

start in PA

## 1. What is Statistics?

→ It is a science of collecting, organizing, summarizing, and analyzing data to� make decision making.

## 2. What is Data?

→ Facts or pieces of information that can be organized

Ex: 1st of class

8, 9, 10, 11, ... ?

age of students

20, 22, 21, ... ?

## # Types of Stats

### 1. Descriptive Statistics :-

It consists of organizing and summarizing data.

### 2. Inferential Statistics :-

Using Data, we can make conclusion using some techniques.

Ex: class → 20 Students

1<sup>st</sup> sem maths 8, 9, 10, 9, 5, ... ?

1. What is the avg of class? (Descriptive)

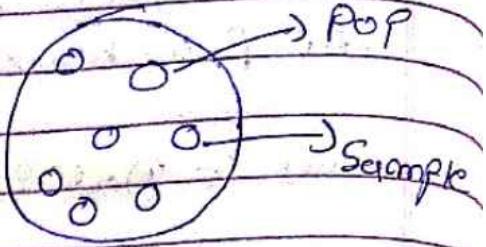
2. 7<sup>th</sup> sem

(Inferential)

## Sample and population

Population =  $N$

whole dataset is known as population



Sample =  $n$

~~small subsets~~

small subsets of data taken from population  
→ Sample

## Sampling Techniques:-

### 1. Simple Random Sampling:-

Every member of population has an equal chance of getting selected in sample ( $n$ )

### 2. Stratified Sampling :-

overlapping groups

Splitting data into non-

Age group

0 - 20

20 - 40

Gender

M

F

### 3. Systematic Sampling:

#### No Convenience Sampling:

Surveying over ( $N'$ ) individuals

where who have knowledge of that particular domain.

1. Exit poll? (Random)

Survey of house hold (with woman) [Convenience]

2.

### Variables

→ It is a property that can hold / take any value

Age: 8, 10, 15, 20, 25, ...

Markes: 70, 80, 95, ...

### Types:

#### 1. Qualitative :- Categorical Values

Based on some characteristic, we can derive categorical values.

IQ :- 0 - 10 → low →

10 - 50 Avg

50 - 70 Good

2. Quantitative:-

Numerical Value (measurable numerically)

Height: {162, 159, 155...}

Weight: {59, 65, 71...}

Discrete (int)

Continuous (float)

Whole no

Decimal NO

1. No. of students

1. Height

{165.2, 167.9...}

2. No. of bank ac.

2. weight

{165.5, 160.9...}

1. Blood pressure → Continuous / discrete

2. Methylated spirit → Qualitative

3. River length → Conti

4. Song length → Conti

5. Gender → Qualitative (category)

## # Variable measurement Scales

~~order matters~~

1. Ordinal : ordered [rank, graduation]
2. Nominal : Categorical values (colors, classes, degrees)
3. Interval : [No zero / absolute point] (order as well as value matter)
4. Ratio : - Zero means nothing

Interval

$$\begin{array}{|c|} \hline 20^\circ\text{C} : 40^\circ\text{C} \\ \hline \cancel{-----} \\ \hline 1^\circ\text{C} : 2^\circ\text{C} \\ \hline \end{array}$$

Ratio

$$\begin{array}{|c|} \hline 20\text{kg} : 40\text{kg} \\ \hline 1 : 2 \\ \hline \end{array}$$

