



Statistic Class (Day - 2)

By Krish Naik Sir



Agenda

- (1) Histogram
- (2) Measure of Central Tendency
- (3) Measures of Dispersion
- (4) Percentiles and Quartiles
- (5) Five Number Summary (Box plot)



Histogram

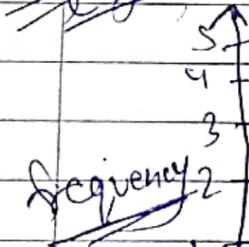
Eg:- Ages = {10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, }
 - , 51, 65, 68, 78, 90, 95, 100 }



Steps

- (a) Sorts the Numbers
- (b) Bins \rightarrow No. of groups
- (c) Bins Size \rightarrow Size of Bins

Eg:-



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

Min = 10

Max = 40

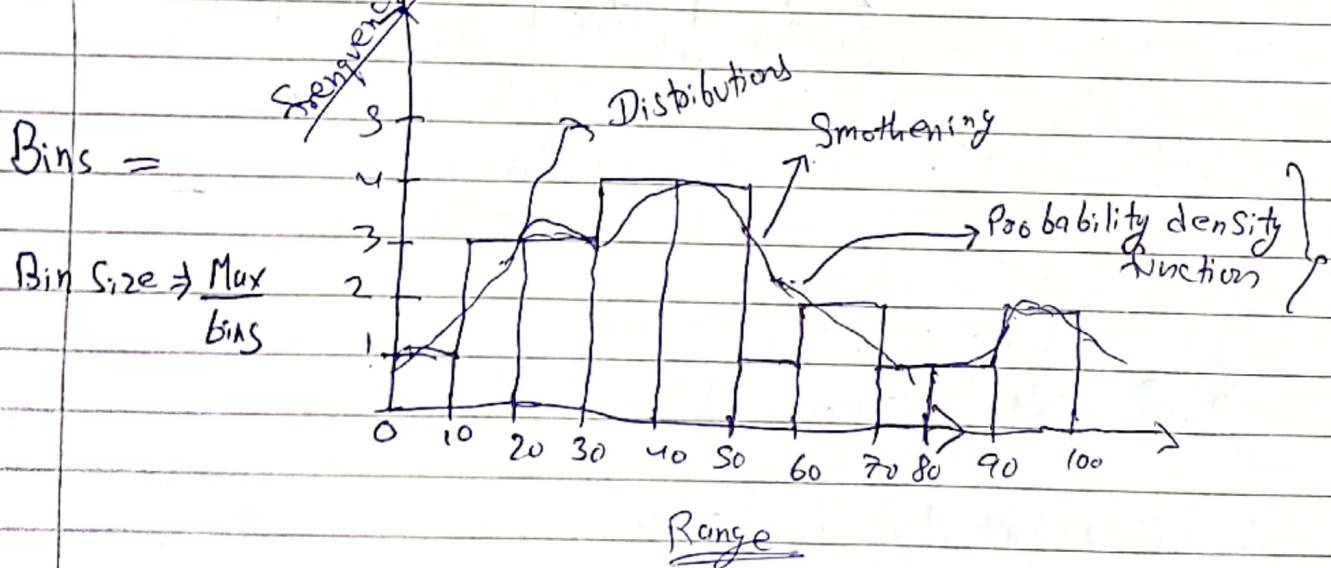
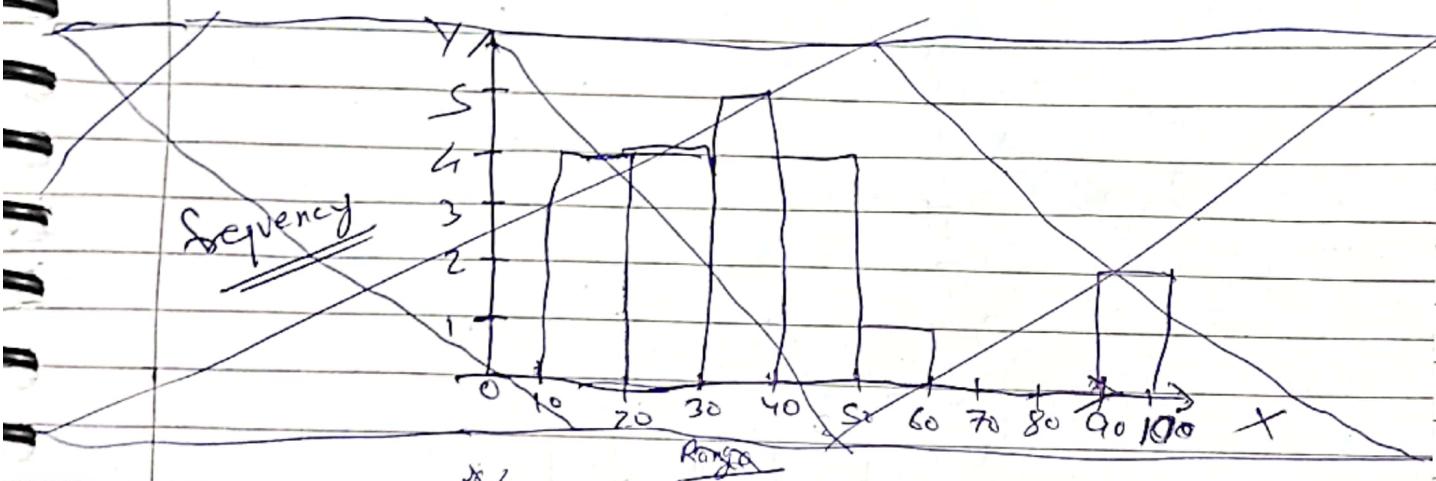
bins = 10

