

Stat

① What is statistics?

→ Collecting, analyzing, interpreting and drawing conclusions.

② What is Data?

③ Type of Stat.

1. Descriptive Stat.

It consists of organizing & summarizing data. (graphs, tables, descriptive)

2. Inferential Stat.

Using Data, we can make conclusions using same techniques.

can

④ Sample and Population

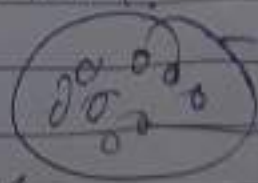
Population - N

→ whole selected is

known as population

Sample - n

→ small subsets of data taken from population



② Sampling technique

① Simple Random Sampling.

↳ Every member of the population has an ~~equal~~ ^{equal} chance of getting selected in sample(n)

② Stratified sampling.

↳ splitting data into non-overlapping groups.
ex: Age, sex, gender

③ Systematic Sampling.

↳ Selects every n^{th} member of the population after a random starting point.

④ Convenience Sampling.

↳ Sample are chosen based on ease of access and availability.

Variables

↳ It is a property that can hold/store / take any value.

Age: $\{8, 10, 15, 20, \dots, 25, \dots\}$
marks: $\{75, 80, 95, \dots\}$

↳ type

1. Qualitative Variable. (Categorical)
value.

→ based on some characteristics can denote categorical values

TG	0 - 10	—	Low
	10 - 50	—	avg
	> 50 - 90	—	Good.

2. Quantitative Variable. (Numerical value)

→ measurable numerically.

→ Height - $\{162, 157, 180, \dots\}$
weight - $\{57, 55, 56, \dots\}$

i) Discrete - Value are countable and finite
(whole numbers) (eg. no. of shoes)

ii) Continuous - Value can take any
value within range.
(Decimal no.) (eg. height)

Variable measurement scales.

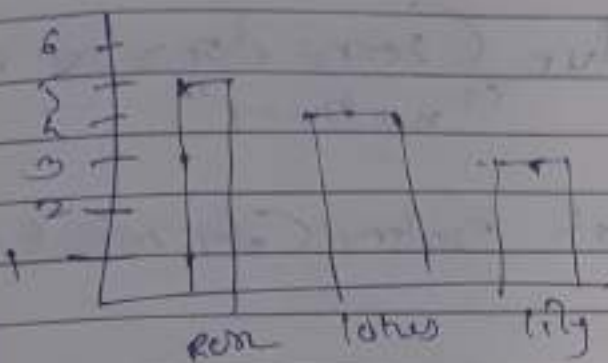
- ① ordinal scale. - (order)
ex - rank, direction.
- ② Nominal scale - (categorical values) (no order)
→ colors, classes, degree.
- ③ Interval Scale : (no zero / absolute point)
ex - temperature.
- ④ Ratio scale. (zero means nothing)
ex - Height, weight, age
BP, Income.

Frequency

data : Flowers

(Rose, lily, lotus, Rose, rose,
rose, lotus, lily, lotus, lily,
lotus)

Flowers	frequency	CF cumulative
Rose	5	5
lily	3	8
lotus	4	12
	<u>12</u>	



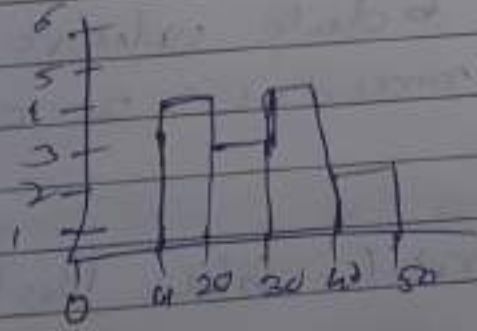
Bar chart

Histogram

num = [12, 15, 13, 15, 21, 21, 27, 35, 34, 36, 39, 41]

bin

10-20	4
20-30	3
30-40	2
40-50	2



Interval size \rightarrow It measures the size of the value but does not have any zero points

→ no absolute value (zero does not mean "nothing")

→ can't calculate other (e.g., 2012) given

ex: temp

0°C does not mean "no temperature"

20°C is not twice as hot as 10°C

relative.

diff. b/w 2000 & 2020 is 20 years but
2020 is not twice as old.

* Ratio scale.

It measures the scale where there is a zero point meaning a zero represents complete absence.

→ Equal interval but value.

→ zero allows ratios

→ can perform all math

measure of Central tendency

Avg \rightarrow mean.

pop.

Sam.

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

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

mean \rightarrow It refers to the measure used to determine the centre of the distribution of the data.