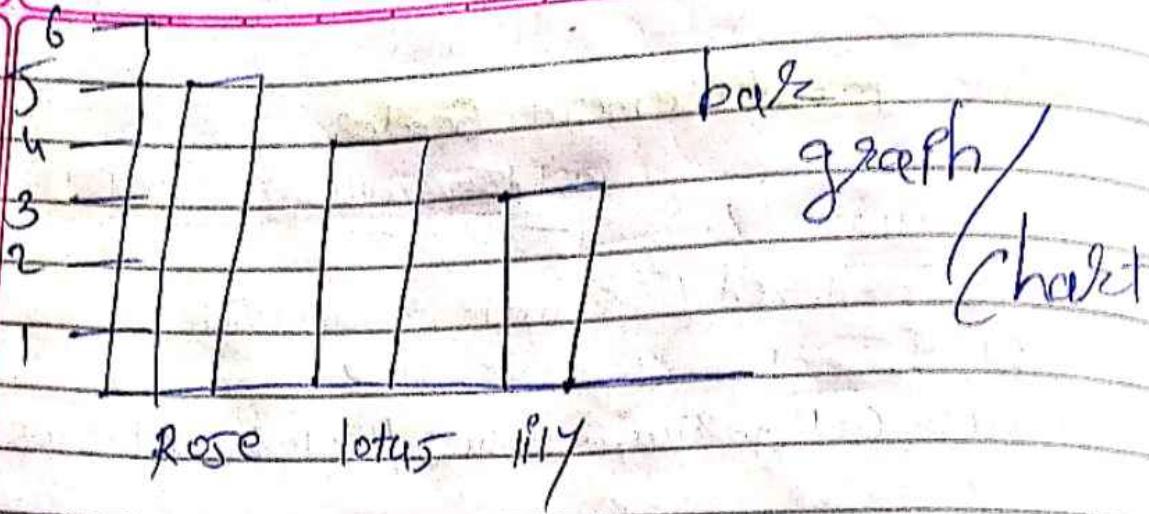
## # Frequency

Data : Flowers

[Rose, lily, lotus, Rose, lily, 8/5P,

Rose, lotus, lily, lotus, lily, lotus

Flower	Frequency	Cumulative
Rose	3	3
lily	3	6
lotus	3	12



bare  
graph/  
chart

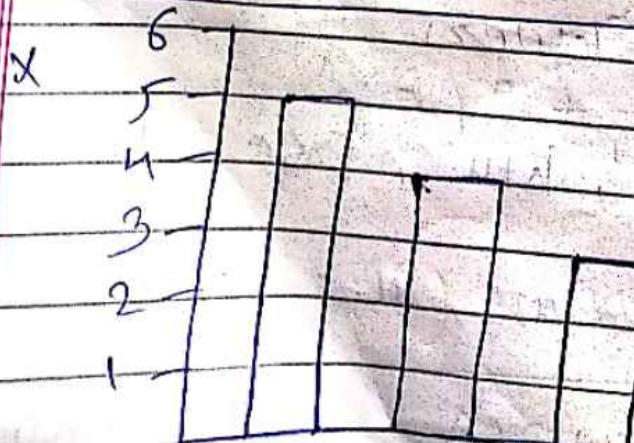
~~# histogrf~~

Marks : [2, 5, 12, 15, 21, 28, 35, 35, 36, 36, 39, 42]

bin

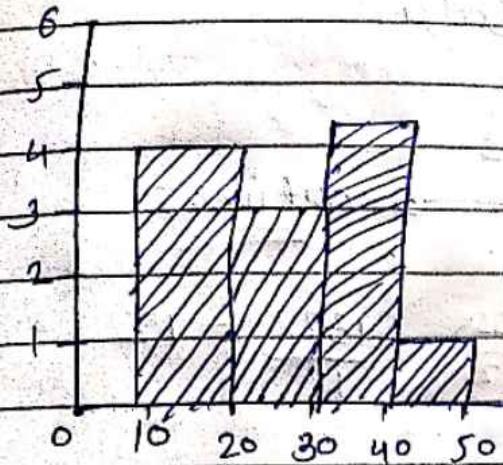
=  $[0-10]$  - 2  
 $[10-20]$  - 2  
 $[20-30]$  - 3  
 $[30-40]$  - n  
 $[40-50]$  - 2

→ Continuous



bare  
graph/  
chart

Rose lotus lily



→ Measure of Central Tendency

→ Measure of Dispersion

→ Distribution

## # Measure of Central Tendency (5)

Avg → mean

Pop

Sum

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

2, 3, 5, 3, 2, 1, 3

$$\frac{2+3+5+3+2+1+3}{7}$$

mean:- It refers to the measure used to represent the distribution.

of the Data.

S 1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100 } L  $\rightarrow$  outliers

$$\frac{32}{10} \rightarrow 3.2 \rightarrow \frac{132}{11} = 12$$

# Medium  $\rightarrow$  middle value

$\hookrightarrow$  in ascending order  
 $\hookrightarrow$  even  
 $\hookrightarrow$  data  $\leftarrow$  odd

Even :-

$$\frac{\left(\frac{n}{2}\right)^{th} + \left(\left(\frac{n}{2}\right)^{th} + 1\right)^{th}}{2}$$

dataset : {11, 12, 13, 14, 15, 16}

$$n = 6$$

$$\frac{\left(\frac{6}{2}\right)^{th} + \left(\frac{6}{2}\right)^{th} + 1}{2} = \frac{3^{th} + 4^{th}}{2}$$

$$\frac{\left(\frac{6}{2}\right)^{th} + \left(\frac{6}{2}\right)^{th} + 1}{2} = \frac{13 + 14}{2}$$

$$\frac{3^{th} + (3 + 1)^{th}}{2} = \frac{27}{2}$$

$$\frac{3^{th} + (3 + 1)^{th}}{2} = 13.5$$

odd:

$$\left(\frac{m+1}{2}\right)^{th} \quad \frac{(m+1)^{th}}{2}$$

$\{11, 12, 13, 14, 15\}$

$$m=5$$

$$\frac{(5+1)^{th}}{2} = \frac{6^{th}}{2} = 3 = 13$$

$\rightarrow \{11, 12, 13, 14, 15, 100\}$

$$\left(\frac{m}{2}\right)^{th} + \left(\left(\frac{m}{2}\right)^{th} + 1\right)^{th}$$

2

$$\frac{(6)^{th}}{2} + \left(\left(\frac{6}{2}\right)^{th} + 1\right)^{th}$$

2

$$\frac{3^{th} + 4^{th}}{2}$$

$$= \frac{13 + 14}{2}$$

$$= 13.5$$

5, 21, 23, 25, 29, 32, 100 } 3

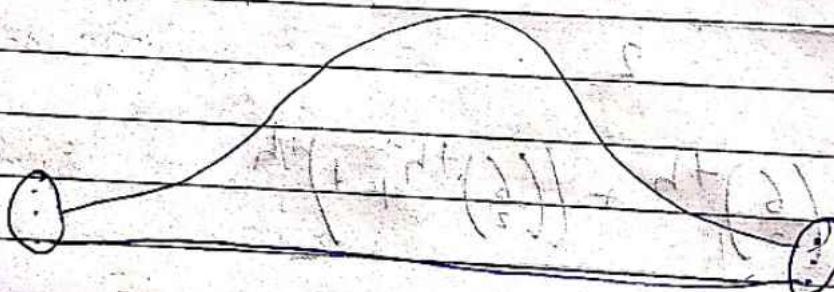
medium  
workers well  
with outside

$$\frac{25 + 29}{2} = \frac{54}{2} = 27$$

outliers:-

a data point who doesn't follow pattern or trend of the data set than it is considered as outlier

[are extreme point]



## # Mode

→ most frequent value (repeated)

