF)}{F} \times h$$

$$\text{for } Q_3 = \left[\frac{3(N+1)}{4} \right]^{\text{th}} \text{ term. } N = 9 \text{ g.}$$

$$Q_3 = \left(\frac{3 \times 100}{4} \right)^{\text{th}} \text{ term} = 75^{\text{th}} \text{ term} = 101 + 14.9 = 115.9$$

$\therefore 75^{\text{th}}$ term falls under CF of 85 and the age against this CF value is 40

Q. calculate Bowley's Measure of skewness for the following dataset representing the ages of a group of people in a sample.

samples: 20, 24, 28, 32, 35, 40, 42, 45, 50, 55

\Rightarrow arrange in ascending order:

20, 24, 28, 32, 35, 40, 42, 45, 50, 55 $N = 9$

Q_2 (median)

$\Rightarrow Q_2 = 35$. (the middle value)

for $Q_1 = \frac{(N+1)}{4}^{\text{th}}$ term = $\frac{10}{4}^{\text{th}}$ term = 2.5th term

for $Q_3 = \frac{3(N+1)}{4}^{\text{th}}$ term = $\frac{30}{4}^{\text{th}}$ term = 7.5th term

so, $Q_1 = 2^{\text{nd}}$ term + 3^{rd} term = $24 + 28 = \frac{52}{2} = 26$

so, $Q_3 = 7^{\text{th}}$ term + 8^{th} term = $42 + 45 = \frac{87}{2} = 43.5$

so, $Q_1 = 26$ & $Q_3 = 43.5$

for $Q_3 = \left(\frac{3(N+1)}{4} \right)^{\text{th}}$ term = $\left[\frac{3(10)}{4} \right]^{\text{th}}$ term = 7.5^{th} term

so, $Q_3 = 7^{\text{th}}$ term + 8^{th} term = $42 + 45 = \frac{87}{2} = 43.5$

$Q_3 = 43.5$

$B = Q_1 + Q_3 - 2Q_2 = 26 + 43.5 - 2(35) = 10.5$

$Q_3 - Q_1 = 43.5 - 26 = 17 = \text{midian}$

$B = Q_1 + Q_3 - 2Q_2 = 26 + 43.5 - 2(35) = 10.5$

$Q_3 - Q_1 = 43.5 - 26 = 17 = \text{midian}$

$$B = \frac{26 + 43.5 - 70}{43.5 - 26} = \frac{69.5 - 70}{17.5} = \frac{-0.5}{17.5} = -\frac{1}{35}$$

$B = -0.028$ Since $B < 0$ the distribution is negatively skewed (left-skewed). This means that the distribution is longer on the left side, indicating that there may be outliers or high values on the right side of the data.

Deciles : $D_1, D_2, D_3, D_4, D_5, D_6, D_7, D_8, D_9$

→ The deciles involve dividing a dataset into ten equal parts based on numerical values. There are therefore nine deciles altogether. Deciles are represented as follows:

$D_1, D_2, D_3, D_4, D_5, D_6, D_7, D_8, D_9$

A decile is used to group big data sets in descriptive statistics further from highest to lowest values or vice versa.

$D_1 = \frac{(N+1)}{10}^{\text{th}} \text{ term}$, $D_2 = \left[\frac{2(N+1)}{10} \right]^{\text{th}} \text{ term} \dots \text{ so on}$

$D_9 = \left[\frac{9(N+1)}{10} \right]^{\text{th}} \text{ term}$.

Q Calculate the D_1, D_5 from the following weights in a family: 25, 17, 32, 11, 40, 36, 13, 5 and 46

→ In ascending order: 5, 11, 13, 17, 25, 32, 36, 40, 46

Here, $N = 9$, $D_1 = \left(\frac{1}{10} \right)^{\text{th}} \text{ term} = 1$, $D_5 = \left(\frac{5}{10} \right)^{\text{th}} \text{ term} = 5$.

$D_1 = \left(\frac{1}{10} \right)^{\text{th}} \text{ term} = \frac{5 \times 10}{10} = 5$, $D_5 = \left(\frac{5}{10} \right)^{\text{th}} \text{ term} = \frac{5 \times 10}{10} = 5$

Percentiles. $\alpha = \text{exp}(\beta)$ or $\alpha = \text{exp}(\beta_0 + \beta_1 x)$

centiles is another form of percentiles.

