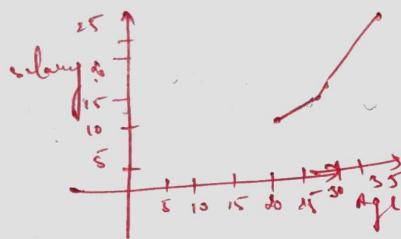


Statistics

- * Statistical analysis is meant to collect & study info ↓ available in large quantities.

→ Branch of maths (where computation is done on bulk of data) using charts & graphs.

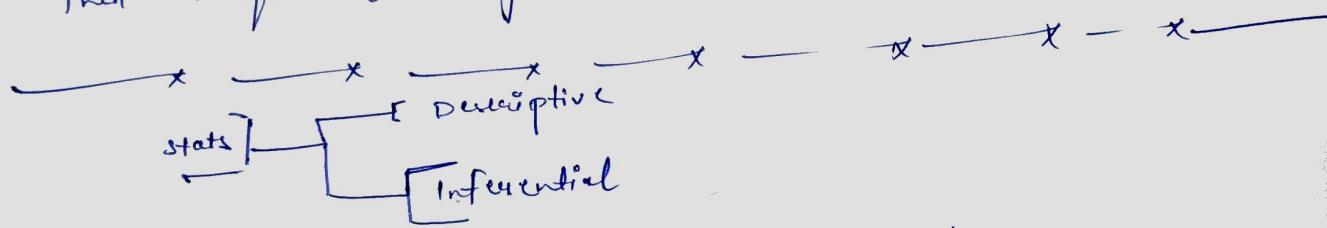
Age	salary
25	10000
29	12000
31	16000
35	25000



→ trend (\uparrow in salary w/ time)

- * Measurements: Data collected for analysis. | Sampling -

↓ Then analysis is done for measurement.



DESCRIPTIVE STATS:- • describe characteristics of data.

(What is data telling you?)

consist 3 measurements

→ central tendency (mean, median, mode)

→ measure of variability (spread) { std, variance }

→ frequency distribution. (cont.)

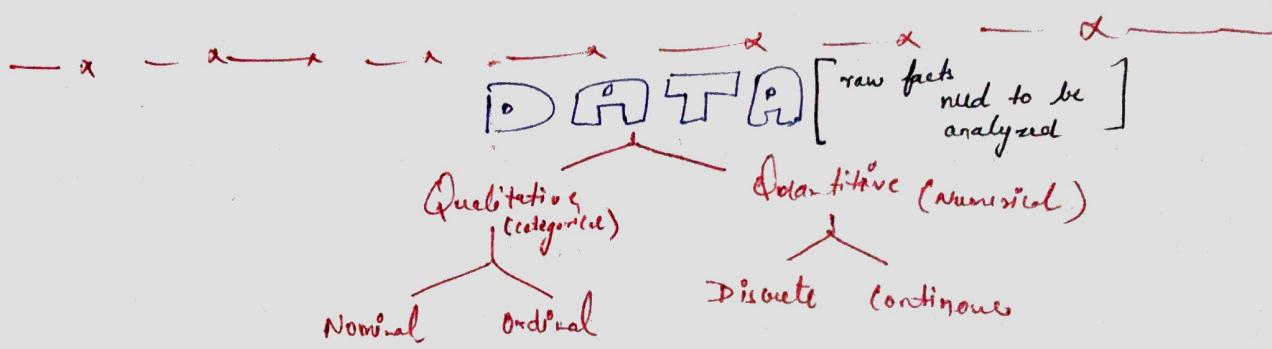
CV → center of data,
MV → deviation, dispersion of
data set
FD → Desirable - occurrence
of data within stand.

$$SD(\sigma) = \sqrt{\frac{\sum (x_i - \mu)^2}{n}} \quad \{ \mu \rightarrow \text{mean} \}$$

INFERENTIAL STATS:- Helps us to make conclusions / predictions about a big group using data from its small sub-group (sample).