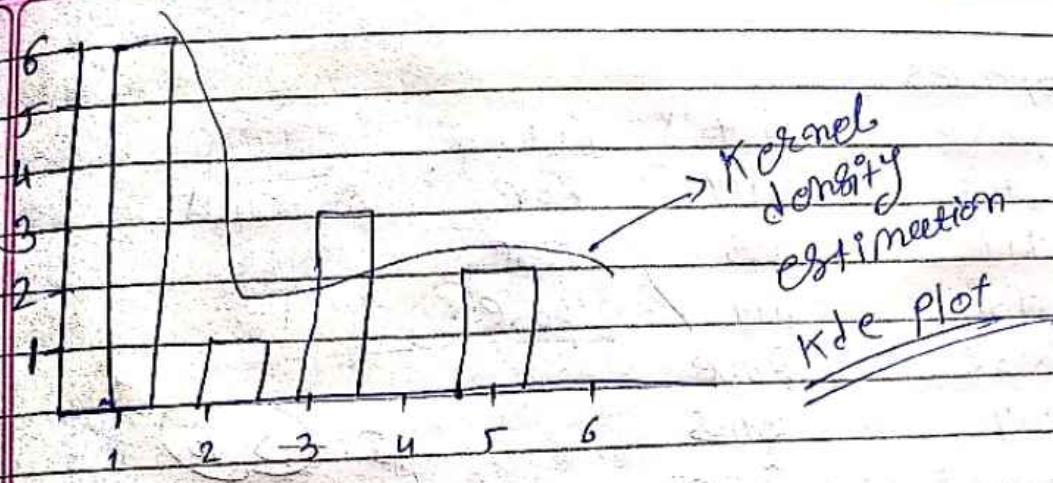
[1, 1, 2, 3, 5, 1, 1, 3, 1, 3, 5, 1, 1]

1 = 6

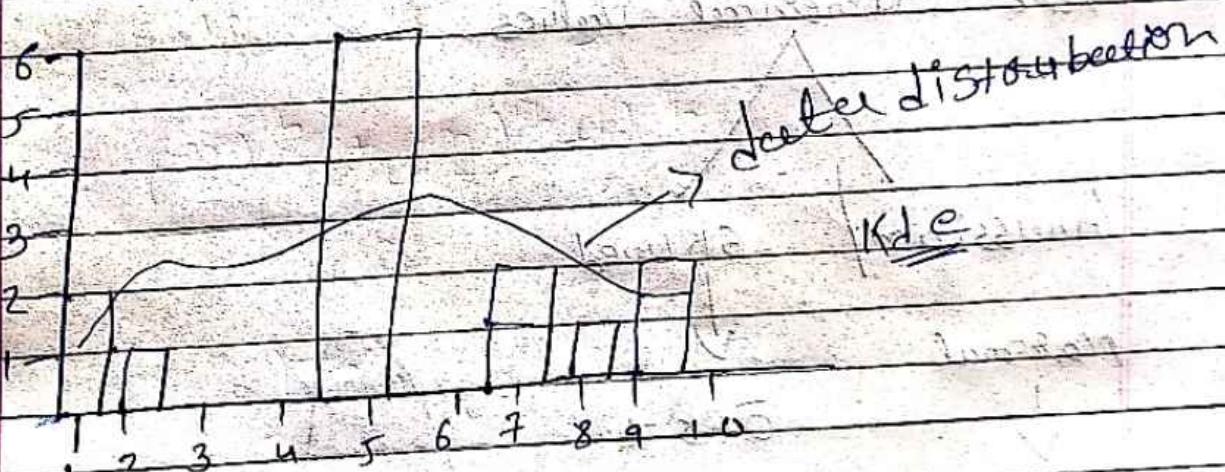
3 = 3

2 = 1

5 = 2



2, 1, 1, 2, 5, 9, 5, 7, 8, 9, 5, 9, 5 }



for categorical missing data

0-5%  $\rightarrow$  mode

$\rightarrow$  new category "missing"  
"unknown"

or

$\rightarrow$  "random"

Species

Rose	1114	L	1114	
lotus	1114		Rose	mult
1114		Moss	Rose	
misi	1114		Rose	
Rose	Rose			
lily	lotus		15 → 10	
Rose	Rose			

False numerical values

Gaussian / Skewed

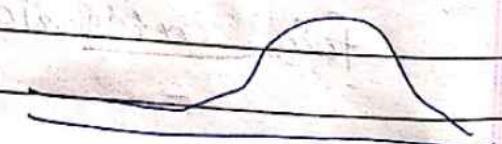
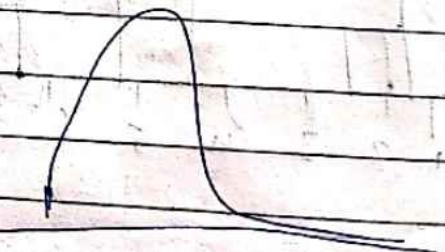
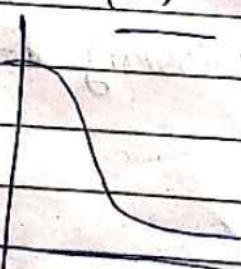
Normal



Mean

0-5-10

10-1-



mean = median = mode

## # measure of Dispersion

↳ spread

$$\text{① } \overline{s} = \frac{\sum |x_i - \bar{x}|}{n}$$

Ex:  $\overline{s} = \frac{|5-1| + |5-1| + |5-1| + |5-1| + |5-1|}{5} = \frac{20}{5} = 4$

$$\text{② } \overline{s} = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n}}$$

Ex:  $\overline{s} = \sqrt{\frac{(2-3)^2 + (2-3)^2 + (2-3)^2 + (2-3)^2 + (2-3)^2}{5}} = \sqrt{\frac{5}{5}} = \sqrt{1} = 1$

## → Variance

$\sigma^2$  measures how far the numbers in a dataset are from the mean (avg).