$\{1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 10\}$

\rightarrow arithmetic

$$\frac{32}{10} = 3.2$$

$$\frac{132}{11} = 12$$

median \rightarrow middle value

even

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

ex - $\{11, 12, 13, 14, 15, 16\}$

$$n = 6$$

$$\therefore \frac{\left(\frac{6}{2}\right)^n + \left(\frac{6}{2}\right)^n}{2} = 11$$

$$\therefore \frac{3^n + (3^n)^4}{2} = 13.5$$

add \therefore $\left(\frac{n+1}{2}\right)$

{ 11, 12, 13, 14, 15 }

$$\left(\frac{5+1}{2}\right)^n = \frac{6^n}{2} \times 3 = 3$$

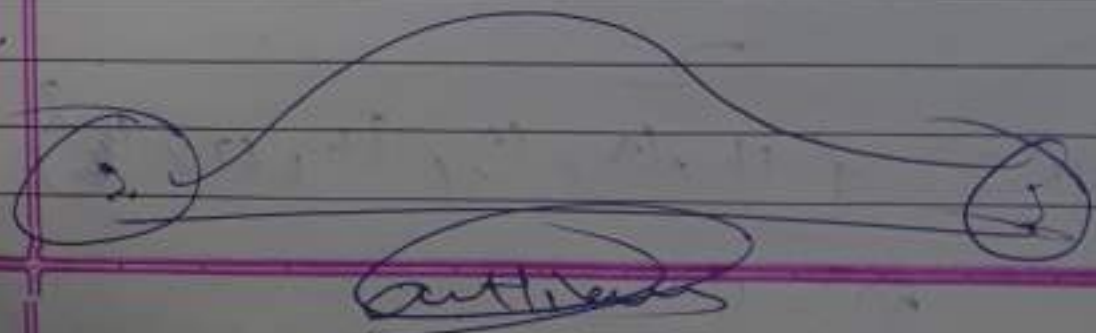
sa actions

{ 21, 23, 25, 27, 32, 100 }

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

unless

→ a player who doesn't follow pattern or trend or doesn't follow is consider as ~~an~~ action

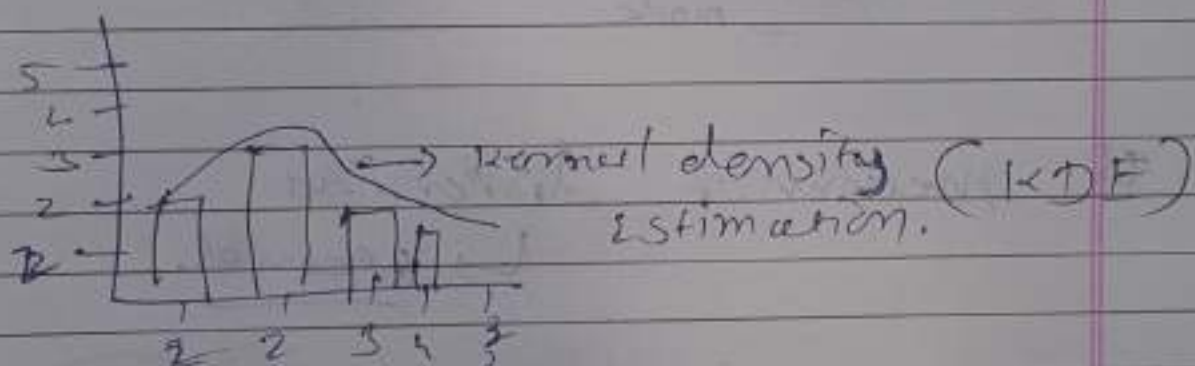


mode

→ most frequent value (repeated)