$\frac{40}{10} \rightarrow 4$

Eg:- Ages $\Rightarrow \{10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51, 65, 68, 78, 90, 95, 100\}$

$$\text{bins} = 10$$

$$\text{bin size} = \frac{100}{10} = 10$$



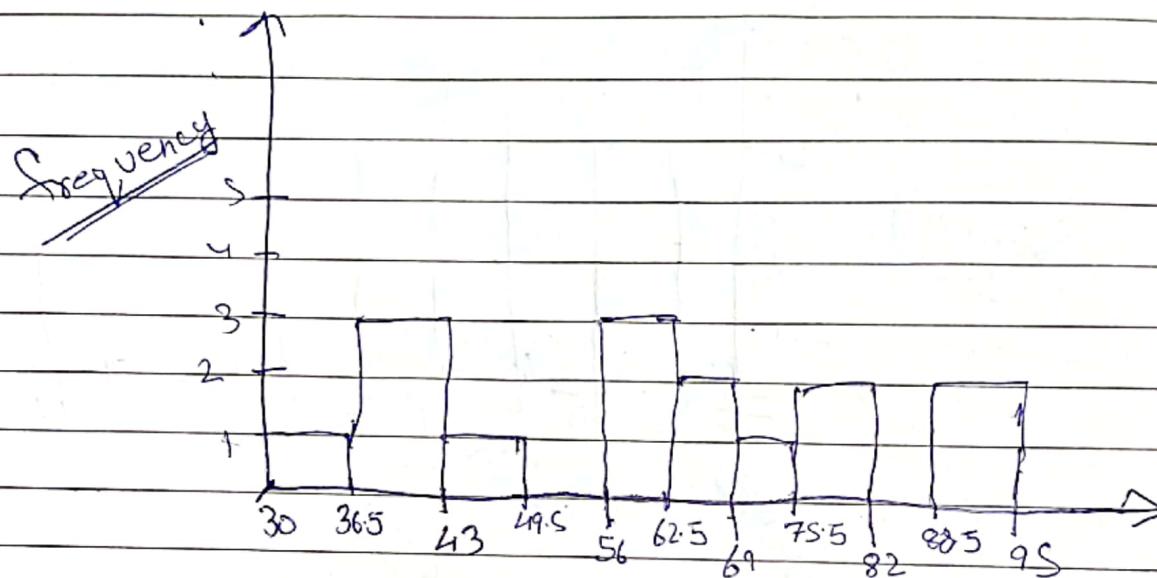
~~Eg:-~~

Weights = $\{30, 35, 38, 42, 46, 58, 59, 62, 63, 68, 75, 77, -80, 90, 95\}$

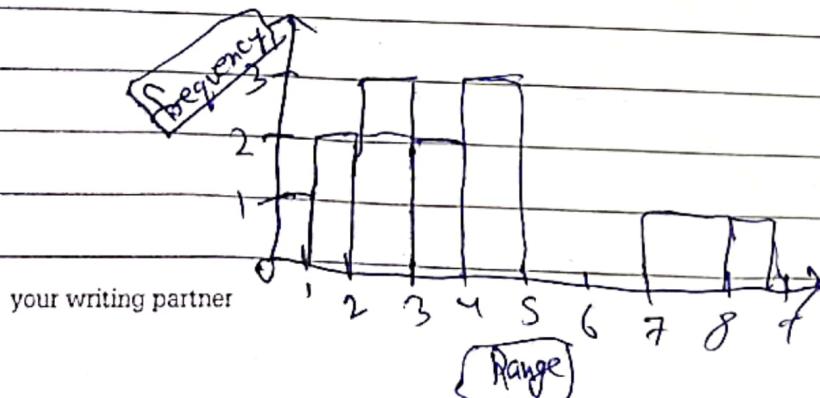
bins = 10

(Continuous Variable)

$$\text{bin size} = \frac{95-30}{10} = \frac{65}{10} = \underline{\underline{6.5}}$$

~~Range~~~~(b) Discrete Continuous:-~~

Eg:- No. of Bank Accounts = $\{2, 3, 5, 1, 4, 5, 3, 7, 8, 3, 2, 4, 5\}$



{ Probability Mass Function }

Pdf \Rightarrow Probability Density function \rightarrow Continuous
Pmf \Rightarrow Probability Mass function \rightarrow Discrete

Subject _____

MON TUE WED THU FRI SAT SUN

② Measure of Central Tendency

① Mean

② Median

③ Mode

A Measure of Central Tendency is a single value that attempt to describe a set of data identifying the central position.

2.1 Mean

Ex. $X = \{1, 2, 3, 4, 5\}$

$$\text{Average / Mean} = \frac{1+2+3+4+5}{5} \Rightarrow \frac{15}{5} \Rightarrow 3$$

~~Sample (n)~~

$N > n$

~~Population (\underline{N})~~

~~Sample (\underline{n})~~

Population mean (μ) =
$$\frac{\sum_{i=1}^N x_i}{N}$$

Sample mean (\bar{x}) =
$$\frac{\sum_{i=1}^n x_i}{n}$$

Subject _____

MON TUE WED THR FRI SAT SUN

