

Day-2 - Statistics

04/09/2022

Histogram:

CONTINUOUS DATA

$$\text{Age} = \{10, 12, 14, 18, 24, 26, 30, 35, 36, 38, 40, 41, 42, 43, 50, 51, 68, 78, 90, 95, 100\}$$

① Sort the Numbers

② Bins \rightarrow No. of groups

③ Bins size \rightarrow size of bins

$$\min = 10$$

$$\max = 100$$

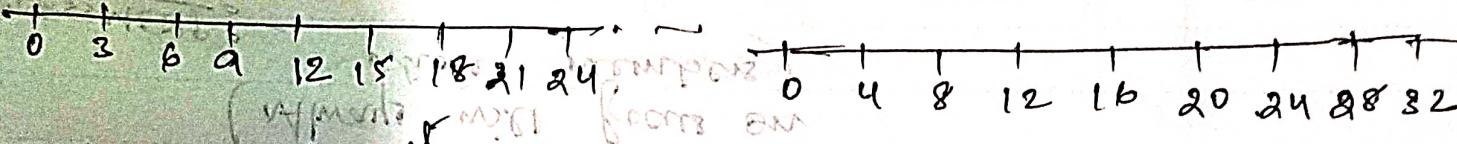
$$[10, 20, 25, 30, 35, 40]$$

Not considering

$$\text{Bin Size} = \frac{\text{Max} - \text{Min}}{\text{Bins}}$$

$$\frac{100 - 10}{10} = 9$$

$$\frac{100 - 10}{10} = 9$$



freq | count

bin 1

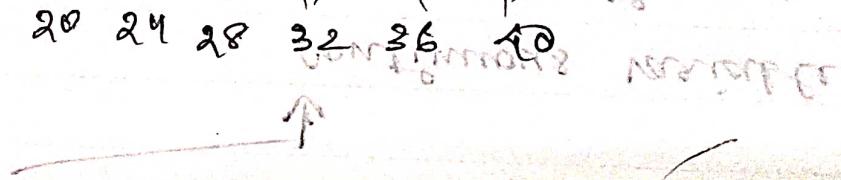
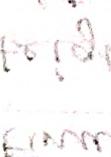
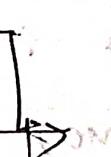
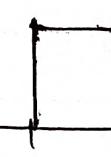
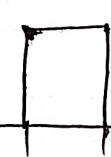
bin 2

