

# Day 2 - Statistics

## Agenda

- ① Histograms ✓
- ② Measure of Central Tendency ✓ } ← 1.15 hrs
- ③ Measure of Dispersion ✓
- ④ Percentiles And Quartiles }
- ⑤ 5 Number Summary (Box plot) . }

## ① Histogram

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

① Sort the Numbers

$$\text{frequency (count)} [10, 20, 25, 30, 35, 40]$$

$$\min = 10$$

$$\max = 40$$

② Bins → No. of groups

$$\text{bins} = 10$$

$$\frac{40}{10} = 4 //$$

③ Bins size → Size of Bins

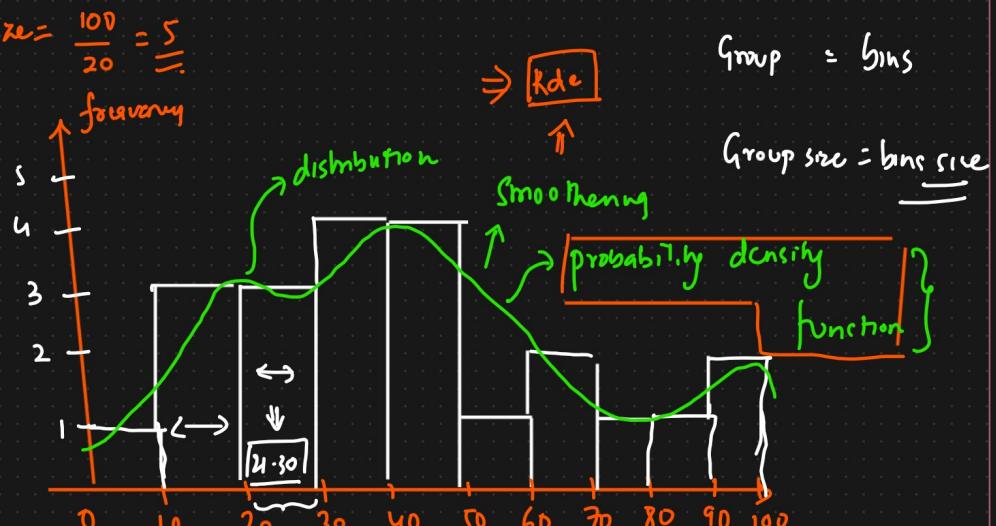


$$\underline{\underline{\text{bins} = 10}} \quad \underline{\underline{\text{bin size}}} = \frac{100}{20} = 5$$

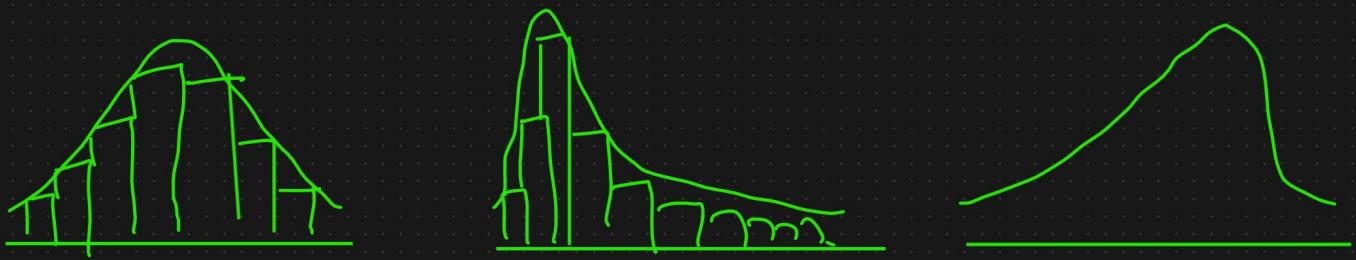
$$\text{bin size} = \frac{100}{10} = 10$$

$$\text{Bin size} = \frac{\text{Max} - \text{Min}}{\text{bins}}$$

$$\text{bins} =$$



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



Assignment  
Weight = { $\boxed{30}, 35, 38, 42, 46, 58, 59, 62, 63, 68, 75, \boxed{80}, 90, \boxed{95}$ }.