~~(S)~~

$$\text{Population mean } (\mu) = \left[\frac{\sum_{i=1}^N x_i}{N} \right]$$

~~Eg:-~~ Population Age = { 24, 23, ~~2~~, 2, 1, 28, 27 }

$$N = 6$$

$$\text{Population mean } (\mu) = \frac{24 + 23 + 2 + 1 + 28 + 27}{6}$$

$$\underline{\mu} = \underline{17.5}$$

~~(A)~~

$$\text{Sample mean } (\bar{x}) = \left[\frac{\sum_{i=1}^n x_i}{n} \right]$$

~~Eg:-~~ Sample Age = { 24, 2, 1, 27 }

$$\boxed{n = 4}$$

$$\text{Sample mean } (\bar{x}) = \frac{24 + 2 + 1 + 27}{4} = \frac{54}{4}$$

$$\underline{\bar{x}} = \underline{13.5}$$

$\text{np.nan} \Rightarrow \text{null value}$

Subject _____

MON TUE WED THR FRI SAT SUN

Eg:

Practical Application [Feature Engineer]

<u>Age</u>	<u>Salary</u>	<u>Family Size</u>
—	—	—
NAN	—	—
—	—	—
—	NAN	—
—	NAN	NAN
NAN	—	—

cols of information

Note: [Where there are NAN value, we will replace]
with mean of that Age & put its →

Eg:

Age | Salary

24	45
28	50
29	NAN
NAN	60
31	75
36	80
NAN	NAN

$$\text{Salary} = 62$$

$$\text{Age} = 29.6$$

$$29.6$$

$$29.6$$

outlier

80

200

$$\text{mean of Age} = \frac{24, 28, 29, 31, 36}{5} = \frac{148}{5} \Rightarrow 29.6$$

AM

your writing partner

2.2

Median

~~(5)~~ $\{1, 2, 3, 4, 5\}$ ~~10~~
 $\bar{x} = 3$

~~(6)~~ $\{1, 2, 3, 4, 5, 100\}$ ^{outlier} ~~10~~
 $\bar{x} = 19.16$

Step to find out Median:-

① Sort the Number

② Find the central Number.