2 MAIN AREAS

- Estimating parameters
- Hypothesis Testing.

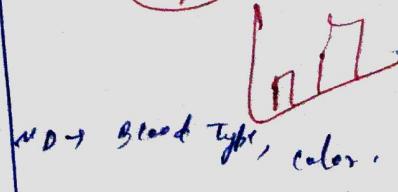


Qualitative/Categorical Data:-
• describes the data that fits into categories.

- cate-info includes cate-variables which describe features or person's gender, hometown, etc.
- Not numerical
- sometimes there have numeric values but has no sense.
e.g. Dob: 17/1/05
Postal code

ex e.g. - Gender
Color
City

- NOMINAL DATA. (Name only)
- info which helps to label the vars w/o numeric values.
- also called nominal scale - can't be measured / scaled / ordered.
- Nom data examined using the grouping method.
 Data grouped into categories
 ↳ frequency / % can be calculated.
- represented by pie charts.



ORDINAL DATA - have order.

e.g. → economic status, Grading (A, B), ratings, This size represented by:- bar chart.

Quantitative Data

- numerical data

→ How much, how many.
→ n , N , w



discrete [finite num. of pos values] can't be subdivided meaningfully. (whole numbers)
e.g. no. of workers.

continuous - data that can be calculated.
has infinite no. of probable values.
can be selected within a range.
ex → Temperature Range.

SAMPLING

Population

The universe of objects that is required to be analyzed.

SAMPLE

subset of population

size of sample & size of population

* Why Sampling Important:-

- > gathering data from population → not possible,
Sampling is applicable in such cases.
- > sampling makes info faster.
- > measuring every d.p.t is not cost effective.
- > can easily analyze the data using samples.
- > smaller data set = smaller data collection cost.

Sampling

Probability Sampling

Non-Probability Sampling

Probability Sampling :-

Every member of population has an equal chance of being selected.

↳ Simple random sampling

→ Stratified Random S.

→ Cluster Sampling

→ Systematic Sampling

Pros

- unbiased & representative
- results are statistically valid

Cons

- Time consuming & costly
- Needs a complete population list.

We use when:-

- to make statistical inferences
- To avoid biases.
- Accuracy & reliability is important.

No. 2 - Probability Sampling

• Not everyone has a chance of being selected.

→ samples are selected on basis of judgement or convenience of accessing data.

We use when,

- money, time, access → limited

- when we need quick insights

Eg

→ Convenience Sampling

→ Purposive Sampling

→ voluntary response sampling

→ snowball sampling

Pros

- fast & cheap
- easy to execute
- useful when less population data available

Cons

- high chance of bias
- results can't be generalized

Cluster Random Sampling :- (PS)

- > Divide population into groups.
- > then randomly select the group from all groups.

Sampling Techniques

Probability Sampling

Random systematic stratified cluster

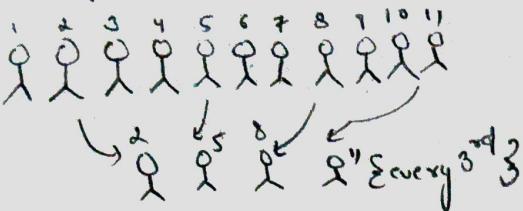
Non-P-Sampling

Convenience Purposive Voluntary Response Snowball

Prob. Sampling

1.) Random Sampling: randomly choose a member from the population.
 > every member or set of members has an equal chance of being selected.

2.) Systematic Sampling:- > Similar to random sampling.
 Pick every k^{th} person. (from a ordered list)



3.) Stratified Sampling:-
 Homogeneous
 > divide population into groups
 > then from each group, select members randomly
 stratum → similar groups.

e.g. students from
each year
(1st, 2nd, 3rd & 4th)

4.) Cluster Random Sampling:-
 Hetero.
 e.g. Diff schools.
 Divide population into groups
 > then randomly select the subgroups from all the groups.

non-p-Sampling

1.) Convenience Sampling: include members who are easy to reach for researcher.