→ 1, 2, 2, 3, 4, 2, 3, 2, 0

mode = 2



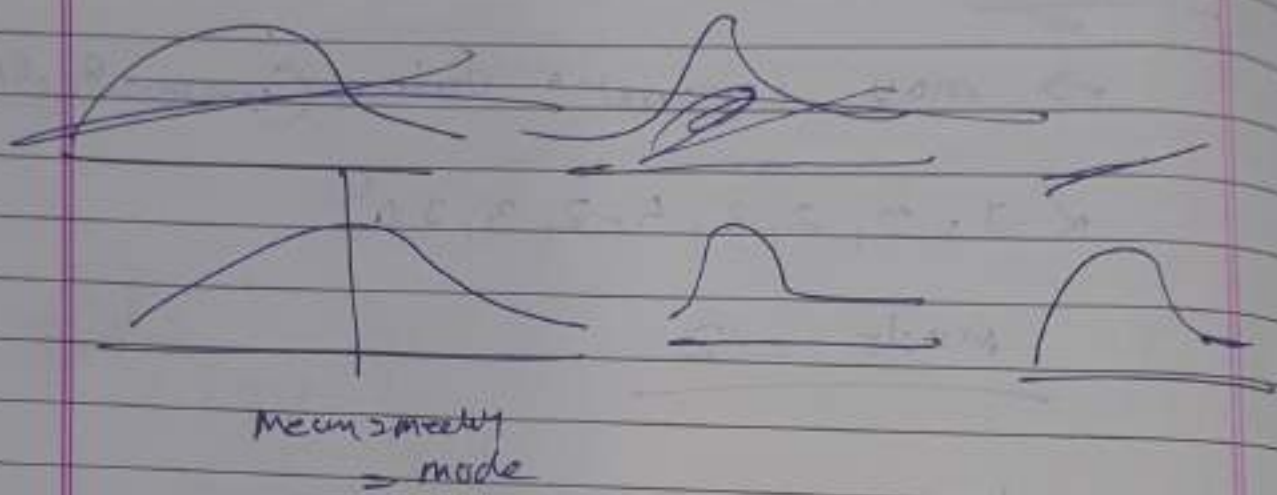
→ for categorized many data

0-5 cr. mode

← New category →
or
→
consider
unlabeled
random

→ for numerical value

Gaussian distribution
Skewed
↓
median.



Measure of Dispersion
↳ spread.

$$[5, 1, 1, 2] \rightarrow \frac{5+1}{3} = 2$$

$$\frac{[2, 2, 2, 2, 2]}{5} = 2$$

→ Variance (σ)

It measures how far the numbers
in a dataset are from the mean.

(how each value differs from a
dataset mean)

High Variance \rightarrow more spread

low " \rightarrow closer between

Pop.

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

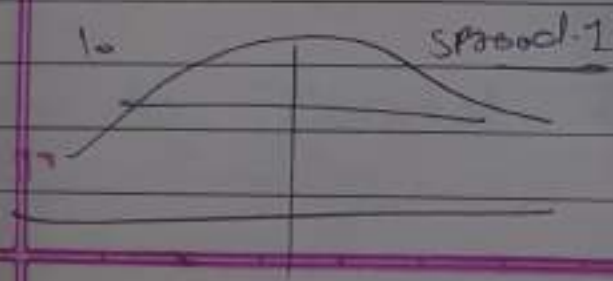
sum

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

i	$x_i - \mu$	$(x_i - \mu)^2$
1	-1.83	3.34
2	-0.83	0.68
3	-0.83	0.68
4	0.17	0.02
5	0.17	1.36
6	2.17	4.70
Sum		10.78

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

$$\sigma = 1.33$$



*

Standard Deviation

- Just a ~~sq~~ squared, mean of Variance
- It gives you a measure of spread that is in the same unit as the original data, making it easier to interpret than Variance

[Same Unit, easily compared]

Pop.

$$\sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{n}}$$

Sample

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

*

(Key point (Variance & SD))

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

* Variance Formula

* Pop.
$$\sum \frac{(x - \mu)^2}{N}$$
 $(N = \text{total no. of data pts})$

When we have data from the entire population, we use 'N' in the denominator. This gives us an exact measure how the data pts vary around the pop. mean (μ).

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

When we are working with a sample, we only have sample mean \bar{x} , which is an estimate of the pop. mean μ .

→ Using sample mean in calculating tend to make the variance slightly smaller than the true pop. variance.

→ To correct this bias (underestimating variability)

→ To correct this bias (under
we divide by 'n-1' instead of 'n'.
This makes variance unbiased
measure of spread. accurate.

→ By Subtracting from n, we
adjust for the fact that the
is not perfectly.

→ Representative of μ this correction
ensures that the sample variance
is an unbiased estimator of the
pop. variance.

* Percentage

1, 2, 3, 4, 5

∴ of the no. 3 of them are odd

$$3/5 = 0.6 \rightarrow \times 100$$

= 60%

* Percentile

→ A Percentile is a value below
which a certain p.c. observations lie.

dataset = [2, 3, 3, 4, 6, 5, 5, 6, 7, 8, 8, 7, 9, 10, 11, 11, 12]

$$P.R = \frac{\text{no of Value below } x}{n} = \frac{16}{20}$$

$$= \frac{16}{20} = 0.80$$

8. what Value exist at percentile ranking of '25'?

