

Intermediate Stats.

Central

- (1) Measure of Central Tendency
- (2) measure of Dispersion
- (3) Gaussian Distribution
- (4) Percentiles and Quartiles
- (5) 5 Number summary (Box Plots).

* Measure of Central Tendency

1) mean

2) median

3) mode

(m) mean

(x̄) x̄ Ba

A measure of Central Tendency is a single value that attempts to describe a set of data by identifying the central position.

① mean (Average)

Population (N)

Population mean (μ)

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

Distribution

$$x = \{1, 2, 3, 4, 5\}$$

Sample (n)

Sample mean (\bar{x}_c) =

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

$$= \frac{1 + 2 + 3 + 4 + 5}{5}$$

$$= \frac{15}{5} = 3$$

Average = 3

Distribution of data = 3

* Population vs Sample Example.

Page No _____
Date _____

* Population data { 24, 23, 2, 1, 28, 27 }
 * median (N)

* Population mean

$$\bar{M} = \frac{24 + 23 + 2 + 1 + 28 + 27}{6}$$

Central distribution = 17.5

* Sample data (n) = { 23, 2, 28, 27 }

* Sample mean. $\bar{x} = 23 + 2 + 28 + 27$ (Random Pick value)

Central distribution = 20.5

* Examples of An mean , median and mode
median

DATASET

Age Salary Family
 NAN = (Empty Value).

—	—	—
—	NAN	—
NAN	—	NAN
—	—	NAN
—	NAN	—
—	—	—
NAN	—	NAN

As, here we have an dataset of an age, salary, and family size, where as there are some empty value (NAN). So here we have to fulfill the data set.

As, If we delete the empty column it may cause a loss of information.

So, we will fill it by Avg Column

$$\text{eg: } \text{Age} = \{24, 26, \text{NAN}, 21, 20, 18\}$$

$$\text{Average}(M) = \frac{24 + 26 + 21 + 20 + 18}{5}$$

$$= \underline{\underline{21.8}}$$

As, here we can replace the NAN value with value of 21.8, So the dataset will be $\{24, 26, 21.8, 21, 20, 18\}$

median :- The middle number in sorted list of Number more descriptive of data set than avg

As, here if we have an dataset $\{1, 2, 3, 4, 5\}$

$$(m) \text{ Average} = \frac{1+2+3+4+5}{5} = \underline{\underline{3}}$$

with same data set $\{1, 2, 3, 4, 5, 100\}$. we add 100

$$(m) = \frac{1+2+3+4+5+100}{6} = \underline{\underline{19}}$$

As Now. $\{1, 2, 3, 4, 5\} \leftrightarrow \{1, 2, 3, 4, 5, 100\}$
we can see an difference between dataset.

$$\text{Range} = (5-1) = \underline{\underline{4}}$$

$$\text{Range} (1-100) = \underline{\underline{99}}$$

outlier

$\{1, 2, 3, 4, 5\}$

start order.

(By Calculating Range and been huge range)

difference over both dataset identify that there is an outlier

* Definition of Outlier.

Outlier is a number that is complete different than entire distribution.

If the data set has an outlier we can replace the value using median.

* Steps to find median. {1, 2, 3, 4, 5, 100}

(i) Sort the Numbers.

(ii) Find the Central Number.

(i) If the no. of element are even we find average.

(ii) If the no. of element are odd we find central number.

Even no.

$$\{1, 2, \boxed{3, 4}, 5, 100\}$$

$$\frac{3+4}{2} = \underline{\underline{3.5}}$$

median 3.5

Odd no.

$$\{1, 2, \boxed{4}, 5, 100\}$$

median 4

MODE

: The Most frequent Occuring element

Dataset.

{ Categorical Variab

Types of flower:

- (1) Lily
- (2) Sunflower
- (3) Rose
- (4) Lily
- (5) NAN
- (6) Rose
- (7) Rose
- (8) Sunflower.

← Replace with something.

As, here most frequent Occuring elem
is an Rose. (so, NAN Replace with pr)

Uses :-

- mean(Avg)
- median → (Outliers)
- mode → (Categorical replacement)

* Assignment

Find mean, median and mode of following distribution.

$$X = \{24, 25, 26, 27, 28, 90, 100, 1000, \\ 1200, 1400, 1400, 1400\}$$

(i) Mean (Average)

$$(M) = \frac{24 + 25 + 26 + 27 + 28 + 90 + 100 + 1000}{12} \\ + 1200 + 1400 + 1400 + 1400 \\ = \frac{6720}{12} = 560.$$

mean = 560.