bins = 10

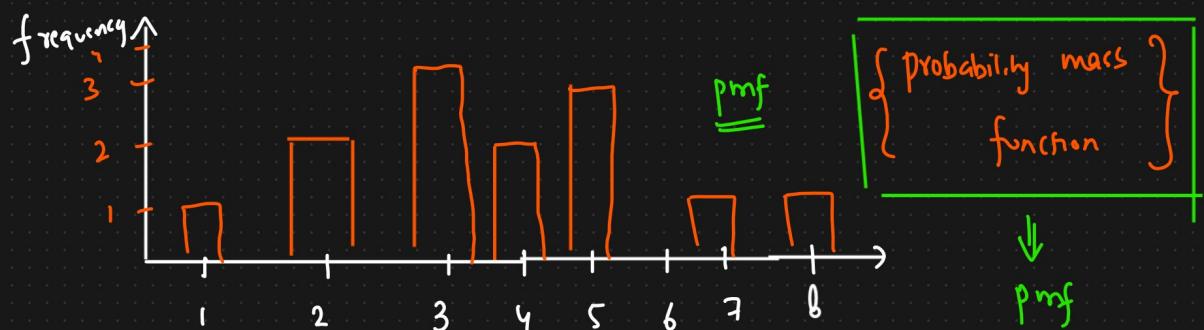
$$\text{bin size} = \frac{95 - 30}{10} = \frac{65}{10} = 6.5$$

{continuous value}

pdf

## ② Discrete

No. of Banks accounts = [2, 3, 5, 1, 4, 5, 3, 7, 8, 3, 2, 4, 5]



pdf : probability density function }  $\rightarrow$  continuous

pmf : probability mass function }  $\rightarrow$  discrete.

## ① Measure of Central Tendency

- ① Mean, ✓      { A measure of CT is a single value that attempts to describe a set of data identifying the central position
- ② Median
- ③ Mode.

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

Population ( $N$ )

$$N \gg n$$

Sample ( $n$ )

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

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

$$N = 6$$

$$N > n$$

$$n = 4$$

$$\{24, 23, 28, 27\} \leftarrow \text{Age}$$

$$\text{Population Age} = \{24, 23, 2, 1, 28, 27\}$$

$$\text{Sample Age} = \{24, 2, 1, 27\}$$

$$\begin{array}{r} 13 \\ 54 \\ \hline 41 \end{array}$$

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

$$\mu = 17.5$$

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

$$\begin{cases} \mu > \bar{x} \\ \bar{x} > \mu \end{cases}$$

$$\bar{x} = 13.5$$

$\boxed{\text{hp-null}}$   $\leftarrow$  Null values

### Practical Application (Feature Engineering)

Age	Salary	Family Size
-	-	-
-	-	-
-	-	-
$\rightarrow \text{NAN}$	-	-
-	-	-
-	$\text{NAN}$	-
-	-	$\text{NAN}$
-	$\text{NAN}$	-
$\rightarrow \text{NAN}$	-	-

$\leftarrow$  loss of Info

$$\boxed{\text{Age} = 29.6}$$

$$\downarrow \downarrow \downarrow$$

$$\boxed{38}$$

Mean

$$\begin{array}{l} \boxed{\text{NULL}} \quad \downarrow \quad \text{val} = 10/6 \\ 10/4 = \boxed{1, 2, 3, 4} \quad \uparrow \quad \boxed{\text{NAN}} \\ \hline \text{mean} \quad \uparrow \quad \boxed{\text{NAN}} \end{array}$$

Age	Salary
24 ✓	45
28 ✓	50
29 ✓	$\boxed{\text{NAN}}$
$\boxed{\text{NAN}}$	$\boxed{\text{NAN}}$
31 ✓	$\boxed{\text{NAN}}$
36 ✓	$\boxed{\text{NAN}}$
$\boxed{\text{NAN}}$	$\boxed{\text{NAN}}$
	$\boxed{62}$
	$\downarrow \downarrow \downarrow$
	$\boxed{85}$

$$\text{Outliers} \leftarrow [80] \quad [200] \leftarrow$$

## ① Median

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

$\bar{x} = 3 \longrightarrow \bar{x} = 19.16$

Outlier  
 $\frac{1+2+3+4+5+100}{6} = \frac{115}{6} = 19.16$

### Steps to find out median

① Sort the Numbers

② Find the central number

① if the no. of elements are even we find the average of central elements

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

Sorted

$$\{0, 1, 2, 3, 4, \boxed{5, 6}, 7, 8, 100, 120\}$$

Mean =  $\frac{25.6}{10}$

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

median = 5

③ Mode : {Most frequent occurring elements}

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

$$\boxed{2, 3}$$

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

## Dataset

Types of flower (categorical variable).

Flower

Sunflower

Rose

NAN ← ROSE

Rose

Sunflower

Rose

NAN ← Rose

40%

FF  
Biased

Under 19

17, 18, 19, 16, 15, 32 → Outlier

## (F) Measure of Dispersion

① Variance ( $\sigma^2$ ) ← Spread of Data.