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

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

$$= 5 + 1 = 6^{th}$$

75 Percentile Value

$$= \frac{75 \times 20}{100} + 1$$

$$= 15 + 1 = 16^{th}$$

* Five number Summary

1. minimum
2. first quartile (Q_1)
3. median
4. third quartile (Q_3)
5. maximum

[1, 7, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 7, 8, 8, 8, 9, 9, 10]
123

lower fence = $Q_1 - 1.5(IQR)$
Higher " = $Q_3 - 1.5(IQR)$

$Q_1 = 25$ Percentile (25th)

$Q_3 = 75$ " (75th)

$IQR = \text{Inter quartile range.}$

$$= Q - Q_1$$

$$Q_1 = \left(\frac{25}{100} \times 20 \right) \text{ \pounds}$$

$$Q_1 = 3$$

$$Q_3 = \left(\frac{75}{100} \times 20 \right) \text{ \pounds}$$

$$= 15$$

$$= 8$$

$$IQR = 8 - 3 = 5$$

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

$$= 3 - 1.5(5)$$

$$= -4.5$$

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

$$= 8 + 1.5(5)$$

$$= 15.5$$

$$= 15.5$$

$$Q_1 = \left(\frac{25}{100} \times 20 \right) \quad \text{H}$$

$$Q_1 = 3$$

$$Q_3 = \left(\frac{75}{100} \times 20 \right) \quad \text{H}$$

$$= 15$$

$$= 8$$

$$IQR = 8 - 3 = 5$$

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

$$= 3 - 1.5(5)$$

$$= -4.5$$

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

$$= 8 + 1.5(5)$$

$$= 8 + 7.5$$

$$= 15.5$$

remaining debt

✓ 1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100

S. No Survey

$$m \cdot m = 1$$

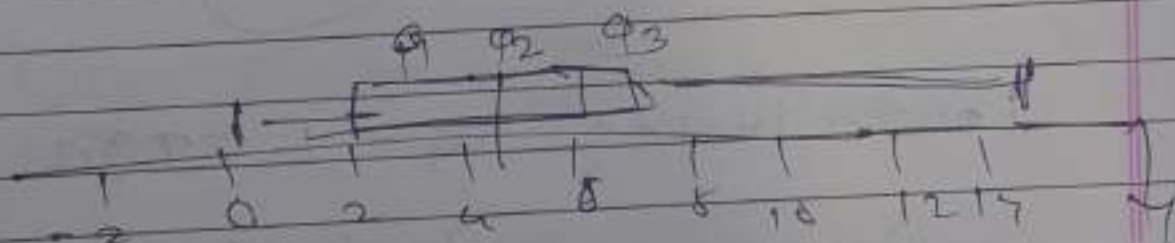
4173

median = 5

28

mean > 15

Box plot



20/1/20

Data Distribution

- It refers to a way in which values or data points are spread or arranged.
- It shows how often different values occur in data set & describes the overall pattern of the data.

→ Center.

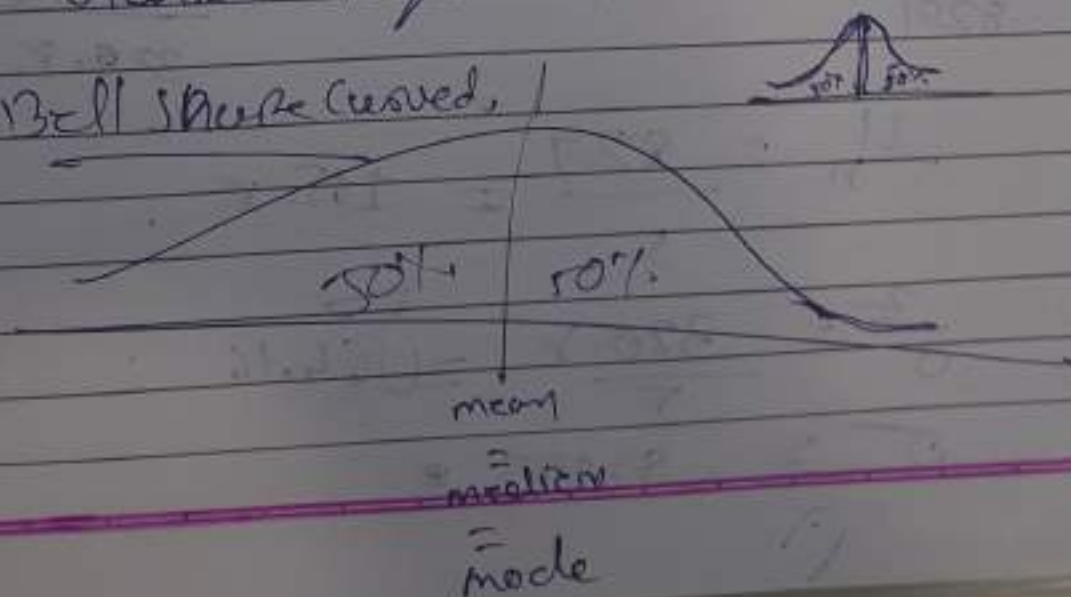
→ Spread

→ Shape.