(How each value differs from a dataset in mean)

High Variance → more spread (far from mean)

low Variance → closer to mean

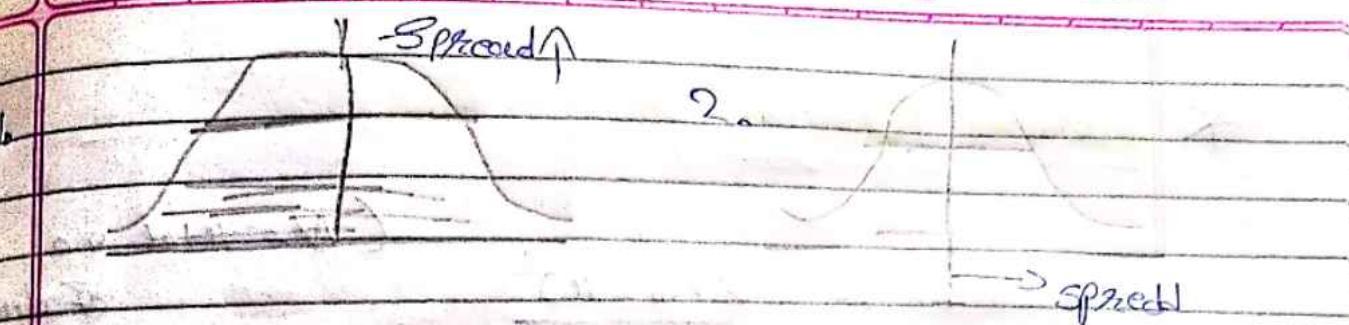
Pop

Sum

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

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

Bessel's  
correction



## Standard Deviation

- Just a squared sort of Variance
- it gives you a measure of spread that is in the same units as the original data, making it easier to interpret than variance.

[Same unit, easily comparable]

Population

Sample

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

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

## Key point (Variance and SD)

- Variance gives you the average squared distance from the mean.
- SD gives you a measure of spread in the same unit

*n* = data point

*N* = total no. of population

*n* = total no. of sample

$\mu$  = mean of population

$\bar{x}$  = mean of sample

Q) Why  $(n-1)$  is used in formula of Sample

	$n$	$x - \bar{x}$	$(x - \bar{x})^2$
1		-1.83	3.34
2		-0.83	0.68
3		0.83	0.68
4		0.17	0.02
5		0.17	0.36
	2.83	2.17	4.70
			10.78

$$S^2 = \frac{10.78}{6} = 1.79$$

$$S = 1.33$$

representative of it. This correction ensures that the sample variance is an unbiased estimator of the population variance.

out + ပို့စ်

~~percentage~~  $\rightarrow \frac{\text{sum of nos. of interest}}{\text{sum of all data pts}} \times 100$

1,2,3,4,5

% of the nos. that are odd?

$$\frac{3}{5} = 0.6 \rightarrow \times 100$$

## Percentile

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

[95, continuing]

50.85 → 6 - intese

Letter Sel = 9, 2, 3, 11, 5, 1, 1, 6, 7, 18, 18, 18, 18, 18, 9, 9

Marks → 10, 11, 11, 12  
15

$$n = 20 \quad \text{ceiling}$$

## \* Variance formula

( $N$  = total no  
of observations)

$$\rightarrow \text{Population} : - \frac{\sum (x_i - \mu)^2}{N}$$

When we have data from the entire population we use ' $N$ ' in the denominator. This gives us an exact measure how the data fits around the population mean ( $\mu$ ).

$$\star \text{ Sample} : - \frac{\sum (x_i - \bar{x})^2}{n-1}$$

$\rightarrow$  When we're working with a sample, only have sample mean  $\bar{x}$ , which is an estimate of the population mean  $\mu$ .

$\rightarrow$  Using sample mean in calculations tend to make the variance slightly smaller than the true population variance.

$\rightarrow$  To correct this bias (underestimating variance) we divide by ' $n-1$ ' instead of  $n$ . This makes variance an estimate larger and accurate.

$\rightarrow$  By subtracting 1 from  $n$ , we account from the fact that the ' $\bar{x}$ ' is not perfectly

Percentile rank =  $\frac{\text{no. of values below } n}{m} \times 100$

$$\frac{16}{20} = \frac{4}{5} \times 100 = 80$$

$$\frac{17}{20} \times 100 = 85$$

Q. What value exists at percentile marking of 75% in a distribution of 20 items?

$$\text{Value} = \left( \frac{\text{Percentile}}{100} \times n \right) + 1$$

$$= \left( \frac{75}{100} \times 20 \right) + 1$$

$$= 5 + 1 = 6$$

$\downarrow$  Index

5

75 percentile

$$\text{Value} = \left( \frac{75}{100} \times 20 \right) + 1$$

$$15 + 1 = 16$$

## \* Fire Number Summary

1. Minimum

2. First Quartile ( $Q_1$ )

3. Median

4. Third Quartile ( $Q_3$ )

5. Maximum

[to detect

outliers]

[1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 15, 27]

[lower fence  $\rightarrow$  higher fence]

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

higher fence =  $Q_3 + 1.5 \text{ IQR}$

$Q_1 = 2.5$  percentile (25%)

$Q_3 = 7.5$  percentile (75%)