2.) Purposive Sampling: > select a sample based on purpose of the research.
 > researcher select sample using their expertise & knowledge.

3.) Voluntary Response Sampling:
 > based on ease of access
 > members volunteers by themselves instead researcher selecting the participants.

4.) Snowball Sampling:
 Existing participants refer or recruit new participants, forming a chain.

Population / Sample	Time	Estimation	Formula
	Time	N	num. of items
add up to mean		\bar{u}	$\sum_{i=1}^N x_i / N$
	Variance	s^2	$\frac{\sum_{i=1}^N (x_i - \bar{u})^2}{N}$
$\sum x_i$ (x_1, x_2, \dots, x_n)	S. size	N	num. of items in sample
	S. mean	\bar{x}	$\sum_{i=1}^N x_i / n$
	S. variance	s^2	$\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$

$x_1 = 5, 10, 15, 20, 30$	Bessel's correction ($n-1$ denominator in S. var formula)
$N = 5$	
$\bar{u} = \frac{50}{5}, \bar{u} = 10$	
$s^2 = \frac{(5-10-1+15+20)^2}{5}$	
$s^2 = (25+25+25)/5$	
$\bar{s} = 12+36+1+16+196/5$	
$\bar{s} = 370/5$	
$s^2 = 74$ (variance)	
$s/\sqrt{t} = \sqrt{74}$	

Des. Stats \rightarrow Measures of Central Tendency.. (Measures the central value of data set.)

> Gives the idea about the concentration of the value in the central part of the distribution.

• Mean

• Median

• Mode

MEAN

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

(sensitive to outliers)

MEDIAN

→ The middle number

→ order all d. plot pick from middle.

→ if there are odd no. in middle,
take mean of those two.

n is odd

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

n is even

$$\frac{\left(\frac{n}{2}\right) + \left(\frac{n}{2} + 1\right)}{2} \left(\frac{m_1 + m_2}{2} \right)$$

MODE The value which occurs more frequently.

- can have more than one mode

- used in categorical dataset.

Measure of Dispersion ..

→ Tells us where the center of data lies. describes, how much data values vary or how they are spread from the avg.

1. Range ($\text{Max} - \text{Min}$) → used to construct control chart is quality assurance helpful when to focus on extreme values of data set.

2. Variance (σ^2) $\sum \frac{(x-\mu)^2}{n}$ } $\mu \rightarrow \text{Mean}$ } measures the dispersion of data around mean
(σ^2) v close to mean \rightarrow low variance [slightly dispersed]
 v far to mean \rightarrow high variance [highly dispersed]

3. Inter Quartile Range : IQR ($Q_3 - Q_1$) → measures the middle 50% of data.
(Q_1 and Q_3).
 Q_1 \downarrow Q_3 \downarrow $75^{\text{th}}/\text{ile}$ \downarrow $25^{\text{th}}/\text{ile}$ → ignores extreme smaller or large values (outliers).
Helpful to detect outliers present in the dataset.

4. Std = $\sqrt{\sigma^2}$
 $\sigma = \sqrt{\frac{\sum (x-\mu)^2}{n}}$

Indicates how far the data points is dispersed from mean. (σ)

5. Mean Deviation → how data is deviated from mean value

$$MD = \frac{\sum |x - \mu|}{n}$$

PROBABILITY

always b/w 0/1

Prob. is measure of how likely is something is to happen.

Trial:- Each performance in random trial.

Event:- outcome of trial $P(E)$

Sample Space:- Set of all possible outcomes in a random experiment $P(S)$.

Random Experiment:- To perform more than once (outcome can't be predicted)

Add" Rule of Probability:- Helps you find the prob. that atleast one of two (or more) events happens.

used for 'OR' questions

CASE - A

Events are mutually exclusive

→ means they cannot occur at the same time.