→ Outliers.

Gaussian / Normal distribution (Bell shape)

Bell Shape Curve,



$$\mu \pm 2\sigma = 88\%$$

$$\mu \pm 2\sigma = 95\%$$

$$\mu \pm 3\sigma = 99.7\%$$

Empirical

rule

Z-Score

$$Z\text{-Score} = \frac{x_i - \mu}{\sigma}$$

Example

Height	weight	$x_i - \mu_n$	
169	60	3.2	10.24
172	65	5.2	27.04
150	45	-15.8	249.64
168	70	2.2	4.84
170	77	4.2	17.64
829			320.8

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

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

$$\sigma = 8.007$$



$$21 \rightarrow \frac{167 - 165.8}{8.007} = \frac{3.2}{8.007} = 0.4$$

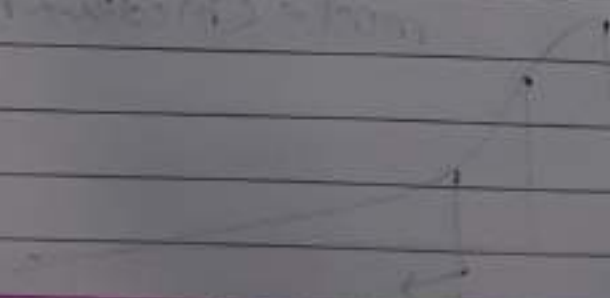
$$22 \rightarrow \frac{6.2}{8.007} = 0.7$$

$$23 \rightarrow \frac{15.2}{8.007} = 1.9$$

$$24 \rightarrow \frac{2.2}{8.007} = 0.27$$

$$25 \rightarrow \frac{4.2}{8.007} = 0.52$$

✓ ⊙



Normalization

$$x = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

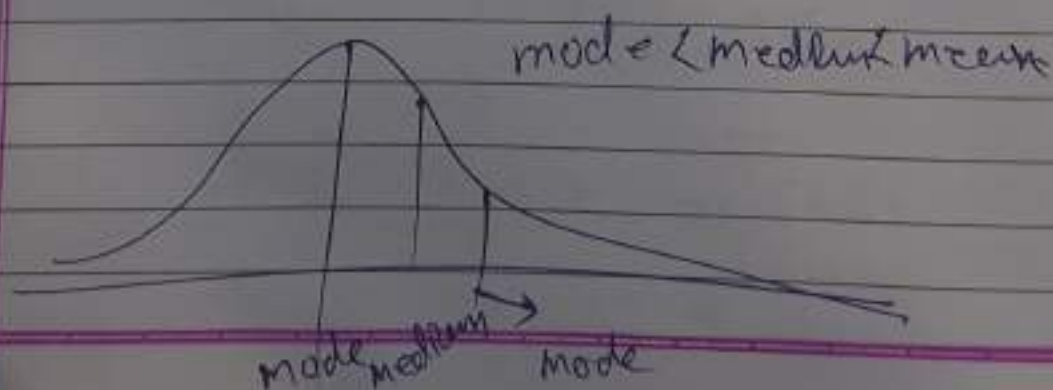
positively skewed distribution

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: \rightarrow A distortion or asymmetry that deviates from the symmetrical bell curve

\rightarrow also called right-skewed or right-tailed distribution.



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,
~~when?~~

→ Why?

Because it makes comparing with data sets easy.

ex 2022

$$D1 \rightarrow \mu = 250 \quad \sigma = 283.8 \approx 10$$

$$D2 \rightarrow \mu = 280 \quad \sigma = 285 \approx 12$$

$$z_1 = \frac{285 - 255}{10} = 3$$

$$z_2 = \frac{263 - 260}{12} = \frac{3}{12} = 0.25$$

$$z_3 = \frac{263 - 255}{10} = \frac{8}{10} = 0.8$$

$$z_4 = \frac{255 - 250}{12} = \frac{5}{12} = 0.41$$

Exponential Distribution. ($\lambda =$ number of events per unit time) (constant rate)

→ It describes the time b/w events in a process where events occur independently 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 distribution extends Bernoulli to n independent trials.

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

$x =$ random variable (no. of successes out of n trials)
 $n =$ total number of trials / $p =$ probability of success (0.5)
 $k =$ number of success ($0 \leq k \leq n$)

$$p = p \text{ (success, my friend)}$$

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

④ $n=5$ $k=3$

$$p(3) = \frac{5!}{2! \cdot 3!} \times 0.3^{(3)} \times 0.7^{(2)}$$

Uniform Distribution

⇒ (Continuous)

$[a, b]$

$$f(x) = \frac{1}{b-a}$$

$$a \leq x \leq b$$

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

→ (discrete) .