Any given observation is essentially divided into a total of 100 equal parts by a percentile or percentile represents as $P_1, P_2, P_3, P_4, \dots, P_{99}$, which is a typical value of peaky for when data is either less than, plus or equal to P_i .

$$\rho_1 = \left[\frac{N+1}{100} \right]^{\text{th term}}, \quad \rho_2 = \left[\frac{g(N+1)}{100} \right]^{\text{th term}} \dots$$

$$P_{\text{avg}} = \frac{1}{N} \sum_{n=1}^N P_n = \frac{1}{N} \sum_{n=1}^N \left[\frac{1}{2} \left(\frac{1}{2} + \frac{1}{2} \cos \left(\frac{2\pi n}{N} \right) \right) \right] = \frac{1}{2} + \frac{1}{2N} \sum_{n=1}^N \cos \left(\frac{2\pi n}{N} \right)$$

PRESBYTERIAN CHURCH OF THE SAVANNAH,
Savannah, Georgia, January 1, 1900.

N = total no. of observations

P_1 is first percentile. P_2 is second percentile and so on.

adults in P₁₀ and P₇₅ for the date reflected to the age

Age (in years) No of members in a housing society.

10 01 20

18 19 20

25 10 10 25

85 20 45

first find cumulative freq

$$P_{1,0} = \left(10(N+1)\right)^{-1} \text{frequency. } 15N = 99 \text{ vibrations. } 0.6$$

1900-1901
1901-1902
1902-1903
1903-1904
1904-1905
1905-1906
1906-1907
1907-1908
1908-1909
1909-1910
1910-1911
1911-1912
1912-1913
1913-1914
1914-1915
1915-1916
1916-1917
1917-1918
1918-1919
1919-1920
1920-1921
1921-1922
1922-1923
1923-1924
1924-1925
1925-1926
1926-1927
1927-1928
1928-1929
1929-1930
1930-1931
1931-1932
1932-1933
1933-1934
1934-1935
1935-1936
1936-1937
1937-1938
1938-1939
1939-1940
1940-1941
1941-1942
1942-1943
1943-1944
1944-1945
1945-1946
1946-1947
1947-1948
1948-1949
1949-1950
1950-1951
1951-1952
1952-1953
1953-1954
1954-1955
1955-1956
1956-1957
1957-1958
1958-1959
1959-1960
1960-1961
1961-1962
1962-1963
1963-1964
1964-1965
1965-1966
1966-1967
1967-1968
1968-1969
1969-1970
1970-1971
1971-1972
1972-1973
1973-1974
1974-1975
1975-1976
1976-1977
1977-1978
1978-1979
1979-1980
1980-1981
1981-1982
1982-1983
1983-1984
1984-1985
1985-1986
1986-1987
1987-1988
1988-1989
1989-1990
1990-1991
1991-1992
1992-1993
1993-1994
1994-1995
1995-1996
1996-1997
1997-1998
1998-1999
1999-2000
2000-2001
2001-2002
2002-2003
2003-2004
2004-2005
2005-2006
2006-2007
2007-2008
2008-2009
2009-2010
2010-2011
2011-2012
2012-2013
2013-2014
2014-2015
2015-2016
2016-2017
2017-2018
2018-2019
2019-2020
2020-2021
2021-2022
2022-2023
2023-2024
2024-2025
2025-2026
2026-2027
2027-2028
2028-2029
2029-2030
2030-2031
2031-2032
2032-2033
2033-2034
2034-2035
2035-2036
2036-2037
2037-2038
2038-2039
2039-2040
2040-2041
2041-2042
2042-2043
2043-2044
2044-2045
2045-2046
2046-2047
2047-2048
2048-2049
2049-2050
2050-2051
2051-2052
2052-2053
2053-2054
2054-2055
2055-2056
2056-2057
2057-2058
2058-2059
2059-2060
2060-2061
2061-2062
2062-2063
2063-2064
2064-2065
2065-2066
2066-2067
2067-2068
2068-2069
2069-2070
2070-2071
2071-2072
2072-2073
2073-2074
2074-2075
2075-2076
2076-2077
2077-2078
2078-2079
2079-2080
2080-2081
2081-2082
2082-2083
2083-2084
2084-2085
2085-2086
2086-2087
2087-2088
2088-2089
2089-2090
2090-2091
2091-2092
2092-2093
2093-2094
2094-2095
2095-2096
2096-2097
2097-2098
2098-2099
2099-20100

$$= \frac{10}{100} \cdot \frac{99+1}{100} = 10^{\text{th}} \text{ term}$$

Now, the 10th term falls under the cumulative frequency of 20 and the age against this CF value is 10.