② Standard deviation ( $\sigma$ ) ←

$$X = \{1, 2, 3, 4, 5\} \quad \mu = 3$$

### Variance

Population Variance ( $\sigma^2$ ) { Degree of freedom }

$$\sigma^2 = \frac{N}{N} \sum_{i=0}^N (x_i - \mu)^2 \quad \{ \text{Bands Correction} \}$$

Sample Variance ( $s^2$ )

$$s^2 = \frac{n}{n-1} \sum_{i=0}^n (x_i - \bar{x})^2$$

$$\begin{aligned} \{1-100\} & \Downarrow \\ & = \Downarrow \text{First one.} = \{1-100\} \end{aligned}$$

Second one.  
↓

↓  
Assignment ←

$$\left\{ 1, 2, 3, 4, 5, 6, 7, 8, 9, 10 \right\} \xrightarrow{\text{Variance}} \boxed{\text{Variance}}$$

$$\left\{ 1, 2, 3, 4, 50, 60, 70, 100 \right\} \xrightarrow{\text{Variance}} \boxed{\text{Variance}}$$

Variance Given

$$\frac{21}{80} = \frac{101}{7}$$

$$\left\{ 1, 2, 3, 4, 5 \right\} \xrightarrow{\text{Mean}} M = 3$$

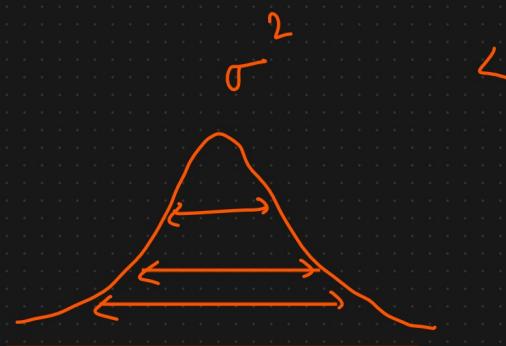
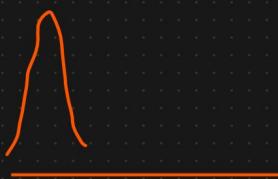
$$\left\{ 1, 2, 3, 4, 5, 6, 80 \right\} \xrightarrow{\text{Mean}} M = 14.4$$

$$\sigma^2 = \frac{(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2}{5} = \frac{4+1+0+1+4}{5} = \frac{10}{5} = 2$$

$$\sigma^2 = \frac{(1-14.4)^2 + (2-14.4)^2 + (3-14.4)^2 + (4-14.4)^2 + (5-14.4)^2 + (6-14.4)^2 + (80-14.4)^2}{7}$$

**Variance ↑↑ Spread ↑↑**

$$\sigma^2 = 719.10$$



④ Standard deviation  $(\sqrt{\sigma^2}) \Rightarrow \boxed{4}$

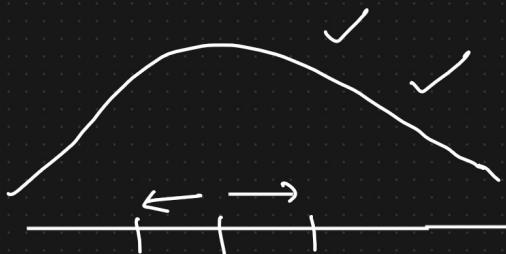
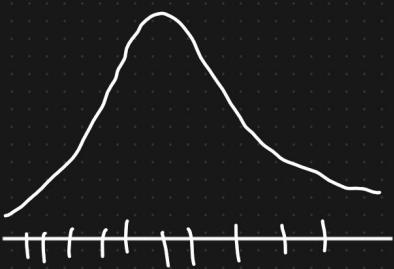
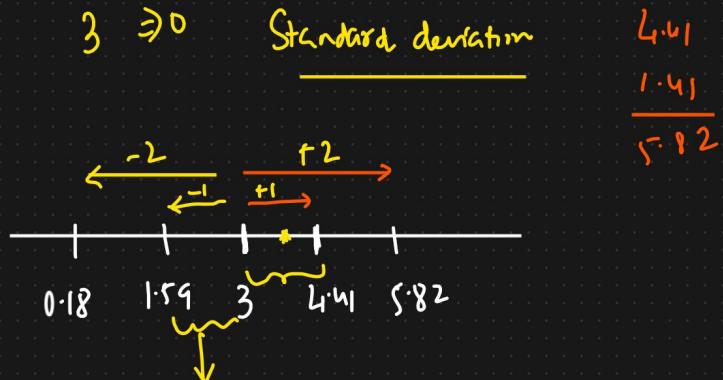
$$\begin{array}{r} 3.00 \\ 1.41 \\ \hline 1.59 \end{array}$$