IQR = Inter Quartile Range

$$= Q_3 - Q_1$$

$$Q_1 = \left( \frac{25}{100} \times 20 \right) + 1$$

8

6

$$Q_1 = 3$$

$$Q_3 = \left( \frac{75}{100} \times 20 \right) + 1$$

7

$$= 16$$

$$\text{IQR} = 8$$

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

$$= 3 - 1.5(5)$$

$$= 3 - 7.5$$

$$= -4.5$$

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

$$= 8 + 1.5 \text{ IQR}$$

$$= 8 + 7.5$$

$$= 15.5$$

Examining data

1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 7, 8, 8, 9, 15

5 No. Summary

min = 1

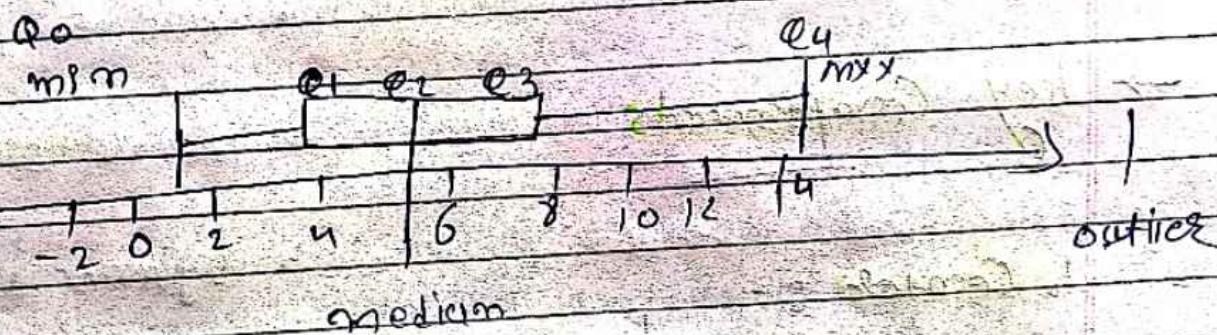
Q1 = 3

median = 5

Q3 = 8

max = 15

Boxplot



## \* Data distribution

- A distribution refers to the way in which values or data pts are spread or arranged.
- It shows how often different values occur in data set and describes the overall pattern of the data.

## → key Components

### 1. Center :-

where the middle data lies  
(mean, median)

### 2. Spread :-

How wide or narrow the data  
(range, variance, SD)

### 3. Shape:-

The overall form of the distribution  
(symmetric, skewed, etc)

### 4. Outliers:-

① data points that are much higher or lower than the rest of the data.

## Types of Distribution

### 1. Gaussian / Normal Distribution

what?

- A symmetric bell shape where most data pts cluster around the mean (center) and fewer values occur as you move away from the mean in both directions.

when?

- when data is evenly distributed around the mean and follows a natural pattern.

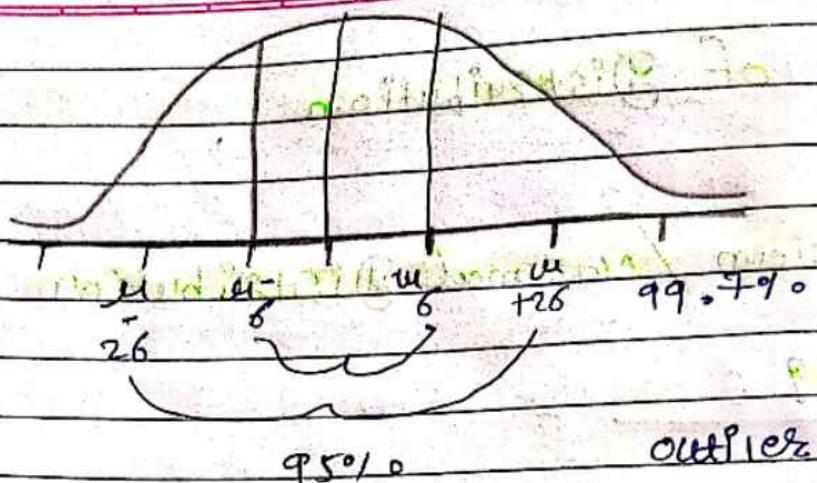
why?

(height, test scores etc)

- It's useful because many natural phenomena follows this pattern and it's a foundation for many statistical methods

where?

- Education, Biology and Quality Control



- Centered
- equally distributed
- Bell shape
- mean = mode = median
- Normal dist
- Empirical formulae
- Symmetric

$$\mu - 1\sigma \quad \mu + 1\sigma \Rightarrow 68\%$$

$$\mu - 2\sigma \quad \mu + 2\sigma \Rightarrow 95\%$$

$$\mu - 3\sigma \quad \mu + 3\sigma \Rightarrow 99.7\%$$

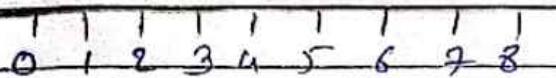
$$\mu = 4$$

$$\sigma = 1$$

$$4 \pm 5$$

$$4.75$$

$$6.95$$



0.256

$$\text{Sample} = \frac{n_1 - \mu}{\sigma / \sqrt{n}}$$

Population

 $\rightarrow z\text{-Score}$ 

How?

$$z = \frac{n_1 - \mu}{\sigma} = \frac{6.95 - 4}{1} = 2.95$$

what?

The  $z$ -score is a statistical method that indicates how many SD a specific data pt is from the mean.

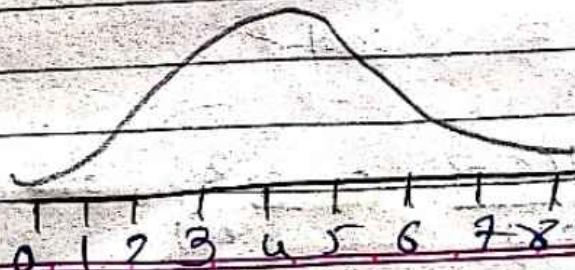
$$a = [1, 2, 3, 4, 5] \quad n = 5$$

$$b = [2, 3, 4, 5, 6] \quad n = 5$$

when?