→ ③ If No. of elements are even we find the average of central element

④ If No. of elements are odd we find the central elements

Eg: $x = \{1, 2, 3, 4, [5, 6], 7, 8, 100, 120\} \rightarrow$ Even

$$\Rightarrow \text{median} = \frac{5+6}{2} = \underline{\underline{5.5}}$$

Eg: $x_1 = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 100, 120\}$

your writing partner

$\lceil \text{mean} \Rightarrow 25.6 \rceil$
 $\lceil \text{median} = 5 \rceil$

~~(A)~~ No outliers \rightarrow Mean
~~(A)~~ With outliers \rightarrow Median

(2) (3) **Mode:-** $\Rightarrow \{$ Most frequent occurring element $\}$

$$\text{Eg} \Rightarrow \{ 1, 2, 2, [3, 3, 3], 4, 5 \}$$

$$\text{Mode} = 3 //$$

$$\text{Eg} \Rightarrow \{ 1, [2, 2, 2], [3, 3, 3], 4, 5 \}$$

$$\text{Mode} = 2, 3 //$$

Eg:- **Dataset**

Type of Flowers {Categorical variables}

Lily
 Sunflower
 Rose
 NAN \leftarrow Rose
 Rose
 Sunflower
 Rose
 NAN \leftarrow Rose

your writing partner

Note:- In Categorical Variable, we follow Mode
 for NAN values.

Q Why Sample Variance σ^2 is divisible by $n-1$??

Subject _____

MON TUE WED THR FRI SAT SUN



Measure of Dispersion

- ① Variance (σ^2) ← [Spread of Data]
- ② Standard Deviation (σ)



Variance

Population Variance (σ^2) :-

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

Sample Variance (s^2)

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



{1, 2, 3, 4, 5, 6, 7, 8, 9, 10}



$$\mu = \frac{1+2+3+4+5+6+7+8+9+10}{10} = \frac{55}{10} \Rightarrow 5.5$$

Eg: $x_1 = \{1, 2, 3, 4, 5\}$

$$\mu = 3$$

Eg: $x_2 = \{1, 2, 3, 4, 5, 6, 8, 0\}$

$$\mu = 4.4$$

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

$$= \frac{4+1+0+1+4}{S} = \frac{10}{S} = 2$$

your writing partner

$$\sigma^2 = \frac{(1-4.4)^2 + (2-4.4)^2 + (3-4.4)^2 + (5-4.4)^2 + (6-4.4)^2 + (8-4.4)^2}{S}$$

$$= \frac{(S-4.4)^2 + (6-4.4)^2 + (8-4.4)^2}{S}$$

Subject _____

MON TUE WED THR FRI SAT SUN

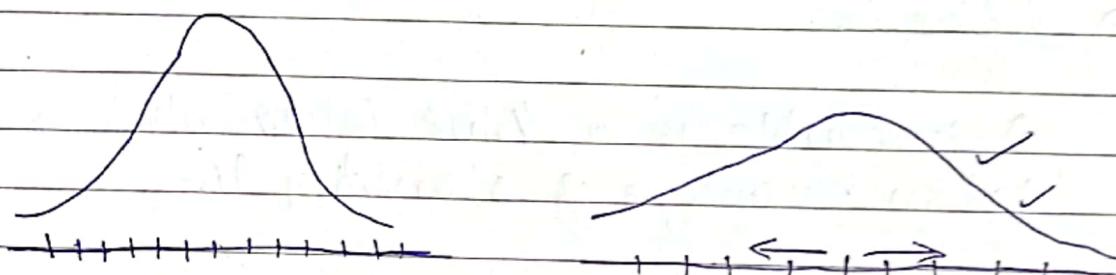
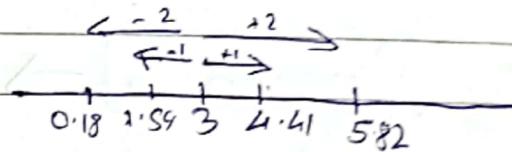
~~Q~~ Standard Deviation ($\sqrt{\sigma^2}$) $\Rightarrow \sigma$

Eg: { 1, 2, 3, 4, 5 }

$$\mu = 3$$

$$\sigma^2 = 2$$

$$\sigma = \sqrt{2} = 1.41$$



~~Q~~

Ans

Percentiles

(A)

Percentiles And Quartiles

(B)

Percentiles \Rightarrow CAT, IIT, JEE, GRE, NEET

Percentile

(C)Definition

A percentile is a value below which a certain percentage of observation lie.