group = bins

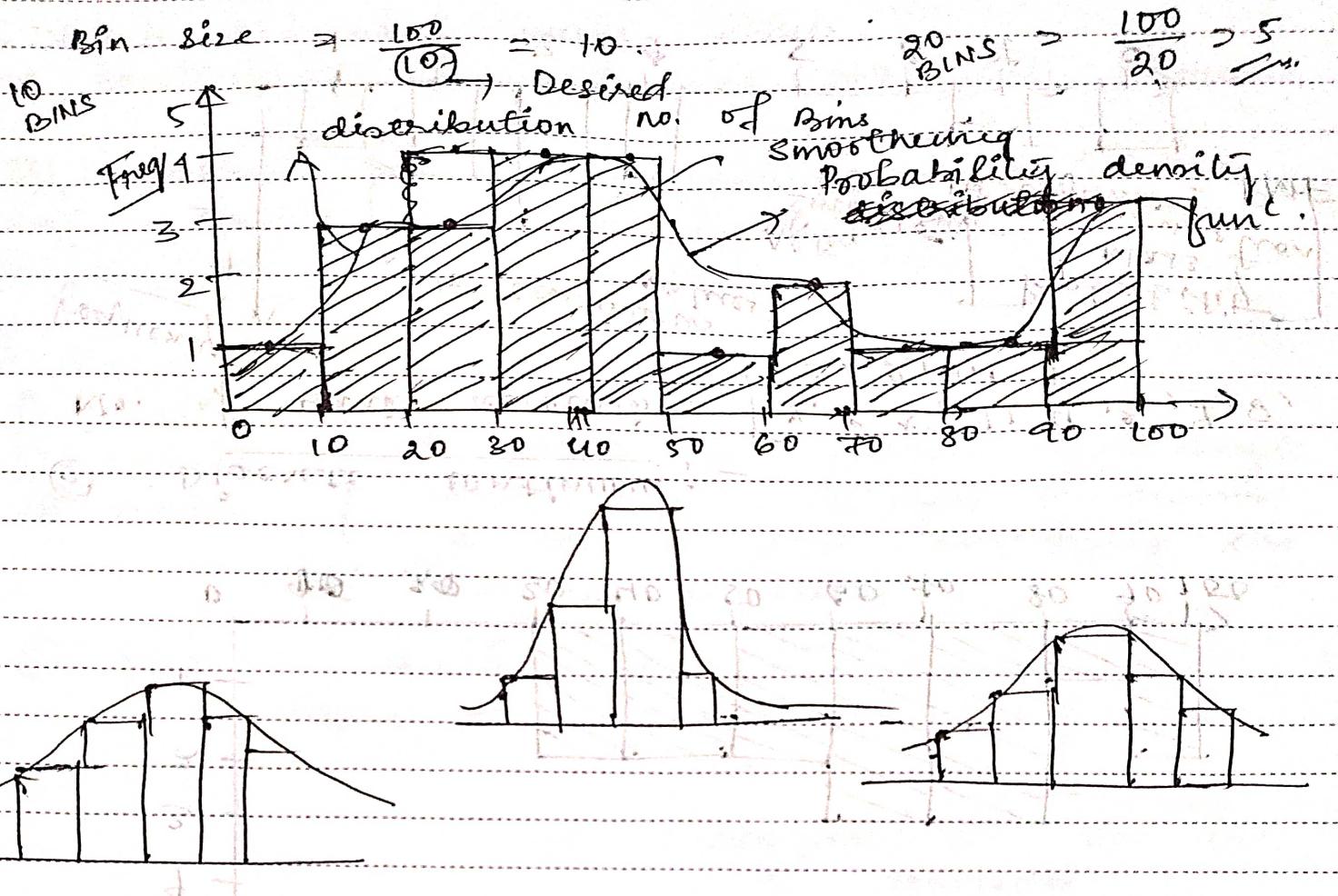
group size = bin size

Bins = 10 No. of Bins = 10

intervals



↑



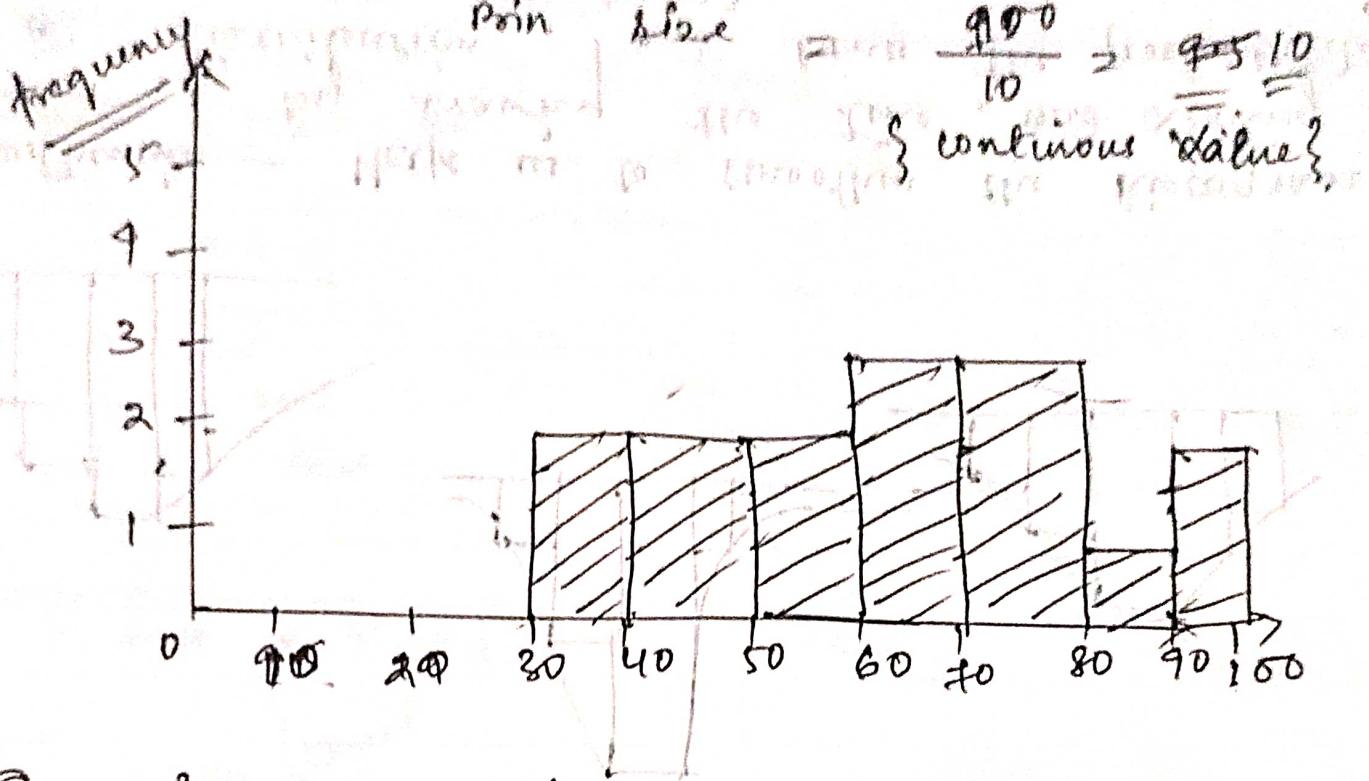
Smoothing — Helps us to smoothen the histogram by drawing the line ~~conforming~~^{connecting} ~~approx~~^{approx} points of distribution, and hence the probability density function (P.d.f).