$P_{10} = 75^{\text{th}}$ term.
Now, the 75th term falls under the CF of 85 and the age against CF value is 40.

Now, the 75th term falls under the cumulative frequency of 85 years. So, the median age is 40 years.

Data Visualization:

1) Histogram.

A histogram is a graphical representation of the distribution of numerical data. It is similar to a bar chart but specifically used for quantitative data that is divided into ranges (also called bins or intervals). Histograms are essential for visualizing the frequency of data points in each range, helping to understand the shape and spread of the data distribution.

2) Box Plot (Box-and-whisker plot).

* outlier

* whisker

* max value of dataset

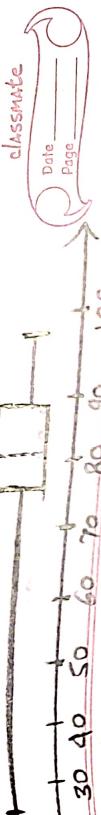


min. value of dataset. Maximum is 91.

max. value of dataset. Minimum is 91.

outlier

Box-and-whisker plot



Test scores for college statistics class held during the evening are:

98, 78, 68, 83, 81, 89, 88, 76, 65, 45, 98, 90, 80, 84.5,

85, 79, 78, 78, 90, 79, 81, 25.5 21 minutes later,

and the smallest and largest values, the median, and the first

and third quartile for these night classes. In what order?

Arrange the data in ascending order.

25.5, 45, 65, 76, 78, 78, 79, 79, 80, 81, 81, 83, 84.5,

85, 88, 89, 90, 90, 96, 96, 98, 98, 98, 98, 98.

$$\text{min value} = 25.5 \quad N=22$$

$$\text{max value} = 98 \quad \text{maximum is } 21 \text{ and } N = 22 \text{ so } Q_1 = 21$$

$$\text{median} = \frac{78+81}{2} = 79.5 \quad \text{since } 21 \text{ is the middle value}$$

$$\text{and } 22 \text{ is even} \quad \text{and } 2 \text{ students have } 78 \text{ and } 79$$

$$\text{so } Q_3 \text{ is the } 16^{\text{th}} \text{ value which is } 90 \text{ and } Q_1 \text{ is } 78$$

$$Q_1 = (N+1)^{\text{th}} \text{ term} = \frac{(22+1)}{4}^{\text{th}} \text{ term} = 23^{\text{th}} \text{ term}$$

$$= 5.75^{\text{th}} \text{ term} = 5^{\text{th}} \text{ term} + 0.75(6^{\text{th}} - 5^{\text{th}}) \text{ term}$$

$$= 76 + 0.75(2) = 76 + 1.50 = 77.5$$

$$Q_3 = (N+1)^{\text{th}} \text{ term} = \frac{(22+1)}{4}^{\text{th}} \text{ term} = 23^{\text{th}} \text{ term}$$

$$= 5.75^{\text{th}} \text{ term} = 5^{\text{th}} \text{ term} + 0.75(6^{\text{th}} - 5^{\text{th}}) \text{ term}$$

$$= 76 + 0.75(2) = 76 + 1.50 = 77.5$$

$$Q_1 \text{ is the median of first half data. } Q_3 \text{ is the median of second half data.}$$

$$\text{for } Q_1, 25.5, 45, 65, 68, 76, 78, 79, 80, 81, Q_1 = 78$$

$$\text{for } Q_2, 81, 83, 84.5, 85, 88, 89, 90, 90, 93, 98, Q_2 = 89.$$



classmate

Date _____

Page _____

#

Data Visualization

1)

Histogram:

A histogram is a graphical representation of the distribution of numerical data. It is similar to a bar chart but specifically used for quantitative data that is divided into ranges (also called bins or intervals). Histograms are useful for visualizing the frequency of data points in each range, helping to understand the shape and spread of the data distribution.

2)

Bins (intervals) - A bin is a continuous range of values

within which data points are grouped together. Each bin represents a specific interval of values, and the height of the bar for each bin shows the no. of data points (or frequency) that fall within that range. $f(14-20) = \text{no. of } (14-20) = 0$

Ex. If you're analyzing the test scores of students (0-100), bins could be set at intervals of 10 (0-10, 11-20, 21-30 etc)

3)

frequency distribution: choosing the no. of bins approx. equal to the square root of the no. of observ. often works well in practice

4)

equal units along the horizontal axis (x-axis, abscissa) reflect the various class intervals of the frequency distribution. Equal units along the vertical axis (the y-axis, ordinate) reflect increased frequency. (The units along the vertical axis do not have to be the same width as those along the horizontal axis.)