(ii) Median

$$\{24, 25, 26, 27, 28, 90, 100, 1000, 1200, 1400, \\ 1400, 1400\}$$

$$\frac{90 + 100}{2} = \frac{190}{2} = 95$$

median = 95

(iii) Mode = 1400

(2)

Measure of Dispersion

(1)

Variance (σ^2) , Sigma²

(2)

Standard Deviation (σ) , Sigma

Dispersion (Distribution of Data).

1.

Variance

Variance is used to understand proper distribution between data.

Example :- If we have Two data As

$$X = \{1, 1, 2, 2, 4\}$$

$$\bar{m} = \underline{2}$$

$$Y = \{0, 2, 2, 3, 3\}$$

$$\bar{m} = \underline{2}$$

As, here we can see that the mean of both data set is equal, do it not? Consider that both data set are same.

* Population Variance (σ^2) , Sigma²

$$\sigma^2 = \frac{N}{\sum_{i=1}^N (x_i - \bar{m})^2} \rightarrow \text{Population mean}$$

(Spread of Distribution)

As example.

$$A = \{1, 2, 2, 3, 4, 5, 10\} \quad (\text{Range} : 10 - 1 = \underline{\underline{9}})$$

\Downarrow Range difference.

$$B = \{2, 2, 4, 6\} \quad (\text{Range} : 6 - 2 = \underline{\underline{4}})$$

As, here with A Dataset we can see that the spread of data distribution is high; with B Dataset distribution is low.

If Variance is Become high the spread of distribution is high.

If Variance is Become low Spread of distribution is lesser.

Bessel's Correction

* Sample Variance.

$$S^2 = \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{n-1} \right)^2$$

Sample mean.

Degree of freedom

As, here we can see that there is some difference in formulas ($n-1$)

Population Standard deviation. [Variance]

formula

$$\sigma = \sqrt{\sigma^2}$$

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

Sample Standard deviation [Variance]

formula

$$s = \sqrt{s^2}$$

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

Dataset : $x = \{23, 21, 20, 19, 24, 27, 28\}$

x	\bar{x}	$(x - \bar{x})$	$(x - \bar{x})^2$
23	23.14	- 0.14	0.196
21	23.14	2.14	4.5796
20	23.14	3.14	9.8596
19	23.14	4.14	17.1396
24	23.14	- 0.86	0.7396
27	23.14	3.86	14.8996
28	23.14	4.86	23.6196
			71.0336

Formula

$$\sigma^2 = \frac{71.0336}{6} = 11.8$$

$$\sigma = \sqrt{11.8} = 3.44$$

Variance = 11.8

Standard deviation = 3.44.

(3) Percentiles and Quartiles.

To understand the percentiles lets get understand percentage.
Basics of example.

Calculating the Percentage of even and odd numbers.

$$\text{Percentage} = \{1, 2, 3, 4, 5, 6, 7, 8\}$$

\therefore Percentage of even no. = $\frac{\text{No. of even number}}{\text{Total no. of numbers}}$



$$50\% = \frac{4}{8} = \frac{1}{2} = 0.5$$

Percentage of odd no. = $\frac{\text{no. of odd number}}{8}$

$$50\%$$

$$= \frac{4}{8} = \frac{1}{2} = 0.5$$

Percentiles.

As, In many of exam we see their is an percentiles such as, SAT, CAT or GRE the result where display as an percentile 99%, 85% or 84%

Definition :- A percentile is a value below which a certain percentage of observations lie.

If a person got an 99 percentile in exam. It means that the person has got better marks than 99% of the entire students.

Examples :-

Dataset : 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 9, 9, [10], 11, [12], 12.

Q. What is the percentile of 10 ?

Formula : Percentile Rank of $x = \frac{\text{number of values} \leq x}{n} \times 100$

$$= \frac{16}{20} \times 100 = \underline{\underline{80}}$$

Percentile = 80

As, here we can say that 10 is greater than 80 percentage of entire distribution.

Q. what is the percentile of 11 ?

$$\text{Percentile Rank of } x = \frac{\text{no. of value below } x}{n} \times 100$$

$$= \frac{17}{20} \times 100$$

$$= \underline{85} \quad \text{Percentile}$$

A.S, If we want to calculate or find
the percentile ..