Weight = { 30, 25, 25, 38, 42, 46, 48, 59, 62, 63, 68, 75, 77, 80, 90, 95 }₁₀₀

Ans. \Rightarrow 10 continuous values \Rightarrow $\frac{95-30}{10} = \frac{65}{10} = 6.5$
 → Not starting from 0 }

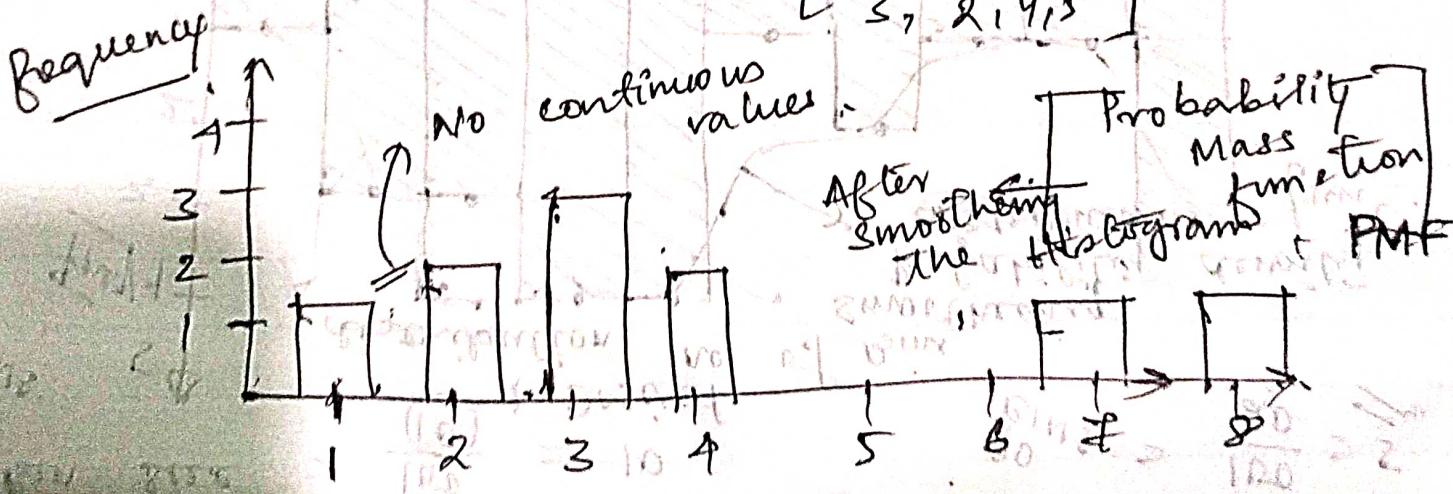
If starting from 0, then,

Point value $\Rightarrow \frac{900}{10} = 90$ \Rightarrow continuous values?



② Discrete vs continuous :-

No. of Bank accounts = [2, 3, 5, 11, 4, 2, 7, 8, 3, 2, 4, 5]



- Discrete has whole number values.
- continuous have decimal and whole numbers.