$z$ -scores are used when comparing individual data pt to the average in normally distributed data set or when converting different types of data to a common scale.

$$\bar{x} = 4 \quad \sigma = 1$$

[ $z$ -S]

#  $\underline{z = \text{Score}}$

$= =$

$z - \text{Score} = \frac{x - \mu}{\sigma}$

$\mu = 0$  and  $\sigma = 1$

### Standard Normal Distribution

Height	Weight	$z = \frac{x - \mu}{\sigma}$	
169	60	3.2	10.24
172	65	6.2	38.44
170	45	-15.8	249.64
168	70	2.2	4.84
170	71	4.2	17.64
829			320.8

$$\mu = \frac{829}{5} = 165.8$$

$$\sigma^2 = \frac{320.8}{5} = 64.16$$

$$\sigma = 8.009$$

$$z_1 = \frac{169 - 165.8}{8.009} = \frac{3.2}{8.009} = 0.3994$$

$\geq 0.4$

$$z_2 = \frac{6.2}{8.009} = 0.7$$

$$z_3 = \frac{-15.8}{8.009} = -1.9$$

$$z_4 = \frac{2.2}{8.009} = 0.27$$

$$z_5 = \frac{4.2}{8.009} = 0.52$$

$$a_1 = -0.01$$

## # Normalization

$$x_n = \frac{x - x_{\min}}{x_{\max} - x_{\min}} \rightarrow \begin{matrix} \text{min max} \\ \text{scaler} \end{matrix}$$

## # Positively skewed distribution

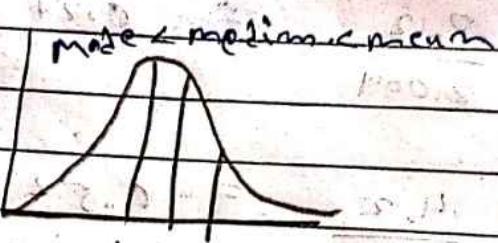
What? :-

In positively skewed distribution, most values are concentrated on the lower end with a long tail extending to the right. A few high values pull the average to the right of the median.

## Skewness :-

A distortion or asymmetry where the distribution deviates from the symmetrical bell curve.

- also called right-skewed for right-skewed distribution



i.e. mean > median > mode

## When? :-

Useful for data with rare but significant high values such as income levels where few individuals earn much more than the rest.

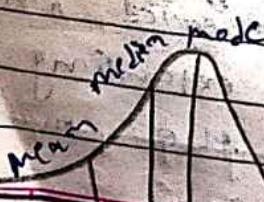
- the mean is greater than the median due to outliers on the higher end.



## Negatively skewed distribution

### What? =

Most values are clustered at higher ends with a few values extending a long tail to the left.



When :-

Used for datasets where values are typically higher, but a few lower values exist  
(returning to age)

- The mean is generally less than median due to leftward tail.
- Also called left skewed or left tailed

Why :-

It helps detect cases where data points are generally higher but occasionally much lower.

# Exponential Distribution  $\Rightarrow$   $X = \text{Constant rate}$

- It describes the time b/w events in a process where events occur independently and at a constant rate.

$$f(x) = \lambda e^{-\lambda x} \quad x \geq 0$$

# Bernoulli distribution (Discrete)

- It models a single experiment with two possible outcomes.

success ( $x=1$ ) , failure ( $x=0$ )

#

## Binomial distribution: Discrete

- Binomial dist extends Bernoulli to n independent trials.

$$P(X=k) \rightarrow \binom{n}{k} p^k (1-p)^{n-k}$$

$X$  = random variable (no. of success out of  $n$ )

$n$  = total no. of trial (100)

$k$  = no. of success ( $0 \leq k \leq n$ )

$p$  = p(success per 1 trial)  $\approx 0.5$

$$\binom{n}{k} = \frac{n!}{(n-k)! \cdot k!}$$

$n=5$        $k=3$

$p=0.5$

$$P(3) = \frac{5!}{3!2!} \times 0.5^{(3)} \times 0.5^{(2)}$$

# Uniform distribution = Continuous

$$[a, b]$$

$$F(n) = \frac{1}{b-a} \quad a \leq n \leq b$$

The probability of any value within the range  $[a, b]$  is same.

# Uniform Distribution (Discrete)

all outcomes are equally likely

$$P(x) = \frac{1}{n} \rightarrow \text{total no.}$$

# Confidence Interval

$\bar{x} = 10 \rightarrow$  point estimate  $\rightarrow$  ll

40-6075k2 pkb $5 - 50 + 5$  $45 - 55 \rightarrow (50k)$  $40 - 50 \quad 60 \quad + 10 - 10$  $\pm 5$  $\pm 10$ 

margin of error

50 → point estimation

 $45 - 55$   
40 - 60

Confidence

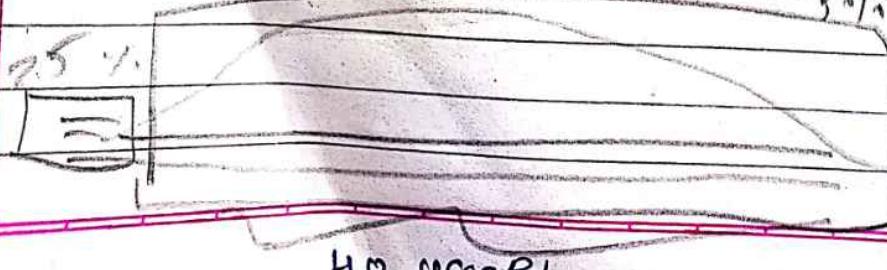
interval

It is a range of values within which we expect a particular population parameter to fall.

Confidence Interval = Point estimate  $\pm$  margin of error

Confidence level

95%



105

25%

H₀ accept

Hypt

## # Hypothesis Testing

Let's → to - loss com

10 1000

- f) Statistical hypothesis test is a method of statistical inference used to decide whether the data at hand is sufficient to support a particular hypothesis.

Hypothesis testing allows us to make probability - size statements about population parameters.