5)

use wiggly lines to highlight breaks in scale

6)

Choosing bins: $\frac{\text{range}}{\text{width}} + 1$. Need about 5-10 bins for 20-50

7)

Too few bins - if there few bins, bin width is too large, the histogram may not show important details about the distribution.

8)

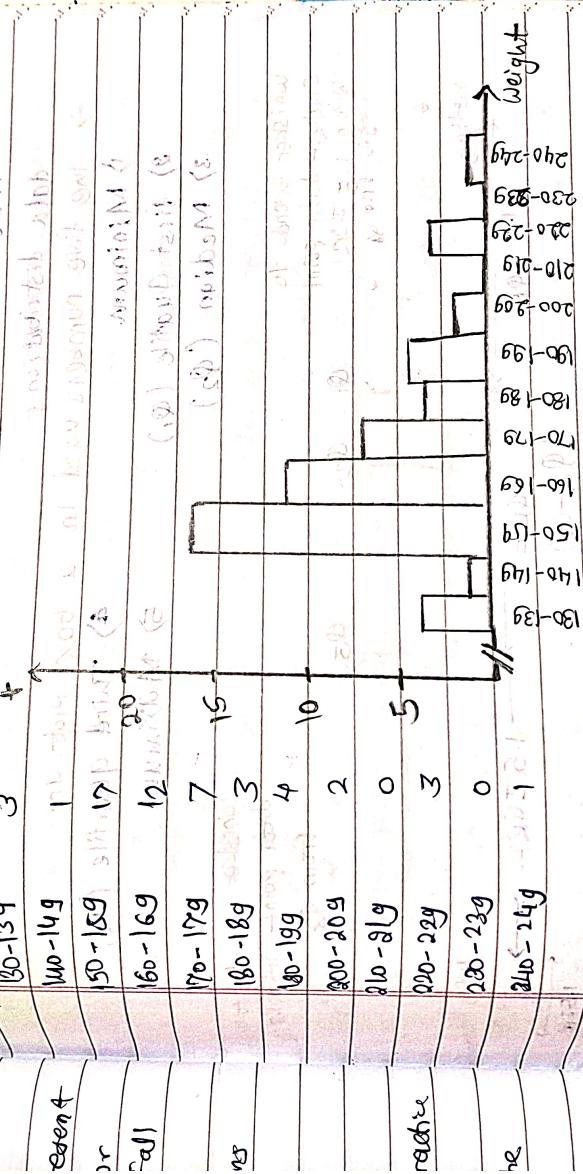
Too many bins - if there are too many bins, the distribution

- Too many Bins = If the bin width is too small, the histogram can be ~~overly~~ detailed and noisy.
- Label the bin (class interval) boundaries on a horizontal scale.
- Mark and label the vertical scale with the frequencies or the relative frequencies.

Q. Draw the histogram information following details.

1) weight of 30 fish kinds
2) Number of fish
3) Weight

100-129	3	F
---------	---	---

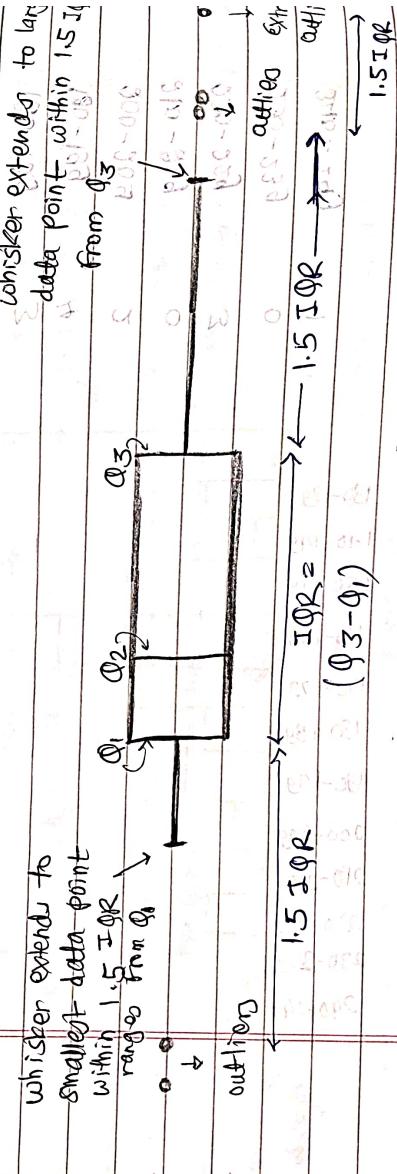


Topic 2) Box Plot (Box-and-whisker plot)

The Box Plot is a graphical representation of a dataset's five-number summary: minimum, first quartile (25th %), median (50th %), third quartile (75th %), and maximum. It is a powerful tool in data analysis because it clearly highlights the dataset's central tendency, dispersion and skewness. It effectively visualizes outliers. This is particularly useful when comparing multiple datasets, as it offers a clear, comparative visualization of the different data distributions.