99 percentile :— It means the person has got better marks than, 99 % of the entire mark.

(D)Dataset

Ascending order

$\Rightarrow \{ 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12 \}$

(E)

What is a percentile rank of 10 ?

(F)

Percentile Rank of $n = \frac{\# \text{ No. of Value below } x}{n}$

$$= \frac{16}{20} = \frac{4}{5} \Rightarrow 80 \text{ percentile}$$

your writing partner

Rank 6 $\Rightarrow \frac{6}{20} = \frac{3}{10} \Rightarrow 35\%$

NOTES

Q) Percentile Rank of 8 $\Rightarrow \frac{8}{20} \Rightarrow 55\%$

Q) What is the value that exists at 25 percentile??

$$\Rightarrow \text{Value} = \frac{\text{Percentile}}{100} \times [n+1] \quad \checkmark$$
$$= \frac{25}{100} \times 20.8 = \underline{5^{\text{th}} \text{ Index}}$$

Q) Output value is 5 //

~~Q) 95 % value = $\frac{95}{100} \times 21 = 19.95$~~



Five Number Summary (Box Plot)

① Minimum

first Quartile (25 percentile) Q_1

Second Quartile ~~Q_2~~ Median "

Third Quartile (75 percentile) (Q_3)

⑤ Maximum

$$\frac{3+3}{2} = 3$$

Eg: $\{1, 2, 2, 2, 2, \boxed{3}, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27\}$

[Lower fence \longleftrightarrow Higher fence]

$$\text{Lower fence} = Q_1 - 1.5(\text{IQR})$$

$$\text{Higher fence} = Q_3 + 1.5(\text{IQR})$$

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

Inter Quantile Range (IQR)

$$25\% Q_1 = \frac{25}{100} * (n+1) = \frac{25}{100} * 21 = 5.25 \Rightarrow \text{Index} \Rightarrow 3$$

$$75\% Q_3 = \frac{75}{100} * (n+1) = \frac{75}{100} * 21 = 15.75 \text{ Index} \Rightarrow \frac{8+7}{2} = 7.5$$

$$\Rightarrow \text{Lower fence} \Rightarrow 3 - (1.5)(4.5) = \boxed{-3.65}$$

$$\Rightarrow \text{Higher fence} \Rightarrow 7.5 + (1.5)(4.5) = \boxed{14.25}$$



Five Number Summary for Box-Plot 3

previous
Page One Pg. 1

- ① Minimum = 1
- ② Q_1 = 3
- ③ Median = 5
- ④ Q_3 = 7.5
- ⑤ Maximum = 9

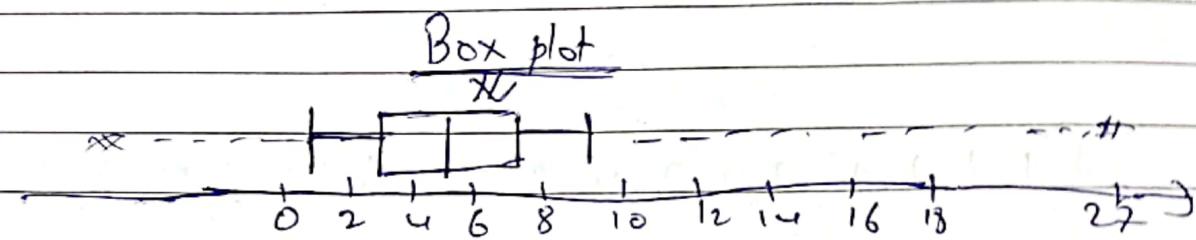


Fig:- Plotting the Box-Plot with above "Five Number summary Information"

X X X X