Pdf : Probability density function → Continuous

Pmf : probability Mass Function → discrete

* Measures of Central Tendency:

{ ① Mean
② Median
③ Mode } (A measure of tendency is a single value that attempts to describe a set of data identifying the central position)

$$\textcircled{1} \text{ Mean} = \bar{x} = \{1, 2, 3, 4, 5\}$$

$$\text{Avg (Mean)} = \frac{1+2+3+4+5}{5} = \frac{15}{5} = 3.$$

Central element in the set of elements

Population (N) : ~~mean~~ Sample (n) mean

Sample mean (\bar{x})

Population mean (μ)

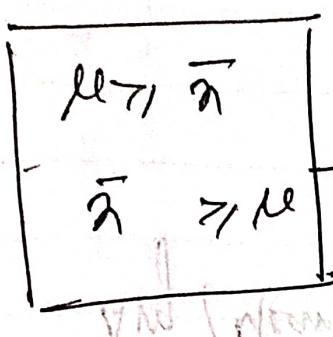
$$\mu_2 = \frac{\sum p_i n_i}{N} = \frac{\sum p_i \bar{x}_i}{N} = \frac{\sum p_i \frac{n_i}{n}}{N} = \frac{\sum p_i n_i}{N} = (\bar{x}) = \frac{\sum n_i}{N}$$

Ex. Population is $\{24, 28, 21, 28, 27\}$

$$\text{Population mean } (\bar{x}) = \frac{(24+23+21+28, 27)}{5} = 27.5$$

$$\text{Sample mean: } (\bar{x}) = \frac{\sum_{i=1}^n x_i}{n} = \frac{24+27+1+2}{4} = 18.5$$

$N > n$



For more all cases, \rightarrow this cannot happen

This can happen

① Practical Application ; Feature Engineering

Age	Salary	Family size
NAN	NAN	NAN

So, we can take the mean of

the columns and substitute the value with

1) if the no. of elements are even,
we find out the avg. of the central
elements?

2) if the no. of elements are odd, we
find the central elements.

Ex. { 0, 1, 2, 3, 4, [5, 6], 7, 8, {100, 120} }
Median = $\frac{5+6}{2} = \underline{\underline{5.5}}$,
Mean = 25.6
No outliers \Rightarrow Mean
with outliers \rightarrow Median

② Mode: { most frequent occurring elements }

Ex. { 1, 2, 1, 2, 3, 3, 4, 3 } \rightarrow Mode = 3

Application

Types of Flowers

dily

Sunflower

Rose

Rose

Sunflower

NAN

Rose

NAN

Rose

NAN

{ Replace the NAN
value with
max. no. of
occurrence
here, Rose }

III : Measures of Dispersion :-

- 1) Variance (σ^2) ← Spread of data
- 2) Standard Deviation (σ)

Variance

population variance (σ^2)

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

$\bar{x}_i - \mu$ = Distance from
the mean.

Sample variance

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

Why $n-1$?
(Assumption)