'all outcomes are equally likely,

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

Confidence Intervals (CI)

$\bar{x} = 50 \rightarrow \text{point estimate} - \mu$

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

Confidence Interval = point estimate \pm measurement error

→ Confidence Level

→ Sampling Statistical Parameters



→ Hypothesis Testing

A statistical hypothesis test is a method

→ Null Hypothesis (H_0)

→ Alternative (H_A)

25/12

Page
Date

Rejection Region method

1. H_0 & H_a

2. $\alpha \Rightarrow$ Value \rightarrow LOS (loss of significance) $\rightarrow 0.05$
 \rightarrow Significance level $\rightarrow 95\%$

3. Assumptions.

4. decide test \rightarrow z-test, t-test

5. Value.

6. test conclusion

7. Reject / Accept

8. ~~state the rejected~~ state Results

* probability \rightarrow possibility chance.

\rightarrow It is a measure of the likelihood of an event.

eg - Roll a die $\{1, 2, 3, 4, 5, 6\} \rightarrow$ sample space

no. of ways an event can occur
total no. of possible outs.

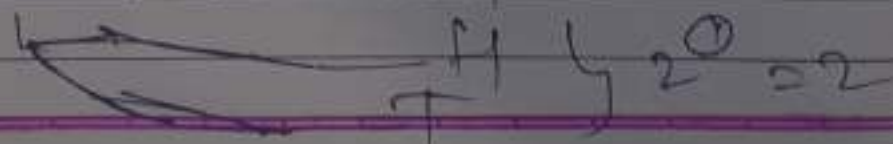
eg even no.:

$P(E) = \frac{\{2, 4, 6\}}{\{1, 2, 3, 4, 5, 6\}} \rightarrow$ favourable.

$$= \frac{3}{6} = \frac{1}{2} = 0.5$$

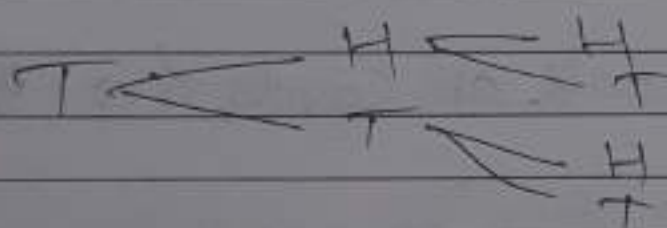
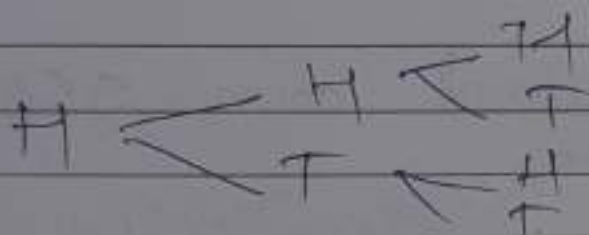
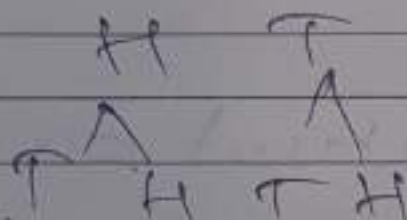
* two events are occurring together

\rightarrow tossing one coin



toss 2 coin

$$2^2 = 4$$



dice roll $\rightarrow 6 \rightarrow 6$

$\{(1,1), (1,2)\}$

* Mutual Exclusive event

Two events are called mutual exclusive "if" they can not occur to occur at the same time.

Roll a dice $\{5, 6\} \rightarrow X$

Roll a coin $\{H, T\} \rightarrow X$

* non-mutually exclusive

multiple events can occur at the same time.

eg - Draw of cards $\{Q, K\}$ in

* Addition rule for probability

Q I toss a coin, what is the probability that the coin comes either on head or tail?

mutually exclusive

$$P(A \cup B) = P(A) + P(B)$$

$$= \frac{1}{2} + \frac{1}{2}$$

$$P(A \cup B) = 1$$

Dice $\rightarrow 1, 3, 6$

$$P(1 \text{ or } 3 \text{ or } 6) = \frac{1}{6} + \frac{1}{6} + \frac{1}{6}$$

$$= \frac{3}{6}$$

$$= \frac{1}{2}$$

$$= 0.50$$

③ Multiplication Rule.