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

CASE - B

Events are not mutually exclusive

→ means they can occur at the same time.

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

Independent Events :- Two events are independent when:

→ the occurrence of one does not affect the occurrence of other.

Independence cond's for events A & B.

$$\text{i)} P\left(\frac{A}{B}\right) = P(A)$$

$$\text{ii)} P\left(\frac{B}{A}\right) = P(B)$$

$$\text{iii)} P(A \& B) = P(A) * P(B)$$

Cumulative Probability: - means the total prob. of all outcomes upto a certain point.

> It tells us the chance that a random variable is less than or equal to a certain value.

$P(a \leq X \leq b)$ a range is given

Conditional Probability: - means the chance of something happening, given that something else has already happened.

$$P(B|A) = P(A \cap B) / P(A)$$

rule of mult $\rightarrow P(A \cap B) = P(A|B) \times P(B)$ or $P(A \cap B) = P(B|A) \times P(A)$

BAYES THEOREM extension of conditional probability

(cond) probability help is: $P(A|B)$

Bayes theorem says, if we know $P(A|B)$, we can determine $P(B|A)$, given that $P(A)$ & $P(B)$ are known to us.

formula: $P(A|B) = \frac{P(A) \times P(B|A)}{P(B)}$

Intuition Behind Conditional Prob. formula

$P(A|B) = \frac{P(A \cap B)}{P(B)}$, is based on the concept of reducing our sample space.

B is occurred & B is our new sample space, we don't consider sample space where B is not occurring.

$P(A \cap B) \rightarrow A$ occurred & B occurred.

BAYES THEOREM

prior

$$P(A|B) = \frac{P(A) \times P(B|A)}{P(B)}$$

likelihood
marginal

BT \rightarrow fundamental concept in field of prob. & statistics. that describes:

how to update the probabilities of hypothesis when given evidence.

. Provides a way to revise existing predictions on theories (update prob.) given new evidence.

Probability Distribution:-

Tell us how the values of data set are spread.
and how likely each value is to occur.

Describes the pattern of data & probabilities of all possible outcomes.

Types:-

1) uniform distribution.

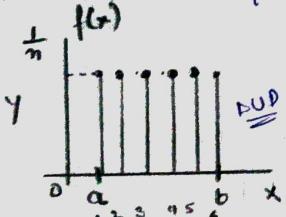
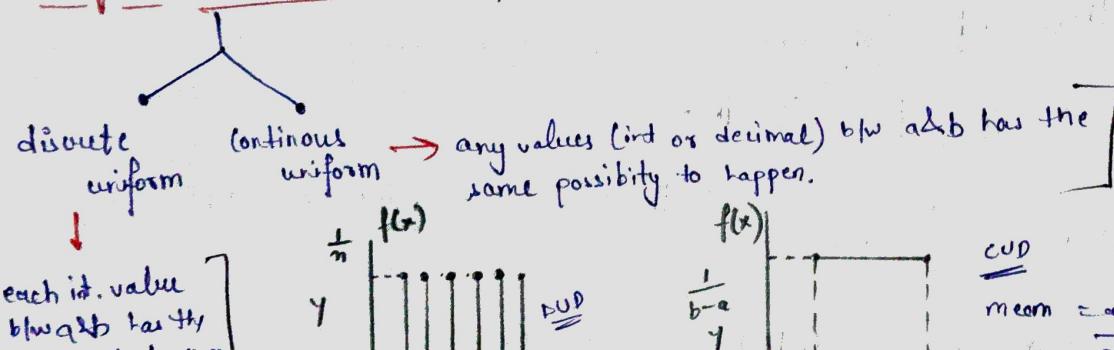
2) binomial

3) Poisson

4) normal / Gaussian

5) calculating prob. with Z-score normal Distribution.

uniform distribution :- all possible outcomes are likely to happen.