Variance

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

$$\bar{x} = 3$$

$$\mu = 3$$

Variance

$$\{1, 2, 3, 4, 5, 6\}$$

$$80$$

$$11$$

$$\bar{x} = 4.4$$

$$\sigma^2$$

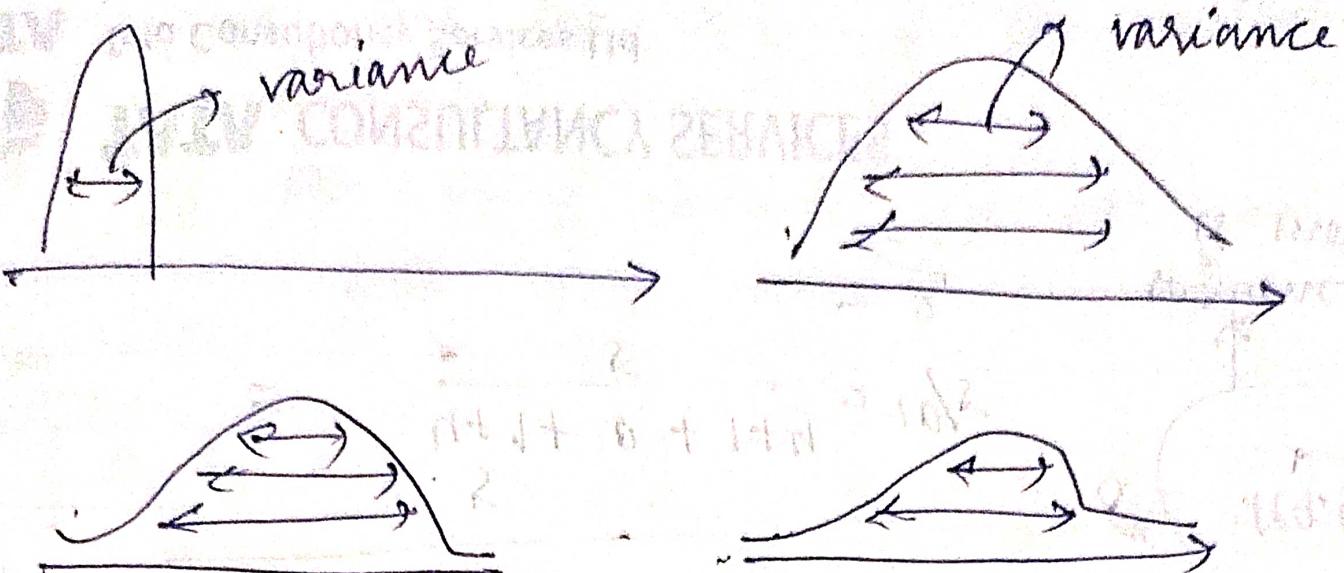
$$\sigma^2 = \frac{(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2}{5}$$

$$\sigma^2 =$$

$$= \frac{4+1+0+1+4}{5} = 10/5$$

$$\sigma^2 = 2.4$$

Variance
is more.

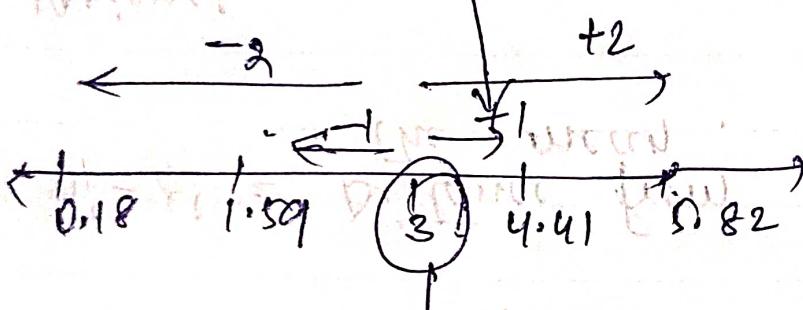


④ Standard Deviation: $(\sqrt{\sigma^2})$

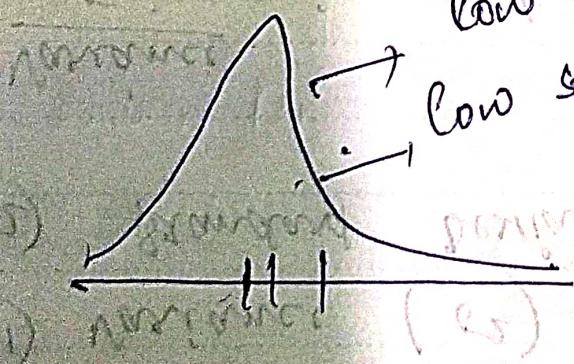
$$\{1, 2, 3, 4, 5\} = \frac{15}{5} = 3$$

$$\sigma^2 = 2$$

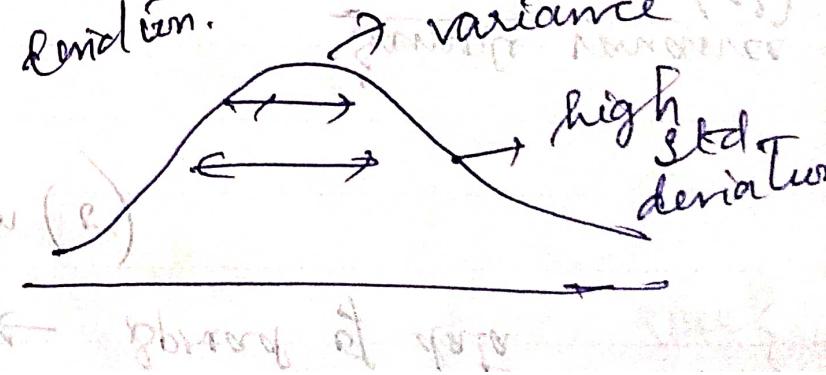
$$\sigma = \sqrt{2} = 1.414$$



How many std. deviation away from the mean? $\Rightarrow +1$



low variance
low std. deviation.



high variance
high std. deviation.