* Independent Event

Dice $\{1, 2, 3, \dots\}$

2 dice \rightarrow ind.

$$\left[\begin{array}{c} \text{1} \\ \text{2} \end{array} \right]$$

$$\frac{1}{6}$$

$$\left[\begin{array}{c} \text{1} \\ \text{2} \end{array} \right]$$

$$\frac{1}{6}$$

* Dependent event



$$P(G) = \frac{3}{5}$$

$$P(P) = \frac{2}{4}$$

* What is the prob of rolling a "5" on a throw.

Roll a dice $\{5, 6\} \rightarrow X$

Roll a coin $\{H, T\} \rightarrow X$

* don't - mutually exclusive.

multiple event can occur at the same time.

eg - Deck of cards $\{Q, K\}$

* Addition rule for probability

Q I toss a coin, what is the probability that the coin comes either on head or tail?

mutually exclusive

$$P(A \cup B) = P(A) + P(B)$$

$$= \frac{1}{2} + \frac{1}{2}$$

$$P(A \cup B)$$

$$= 1$$

1 Die $\rightarrow 1, 3, 6$

$$P(1 \text{ or } 3 \text{ or } 6) = \frac{1}{6} + \frac{1}{6} + \frac{1}{6}$$

$$= \frac{3}{6}$$

$$= \frac{1}{2}$$

0.50

Multiplication Rule,

* Independent Event

Die $\{1, 2, 3, \dots\}$

2 dice \rightarrow ind.

(1)

$$\frac{1}{6}$$

(2)

$$\frac{1}{6}$$

* Dependent event



$$P(G) = \frac{3}{5}$$

$$P(P) = \frac{2}{4}$$

What is the prob of getting a "5"

Hypothesis testing

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

Null Hypothesis H_0 .

The null hypothesis assumes that there is no significant relationship or effect b/w two variables.

It serves as a starting point for HT & ~~test~~ ^{assumptions} 'statistical' or the assumption of no effect until proven otherwise.

The purpose of HT is to gather evidence to reject or fail null hypothesis in favour of alternatives.

hypothesis, which claims there
is significant effect or relation.

Alternate Hypothesis

It is a statement, that contradicts the H_0 & claims
there is significant effect or
relation.

Rejection Region Method

1. H_0 & H_a

2. $\alpha \rightarrow$ value. \leftrightarrow LOS

\hookrightarrow significance level $\rightarrow 0.05$

3. ~~ass~~ assumptions.