$$P(X=1) = 1/6$$

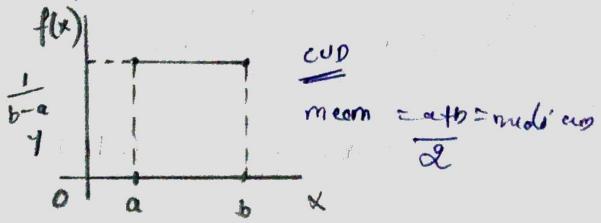
$$P(X=2) = 1/6$$

$$P(X=3) = 1/6$$

$$P(X=4) = 1/6$$

$$P(X=5) = 1/6$$

$$P(X=6) = 1/6$$



$$P(X=x) = \frac{1}{b-a} \text{ if } a \leq x \leq b$$

lower limit upper limit

binomial distribution :- tell the prob. of getting a certain num. of success in a fixed number of independent trials.

where, each trial has only two outcomes - success or failure.

what's models, how many times success happens out of fixed num. of tries.

> N-trials

> trials are independent

> $P(\text{success}) = p$

> $P(\text{failure}) = 1-p$

- random var X : total num. of success in 'n' triol.

$$P(X=x) = nC_x p^x (1-p)^{n-x}$$

$$nC_x \cdot p^x (1-p)^{n-x}$$

Poisson Distribution:-

Tells you the prob. of a certain number of events happening in a fixed period of time - space.

In short:- how many times will something happen in a fixed time or area.

Def: Dis. prob dist. of num. of events occurring in given t-period, given any no. time the event occurs over t period.

Let $x \rightarrow$ random discrete var [num. of events over a time]
ex → cell/hr at cell under.

$\lambda \rightarrow$ expected value (avg) of x , eg - avg. num. of cell in cell.

If x follows a poisson dist. $P(x=a) = \frac{\lambda^a e^{-\lambda}}{a!}$

[$e = 2.718$] math constant.
Euler's Num

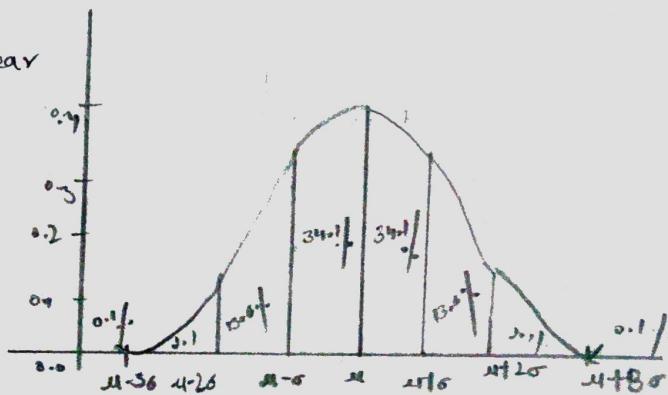
normal distribution:- a prob. distribution that is symmetric about the mean.

Gaussian showing the data around the mean, are more frequent in occurrence than data far from mean.

for all normal distribut's:

- 68.2% of observat's will appear within $+/- 1\text{std}$ of mean
- 95% of observat's will fall within $+/- 2\text{std}$
- 99.7 with $+/- 3\sigma(3\text{std})$.

(0.3%) → outliers



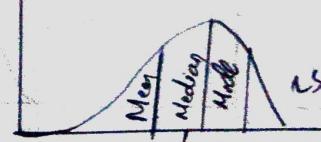
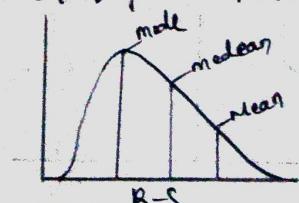
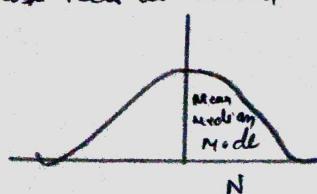
SKEWNESS

measures the distribution of data. Indicates whether data is distributed symmetric or not.