III Percentiles and Quartiles:

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

Percentage of even no. = $\frac{\text{No. of even nos.}}{\text{Total no. of nos.}} = \frac{4}{8} = 50\%$

Percentile = Grade, CAT, IELTS, etc., Num. = 0.5
= 50%

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

Ex:- 99 percentile: It means the person has got better marks than 99 other students.

Dataset: 1, 2, 3, 4, 5, 6, 7, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12 → Should be ascending.

What is the percentage percentile of rank 10?

Percentile rank of $x^{(10)}$ $\Rightarrow \frac{\text{No. of values below } x^{(10)}}{n} = \frac{16}{20} = 80$ percentile

Q: What is the value that exists at 95 percentile?

$$\text{Value} = \frac{\text{Percentile}}{100} \times n$$
$$= \frac{95}{100} \times 20$$

Index starts from 0,

value is at the $= 5^{\text{th}}$ index.

~~for even $= \alpha(n+1)$~~
~~for odd $= \alpha(n)$~~

$$\frac{95}{100} \times 20 = 19^{\text{th}} \text{ index}$$

Number Summary:

- 1) Minimum
- 2) first Quartile
- 3) Median
- 4) Third Quartile
- 5) Maximum.

Use to remove the outliers.

Example: {1, 2, 2, 2, 3, 3, 3, 4, 5, 5, 6, 6, 6, 7, 8, 9, 27}

Thus, we try to create a fence.

[Lower fence \leftrightarrow Higher fence]

lower fence $= Q1 - 1.5 \cdot (\text{IQR})$

$$[-3.65 \leftarrow \rightarrow 14.25]$$

is not in the fence, so this is outlier.

$$\text{Lower Fence} = Q_1 - 1.5(\text{IQR}) \quad \text{IQR} = Q_3 - Q_1$$

$$\text{Higher Fence} = Q_3 + 1.5(\text{IQR}) \quad \downarrow \quad 75 - 25 \\ \Rightarrow 7.5 - 2.5 \\ = 4.5$$

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

$$= \frac{25}{100} \times 21$$

$$= 5.25 \Rightarrow \text{Index}$$

\rightarrow i.e., there are no 5.25 index, so avg

$$\left(\frac{3+3}{2} = 3 \right)$$

$$Q_3 = \frac{75}{100} \times 21 = 15.75$$

15th index

$$+ \frac{8+9}{2} = 8.5$$

Thus,

$$\text{lower fence} = 3 - 1.5 \times 4.5 = -3.65$$

$$\text{higher fence} = 7.5 + 1.5(4.5) = 14.25$$

Thus, 5 number summary:

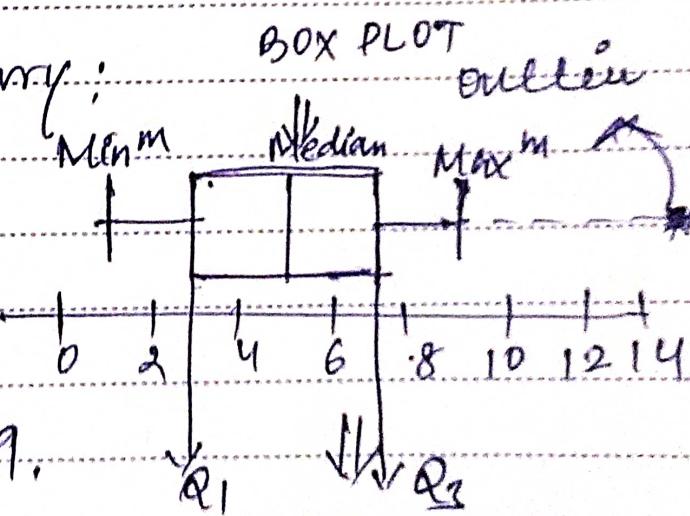
$$1) \text{ Minimum} = 1$$

$$2) Q_1 = 3$$

$$3) \text{ Median} = 5$$

$$4) Q_3 = 7.5$$

$$5) \text{ Maximum} = 9.$$



To treat outliers,