## # Null Hypothesis $H_0$

The null hypothesis assumes that there is no significant relationship or effect b/w two variables. In simpler terms  $\rightarrow$  it says "nothing new is happening".

It serves as a starting point for HT and represents a critique of the assumption of effect until proven otherwise.

- The purpose of  $H_1$  is to gather evidence to reject or fail null hypothesis in favor of alternative hypothesis, which claims there is significant effect or relationship.

# Alternate Hypothesis  $H_1$  or  $H_a$

It is a statement, that contradicts the  $H_0$  and claims there is significant effect or relationship.

# Rejection Region Method

1.  $H_0$  and  $H_1$

5%

2.  $\alpha \rightarrow$  value

10.5

Loss of

Significance

Significance level  $\rightarrow 95\%$

3. assumptions

$N \geq n \geq 30$

4. decide test  $Z$ -test

6-

+ - test

5. value

?

6. Test Conduct

7. Rejection of  $H_0$

$$8. \underline{\underline{z_1 > z_0}}$$

~~1.~~ 2. Avg - 50g  $\log \rightarrow$   
 $\bar{x} - u$   
 $s - u$

10.  $H_0: \bar{x} = 50g$        $H_a: \bar{x} \neq 50g$

$$2. \alpha = 0.05$$

$$3. n \geq 30 \quad z = +0.84$$

4. Z-test

$$S = \frac{u_9 - 50}{s} \times \sqrt{n}$$

$$= \frac{-1}{s} \times \sqrt{n}$$

$$= \frac{-\sqrt{40}}{4}$$

$$Z = 1.058$$

7. Reject / Accept

8. State result

1.  $\mu_0 \rightarrow$  units per day

$$\sigma = 5$$

$\mu_0 = 55$

Training

30 emp  $\rightarrow$  53 units per day

1.  $H_0: \mu = 50 \rightarrow H_0$

$H_a: \mu > 50$

one-tailed

test

2.  $\alpha = 0.05$  15%

3. Data normal, 6 random

$$n = 30$$

4. Z-test

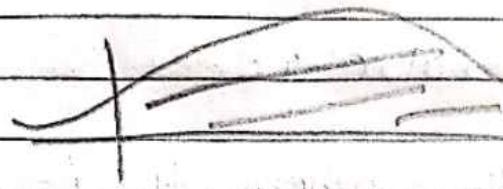
$$5. Z\text{-score} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} - POF$$

$$\frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

$$\frac{53 - 50}{5/\sqrt{30}}$$

$$= \frac{3}{5} = \sqrt{30} - \frac{9}{20} \times 30$$

6.



1.95

1.99

T+6

T-96

-P96  
2017  
0.75

1.96 (60)

2012  
0.025

0.975

97.5

7. Null type Pcep+

Type I error = Type II error

H = 50g

Type I error = Type II error

# Errors

Type I

Type II

reject H<sub>0</sub>

Type I

Correct

accept H<sub>0</sub>

Correct

Type II

Type I

False H<sub>0</sub>

L

No reject → H<sub>0</sub> is correctrejecting H<sub>0</sub> when H<sub>0</sub> is actually correct

Type - I  $\rightarrow$  False - Ic

Accept  $H_0$  when  $H_0$  is actually incorrect.

$P > \frac{0.05}{\alpha} \rightarrow$  Null accept

$P < 0.05 \rightarrow$  Null reject

$P < 0.01 \rightarrow$  Strong evidence

$0.01 \leq P < 0.05 \rightarrow$  Moderate evidence

$0.05 \leq P < 0.1 \rightarrow$  Weak evidence

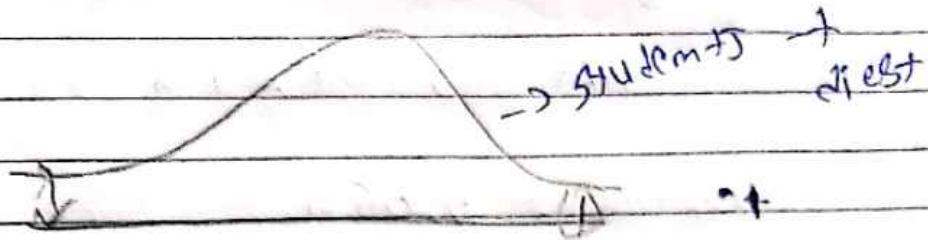
$P \geq 0.1 \rightarrow$  No evidence

$P < \alpha \rightarrow H_0$  Reject

$P > \alpha \rightarrow H_0$  accept

P-value :-

It is a measure of the strength of the evidence against the null hypothesis.

1 - t-test3 types

## 1. one Sample t-test

Compares the mean of a single sample to a known  $\mu$

$$\mu = 50 \quad n = 30 \quad \bar{m} = 49.79 \\ s = 1.2$$

when we get  $\rightarrow z$ -test ( $\text{pop} \rightarrow \text{standard dev}$ )

$\rightarrow \rightarrow t\text{-test } ] (\text{sample - std})$

$$\mu = 50 \quad n = 25 \quad \text{std} = 1.29 \\ \bar{m} = 49.79$$

$$H_0: \mu = 50 \quad H_1: \mu \neq 50$$

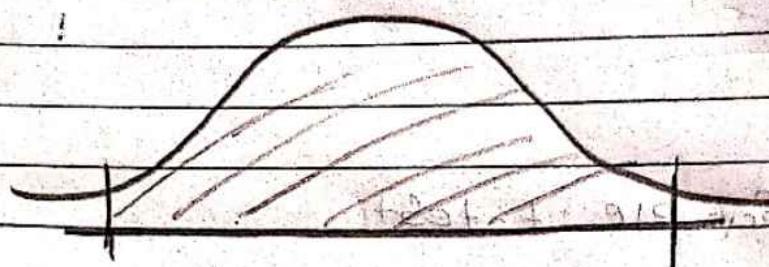
$$t = \frac{\bar{m} - \mu}{s/\sqrt{n}} = \frac{49.79 - 50}{1.2 / \sqrt{25}} = \frac{0.3}{1.2} \times 5$$

$$= \frac{1.5}{1.2} = -1.25$$

$df = \text{degree of freedom}$

$$= n-1 \rightarrow df = 24$$

$$t_{\text{crit/FD}} = 2.064$$



$$-2.064 \quad 0 \quad 2.064$$

$H_0$  accept  $\rightarrow$  sig

2. Independent Two Sample  $t$ -test

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

$$\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \rightarrow \text{standard error}$$