Q. find the 25% percentile ?

$$\text{formula : Value} = \frac{\text{Percentile}}{100} \times n + 1$$

$$= \frac{25}{100} \times 20 + 1$$

$$= 5.25 \quad (\text{Index})$$

(i)

Minimum

First Quartile (25%) (Q₁){ Remaining
Outliers . }

Median

Third Quartile (75%) (Q₃)

Maximum.

5 number summary (Box Plot)

As we understand the concept with an example.

$$\left\{ \begin{array}{l} \boxed{1}, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 7, 8, 8 \\ \boxed{27} \end{array} \right\}$$

$\xrightarrow{\text{L}_{15}}$

$\xrightarrow{\text{Outlier}}$

As here we can see that 27 is an outlier but also consider that it is also an outlier. As in this condition we create a fence to understand / determine the outlier.

Standard formula based on data distribution

$$\begin{aligned} \text{Lower fence} &= Q_1 - 1.5 (\text{IQR}) \\ \text{Higher fence} &= Q_3 + 1.5 (\text{IQR}) \end{aligned}$$

$$\text{IQR} = (\text{Upper Quartile range})$$

As first find Q₁ and Q₃ of data set.

Lower fence (25 Percentile)

$$Q_1 = \frac{25}{100} \times (n+1)$$

$$\therefore \frac{25}{100} \times (20) = \underline{\underline{5}}$$

So the 5th Order of dataset

Higher fence (95 Percentile)

$$= \frac{95}{100} \times (n+1)$$

$$\therefore \frac{95}{100} \times (20) = \underline{\underline{19}}$$

So the 15th Order of dataset is $\boxed{19}$

$$IQR = Q_3 - Q_1$$

$$= 4 - 3 \\ \underline{\underline{1}} \\ (IQR)$$

Now we will calculate an Lower and higher fence.

Lower fence = $Q_1 - 1.5 \times IQR$

$$= 3 - 1.5 \times 4$$

$$= 3 - 6$$

$$= -3$$

$$\text{Upper fence} = Q_3 + 1.5 \times IQR$$

$$= 7 + 1.5 \times 4$$

$$= 13$$

As after getting all the values now just put the value using 5-number summary

$$\text{minimum} = 1$$

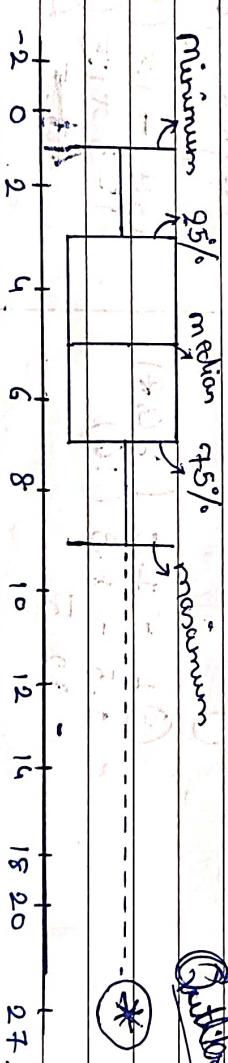
$$\text{First Quartile (Q1)} = 3$$

$$\text{median} = 5$$

$$\text{Third Quartile (Q3)} = 7$$

$$\text{maximum} = 9$$

Box Plot Identifying Outlier



Assignment

Page No.
Date

* Find minimum, maximum, Q_1 , Q_2 , Q_3 median of following data set

$\{ -8, 1, 2, 4, 5, 6, 8, 15, 20, 120 \}$

Outlier

$$Q_1 = \frac{25}{100} * 11 = 2.75 \quad (\text{Index})$$

$$Q_3 = \frac{75}{100} * 11 = 8.25 \quad (\text{Index})$$

1 + 2

$$Q_1 = 1.5 \\ Q_3 = 14.5 \\ \frac{Q_3 - Q_1}{2} = \frac{14.5 - 1.5}{2} = 6.5$$

$$\frac{35}{2} = 17.5$$

$$IQR = Q_3 - Q_1 \\ = 14.5 - 1.5 \\ = 13$$

Lower fence.

Higher fence

$$Q_1 = 1.5 \quad (IQR) \\ 1.5 - 1.5 \quad (16) \\ - 22.5 \\ \boxed{41.5}$$

$$\{ -8, 1, 2, 4, 5, 6, 8, 15, 20, 120 \}$$

minimum = -8

Q_1 = 2.5

median = 5.5

Q_3 = 17.5

maximum = 20