→ The five numbers used in a box plot are:

- 1) Minimum
- 2) First Quartile (Q_1)
- 3) Median (Q_2)
- 4) Third Quartile (Q_3)
- 5) Maximum



- The Q_2 median is the middle value that separates the data into two halves. It measures central tendency, providing a snapshot of the data's center.
- Q_1 and Q_3 , marking the box ends, reflect the dataset's dispersion. Q_1 represents 25th % of the dataset and Q_3 represents 75th % of the dataset.
- The whiskers are lines extending from the box, reaching the min. and max. non-outlier data points.



- $Q_1 \rightarrow$ median of first half of the data.
- $Q_3 \rightarrow$ median of second half.

Usually, the lower whisker extends from Q_1 to the smallest non-outlier data point and the upper whisker extends from Q_3 to the largest non-outlier data point.

The length of the box is the IQR.

$$\boxed{\text{IQR} = Q_3 - Q_1}$$

- The IQR measures the middle 50% of the data, measuring dispersion or spread.
- Outliers are typically calculated as data points that fall below $(Q_1 - 1.5 \text{ IQR})$ or above $(Q_3 + 1.5 \text{ IQR})$.

These outliers are represented as individually pointed outside the whiskers in the box plot.

- A point more than 3 IQR from the box edge is called an extreme outlier.

~~if box contains 100% of data, then it's not an outlier~~

How to Interpret a Box Plot:

- 1) Length of the Box: The length of the box (between Q_1 and Q_3) represents the IQR, showing the spread of the middle 50% of the data.

~~if IQR = 0, then no spread~~

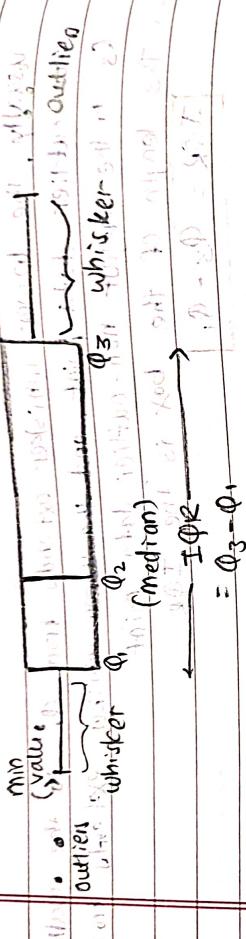
~~if IQR > 0, then spread~~

~~if IQR < 0, then no spread~~

- 2) Median line: If the median line is closer to Q_1 , the data is left skewed (longer whisker on the right).

- 3) Whiskers: The whiskers represent the range of the data, excluding outliers.

* Box plot (5-Number summary) \rightarrow max value



Q. 22,25, 23, 17, 19, 33, 64, 117, 2018, ~~Half Box~~ lot. 2

as cending order

19	18	19, 20	22, 23, 25	33, 64, 83
First part	Second part	Third part	Fourth part	Fifth part

$$\varphi_2 = \frac{20 + 22}{2} \quad 21 \rightarrow \varphi_2$$

$\beta \rightarrow 0$, and $\beta \rightarrow \infty$ are similar.

If two terms are present at center, then it would take

Average of these two terms for finding ϕ and β_3

$$JQR = \varnothing_3 - \varnothing_1 = 125 - 118 = 7$$

New check for outliers

- ① Higher outlier
- ② Lower outlier

$$= Q_3 + (1.5 \times I(Q)) = Q_1 - (1.5 \times I(Q))$$

$$= 1.25 + (1.5 \times 7) = 1.25 + 10.5 = 18 - 10.5 = 7.5$$

$$= 35.5 \text{ (high) on 99.5\% = 9.5}$$

Any data > 3S, S, or 1 S will tell if any data < 7.5% of

is lower outlier.

$$\phi_1 = 18^\circ, \phi_2 = 91^\circ, \phi_3 = 25^\circ$$

min = 17 max = 33 address = 64 offset = 111

not 64 bcoz, 64 is outlier