3. paired  $t$ -test

$$t = \frac{\bar{d}}{\frac{s_d}{\sqrt{n}}}$$

D - 2  $\rightarrow$  mean diff

$s_d = S_{d\bar{d}} \rightarrow$  diff

# Chi-Square test

$$= = = = =$$

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

$O \rightarrow$  observed frequency

$E \rightarrow$  expected frequency

$E = \frac{\text{Row} \times \text{Column Total}}{\text{Grand Total}}$

$$df = (r-1) \times (c-1)$$

$r = \text{no. of rows}$

$c = \text{no. of cols}$

Q

	Satisfied	Not Satisfied	Total
High school			
College	50	70	120
Pg	90	60	150

$$\begin{array}{r} 20 \\ \hline 16 \\ \hline 30 \\ \hline 140 \\ \hline 300 \end{array}$$

EF	S	NT	Total
	5		

$$\begin{array}{r} 120 \times 160 = 64 \\ \hline 120 \times 140 = 56 \\ \hline 300 \end{array}$$

$$\begin{array}{r} 150 \times 160 = 80 \\ \hline 150 \times 140 = 70 \\ \hline 300 \end{array}$$

$$\begin{array}{r} 160 \times 30 = 16 \\ \hline 140 \times 30 = 14 \\ \hline 300 \end{array}$$

$$S \quad 64 \quad 56 \quad 120$$

$$C \quad 80 \quad 70 \quad 150$$

$$P.U \quad 16 \quad 14 \quad 30$$

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

$$S.S = \frac{(50-64)^2}{64} = \frac{196}{64} = 3.06$$

$$S.NS = \frac{(70-56)^2}{56} = \frac{196}{56} = 3.5$$

$$C.S = \frac{(90-80)^2}{80} = \frac{100}{80} = 1.25$$

$$C.NS = \frac{(60-70)^2}{70} = \frac{100}{70} = 1.42$$

$$P_{14.05} = \frac{(20-16)^2}{16} = \frac{16}{16} = 1$$

$$P_{4.15} = \frac{(10-14)^2}{14} = \frac{16}{14} = 1.14$$

$$\begin{aligned} \chi^2 &= 3.06 + 3.5 + 1.25 + 1.42 + 1 + 1.14 \\ &= 11.37 \end{aligned}$$

$$df = (3-1) \times (2-1)$$

$$= 2$$

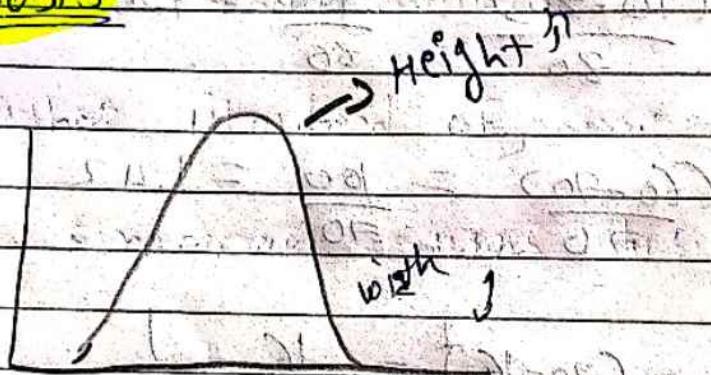
$$\chi^2_{\text{critical}} = 5.991$$

$$\chi^2 = 11.37$$

Rejected  $H_0$

~~RR~~ ~~RR~~

## ~~#~~ Kurtosis



- Kurtosis measures "skewedness" of a distribution or how extreme the outliers are.

## 1. Mesokurtic

- Tails are similar to N(0,1)

- Distribution with kurtosis  $\approx 3$

e.g. Standard Normal Distribution

## 2. Leptokurtic

- A distribution with  $Kurtosis > 3$
- Heavy-tails (with more extreme outliers)

e.g. - t-distribution with very small df  
(distribution)

## 3. Platykurtic

Smooth distributions with  $Kurtosis < 3$

- Light tails (few extreme outliers)

e.g. Uniform dist

Excess kurtosis  $\rightarrow EK = 3$

$EK = 0 \rightarrow$  Leptokurtic

$EK < 0 \rightarrow$  Platykurtic

$$EK = \frac{n}{n-1} \cdot \frac{\sum (x_i - \bar{x})^4}{\left( \sum (x_i - \bar{x})^2 \right)^2}$$

$n = \text{no of observation}$

$x_i = \text{each data point}$

$\bar{x} = \text{mean}$

$\Rightarrow$  High kurtosis (Leptokurtic)

- more extreme outliers
- higher likelihood of such extreme values
- e.g. Financial returns during market crisis / crashes.

# Lower kurtosis (Platykurtic)

- Fewer extreme outliers
- Data is evenly spread

# Kurtosis near 3 (mesokurtic)

- Similar to normal distribution,

*dates  
years*

Left skewed distribution

Normal distribution

Right skewed distribution

Date  
Page