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

$$M = 3$$

$$\sigma^2 = 2$$

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



## ④ Percentiles And Quartiles

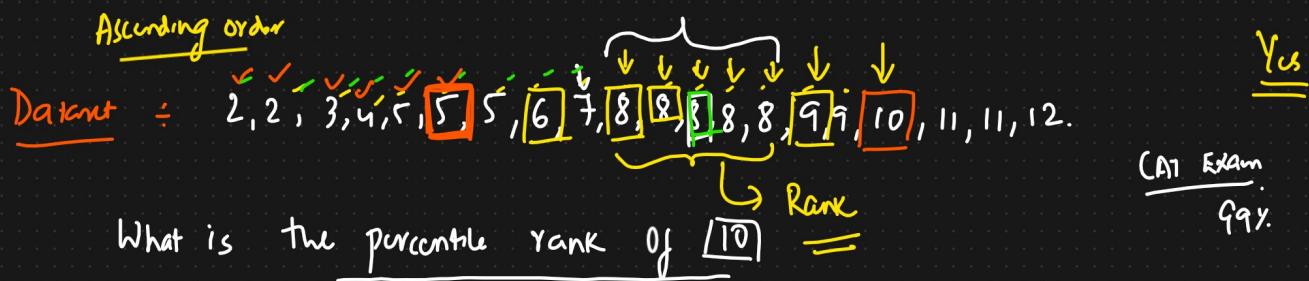
$$\text{Percentage} = \{ 1, 2, 3, 4, 5, 6, 7, 8 \}$$

$$\text{Percentage of Even Number} = \frac{\text{No. of even Numbers}}{\text{Total no. of Number}} = \frac{4}{8} = 0.5 = 50\%$$

Percentiles : CAT, IAT, JEEL, SAT, GRE, JEE, NEET  $\Rightarrow$  Percentiles

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

99 percentile = It means the person has got better marks than 99% of the entire students



(CAT Exam)  
99%

= 0.8

Next item

$$\text{Percentile Rank of } x = \frac{\# \text{No. of Value below } x}{n} = \frac{16}{20} = 80 \text{ percent.}$$

45 percentile

$$= \frac{14}{20} = 70 \text{ percent.}$$

④ What is the value that exists at 25 percentile

75%

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

$$= \frac{25}{100} \times 20 = \frac{5^{\text{th}} \text{ Index}}{20}$$

$$\text{Or } p = 5$$

$$= \frac{95}{100} \times 21$$

⑥ 5 number Summary

① Minimum

② First Quartile (25 percentile) (Q1)

③ Median

④ Third Quartile (75 percentile) (Q3).

Box plot

⇒ Remove the outliers.

### ⑤ Maximum

$$\{1, 2, 2, 2, 3, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, \underline{\underline{12, 14}}\} \quad \downarrow \text{outlier}$$

$\frac{\downarrow}{15} \frac{\downarrow}{16} =$   
 $\frac{\downarrow}{5.25}$



[Lower Fence  $\longleftrightarrow$  Higher Fence]

$\underline{\underline{}}$

$$\downarrow [-3.65 \longleftrightarrow 14.25]$$

$$\leftarrow \text{lower Fence} = Q_1 - 1.5(IQR) \leftarrow$$

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

75 25

$$IQR = Q_3 - Q_1$$

$\downarrow = =$

Inter Quartile Range (IQR)

$$Q_1 = \frac{25}{100} \times 21 = 5.25 \quad \text{Index} = 3 =$$

$$Q_3 = \frac{75}{100} \times 21 = 15.75 \quad \text{Index} = \frac{8+7}{2} = \underline{\underline{7.5}} =$$

$$\text{lower Fence} = 3 - (1.5)(4.5) = \underline{\underline{-3.65}}$$

$$\text{Higher Fence} = 7.5 + (1.5)(4.5) = \underline{\underline{14.25}}$$

$$\{1, 2, 2, 2, 3, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, \underline{\underline{12, 14}}\}.$$

-5

5 Number Summary.

① Minimum = 1 ✓

②  $Q_1 = 3$  ✓

③ Median = 5 ✓

④  $Q_3 = 7.5$  ✓

⑤ Maximum = 9 ✓

Box Plot

$\downarrow$



$\downarrow$

To Treat Outliers