4. decide test \rightarrow z-test
t-test

5. Value. \downarrow

6. test conduct

7. Reject / Accept / \hookrightarrow final
Result

2.1

50 \rightarrow Unit per day

$$\sigma = 5$$

30 ems \rightarrow 53 unit per day

2. $\mu = 50 \rightarrow H_0: \mu = 50 \quad H_a: \mu > 50$

2. $\alpha = 0.05, 15\%$ on tail test

3. Data - Normal, σ , random

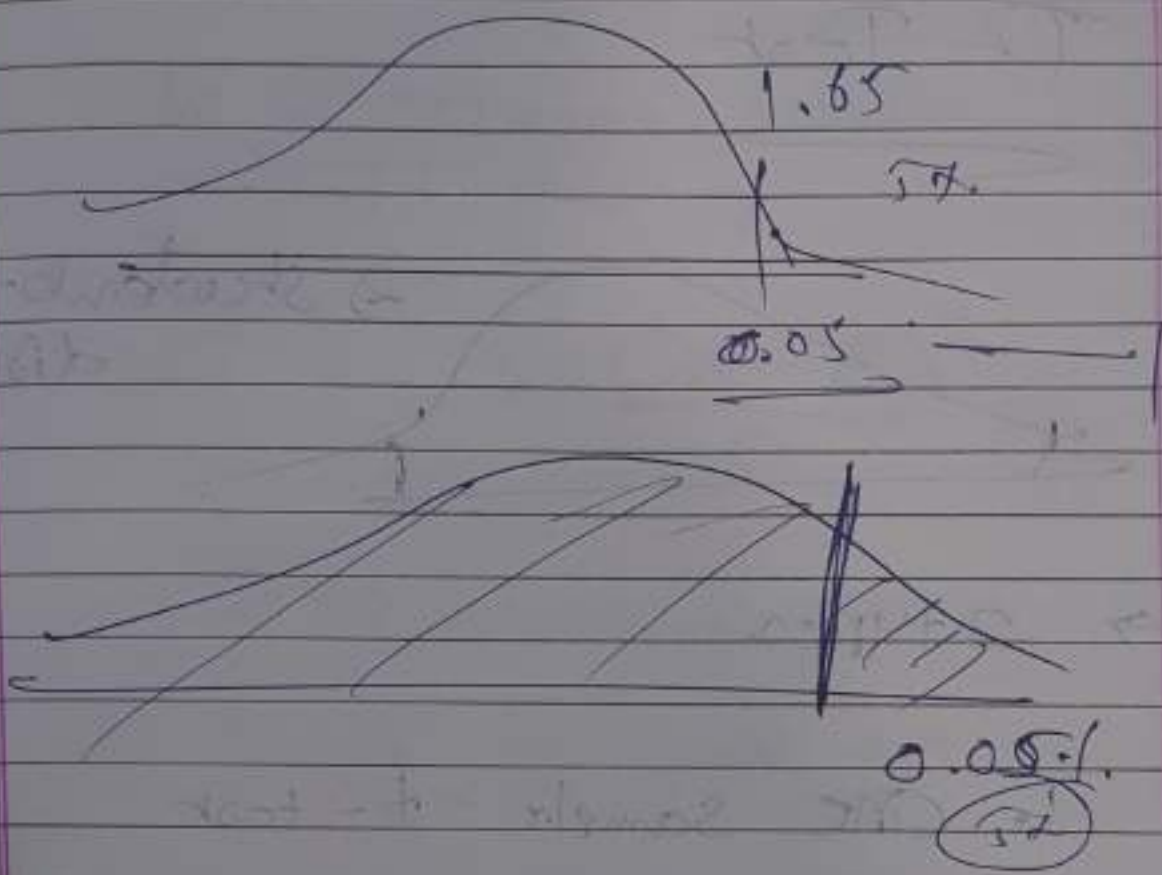
$$n = 30$$

4. Z-test

5. Z-score $\rightarrow \frac{n_1 - \mu}{\sigma}$ - pop.

$$\frac{n_1 - \mu}{\sigma/\sqrt{n}} = \frac{53 - 50}{5/\sqrt{30}}$$

$$2 \approx 3.28$$



Reject H_0

$$\mu > 50$$

T-Test



→ 3 types

1. One Sample - t - test

Compares the mean of sample
to a known μ

2. Independent two sample t - test

or paired t - test

Chi - Square test

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

O \rightarrow Observed freq.

E \rightarrow Expected freq.

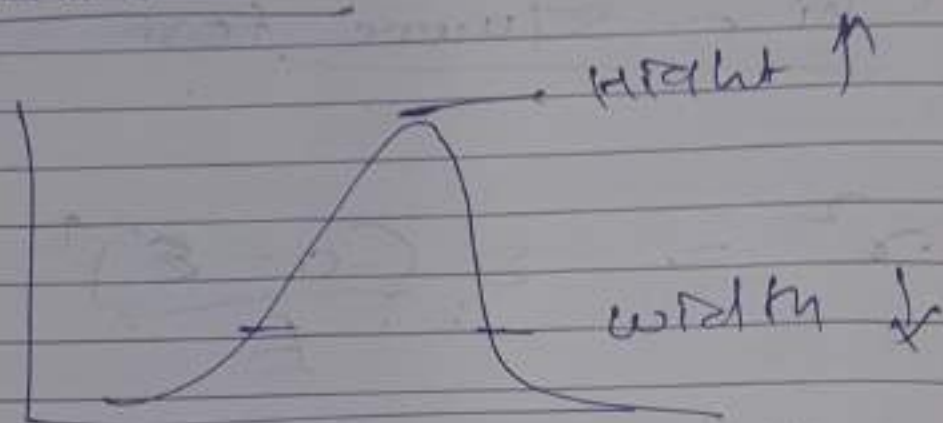
$$E = \frac{\text{Row total} \times \text{coll. total}}{\text{Grand total}}$$

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

r = no. of rows

c = no. of coll.

#

kurtosis

→ kurtosis measures "tailedness" of a distribution or how extreme the outliers are.

→ 1. mesokurtic

→ tails are similar to $N(0,1)$

→ distribution with kurtosis ≈ 3 .

→ ex. standard normal distribution

→ 2. leptokurtic

→ A distribution with kurtosis > 3

→ heavy tails (with more extreme outliers)

Platykurtic

- distribution with kurtosis < 3
- lighter tailed.
- eg - Uniform dist.

Excess kurtosis $\rightarrow K - 3$

$K = \text{kurtosis}$

$EK > 0 \rightarrow \text{leptokurtic}$

$EK < 0 \rightarrow \text{Platy}$

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

$n = \text{no. of observations}$

$x_i = \text{each data points}$

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

→ High kurtosis (Leptokurtic)

- More extreme outliers.
- Higher likelihood of rare, extreme values.
- Higher ρ .
- Financial return during market crisis / recessions.

• Lower kurtosis (Platykurtic)

- Fewer extreme outliers.
- Data is evenly spread.

• Kurtosis near 3 (mesokurtic)

- Similar to normal distribution.