→ if symmetric → data is distributed normally.

→ Types → Symmetric → no skewness → -ve skewness

Skewness tells us which side the tail of the data is longer → right / left.



11

Transform Techniques (\rightarrow Normal Dist)

(12)

- log Transform
 - square
 - sq root
 - cube
 - cube root
 - Reciprocal
- } real polar / numpy

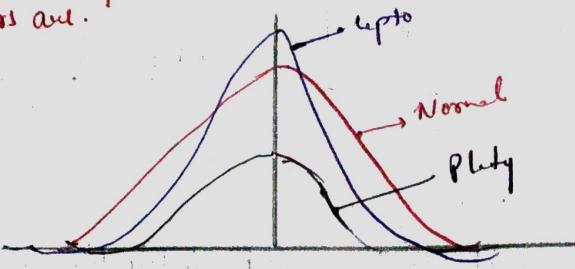
Kurtosis & measures the thickness of the distribution of data.

- degree of tailedness of the data measured by kurtosis.

- Tells, extent to which the distribution is more / less prone than normal dist.

Simply:- It measures how concentrated your data around the mean & how extreme your outliers are.

- Types
- PlatyKurtic
 - Mero Kurtic
 - LeptoKurtic

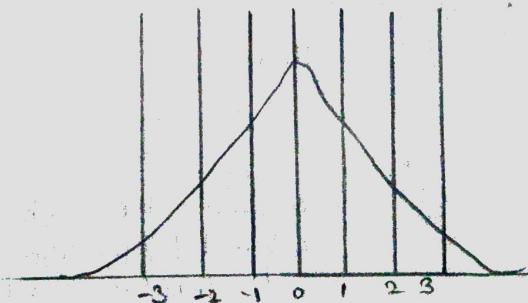
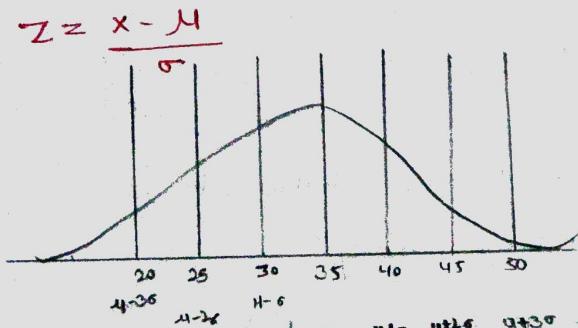


feature Scaling!:-

Data preprocessing step where we bring all numerical features to a similar scale (range or distributn). so that no single feature dominates others due to its Magnitude.

Calculating prob. with z-score for Normal Distributn. (explain Σ)

- i) find out how far the value of x is from μ . (in terms of σ)
- ii) The number named as "z-score", standard score or z-value.



Normal Distri (X)

$$① z = \frac{x-\mu}{\sigma}$$

$$z = \frac{\mu-3\sigma-\mu}{\sigma} = -3$$

$$z = \frac{\mu-2\sigma-\mu}{\sigma} = -2$$

$$z = \frac{\mu-\sigma-\mu}{\sigma} = -1$$

Standard Normal Dist (z)

$$z = \frac{\mu+\sigma-\mu}{\sigma} = 1$$

$$z = \frac{\mu+2\sigma-\mu}{\sigma} = 2$$

$$z = \frac{\mu+3\sigma-\mu}{\sigma} = 3$$

$$z = \frac{\mu+2\sigma-\mu}{\sigma} = 2$$

$$z = \frac{\mu+3\sigma-\mu}{\sigma} = 3$$

Normalization (Min-Max Scaling) :-

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

(13)

When to use:-

- need bounded values e.g. (neural networks)
- sensitive to outliers.

Covariance: statistical term that refers to a relationship b/w two random variables, on how one var. changes would impact the other one.

- covariance value can change from $-\infty$ to ∞ .
- +ve co-var shows direct relationship & represented by +ve num.
if one var. is \uparrow then also other is \uparrow . $A \uparrow \rightarrow B \uparrow$
- -ve co-var shows inverse relatn. $A \uparrow \rightarrow B \downarrow$
- value doesn't represent the magnitude of the relationship. only direct matters.

$$\text{covariance}(x, y) = \frac{1}{n} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$$

i.e. for population

where

$x_i \rightarrow$ data of value x

$y_i \rightarrow$ ————— y

$\bar{x} \rightarrow$ mean of x

$\bar{y} \rightarrow$ ————— of y

$N \rightarrow$ number of data values

for sample,
 $\therefore N-3 \approx n-1$

Correlation

Similar to covariance, also for measuring how 2 vars are related or moving together.

- ranges from $+1$ to -1 .
- +ve correlatn $\rightarrow A \uparrow \rightarrow B \uparrow$
- -ve correlatn $\rightarrow A \uparrow \rightarrow B \downarrow$ (inverse)
- Both direct & magnitude are useful, bigger num \rightarrow strong relatn.

$$P_{xy} = \text{correlatn}(x, y) = \frac{\text{cov}(x, y)}{\sqrt{\text{var}(x)} \sqrt{\text{var}(y)}} = \frac{\text{cov}(x, y)}{(\text{std}(x)) (\text{std}(y))}$$

Correlatn value	Meaning
0-0.3	weak linear relatn
0.3-0.7	Moderate L. relatn
0.7-1	Strong linear relatn

Hypothesis Testing:-

A statistical method that is used to making statistical decisions using experimental data.

- Hypothesis Testing is basically an assumption that we make about the population parameter. "Used to make decisions about population using sample Data".
ii) Hyp-Test helps to decide, based on sample Data, whether claim about popu. is likely T/F?
Why we use it? → an essential procedure in stats.
It evaluates two mutually exclusive statements about a population to determine which statement is better supported by sample data.

Hypo. Testing

Null Hypothesis (H_0)

Its a general statement or default position that there is no relation b/w two measured phenomenon among groups.
In other words its a basic assumption made based on domain or problem knowledge
→ assumes no diff / effect
ex - a company's product's ! = 50 unit/day.

Alternative Hypothesis. (H_1 , H_a)

In the hypothesis used in hypothesis testing that is contrary to null hypo. usually taken to be that observations are the result of a real effect (with some chance it may atⁿ happen)
→ assumes some effect / some difference

ex - company's product's ! = 50 unit/day

Tailed Test: tail → ends of prop. distribution curve.

When we perform Hypothesis testing, we check where the sample statistic falls. - in the left or right tail, or both tails.

Depending on the type of hypothesis , we use

- one tailed Test
- two tailed Test

One tailed Test: where, region of reject is on only one side of sampling distribution.
eg → marks > 800

Two tailed Test: the critical area of distribution is two sided. whether a sample is $>$ OR $<$ than a certain range of values.

eg → marks ! \neq 800

p-value :- p-value or calculated prob. of finding the observed or more extreme results when the null hypothesis (H_0) of study question is true.

Simply: It measures the prob. of that your data happened by random chance.

Significance value (α) :- A threshold you choose before testing to decide how shall the p-value must be to reject null hypoth.

$$\alpha = 0.5 \text{ (most common)}$$

$$\alpha = 0.01 \text{ (strict)}$$

$$\alpha = 0.10 \text{ (less strict)}$$

Data - Type	What we see	Type of test
1 categorical (c)		1-sample t-test
1-Numerical (n)	---	t-test
1N + 1C	---	t-test / ANOVA
2N	---	Correlation-test
2C		Chi Square test.

z-test and to compare the mean of the two given sample & infer whether they are from the same distribution or not.

We do implement z-test when the sample size less than 30.

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - (u_1 - u_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$