

Inferential Stats

* Basics of Probability

- Many a time you may require large amount of data for analysis which may need too much time and resources to acquire.
- In such situations we are forced to work with a "smaller sample of data", instead of having the entire data to work with.
- Situations like these arise all the time at big companies like Amazon, for eg Amazon QC dept wants to know what proportion of the products in its warehouse are defective, so instead of going through all their products, the Amazon QC team can just check a small sample of 1000 products and then find for this sample the defect rate (i.e. the proportion of defective rate)
- ⇒ then based on the sample's defect rate, the team can 'infer' what the defect rate is for all the products in the warehouse.
- "the process of "drawing" insights from sample data is called Inferential Statistics"

So In easy words "Inferential statistics" is a way of making guesses or conclusions about a larger group (called as Population) by looking at the data from a smaller group (called as sample).

e.g Instead of Asking 100000s of people in the city what their favourite fruit is, we could ask 100 people (a sample) and use that info to predict what the whole city (population) might prefer.

Probability

↳ How likely something is to happen.

Many events can't be predicted with total certainty, the best we can say is how likely they are to happen, using the idea of probability.

⇒ Tossing a coin : When a coin is tossed, there are two possible outcomes Heads (H) or Tails (T)

Also:

→ the probability of the coin landing H is $\frac{1}{2}$

→ the probability of the coin landing T is $\frac{1}{2}$

⇒ Throwing a Dice : When a single dice is thrown there are 6 possible outcomes 1, 2, 3, 4, 5, 6
the prob. of any 1 of them is $\frac{1}{6}$

Probability in general : $\frac{\text{No. of ways it can happen}}{\text{total No. of outcomes}}$

e.g. chance of rolling a 4 with dice

↙ no. of ways it can happen : 1 (there is only 1 face with a 4 on it)

total no. of outcomes : 6 (\because there are 6 faces altogether)

$$\Rightarrow \text{Probability} = \frac{1}{6}$$

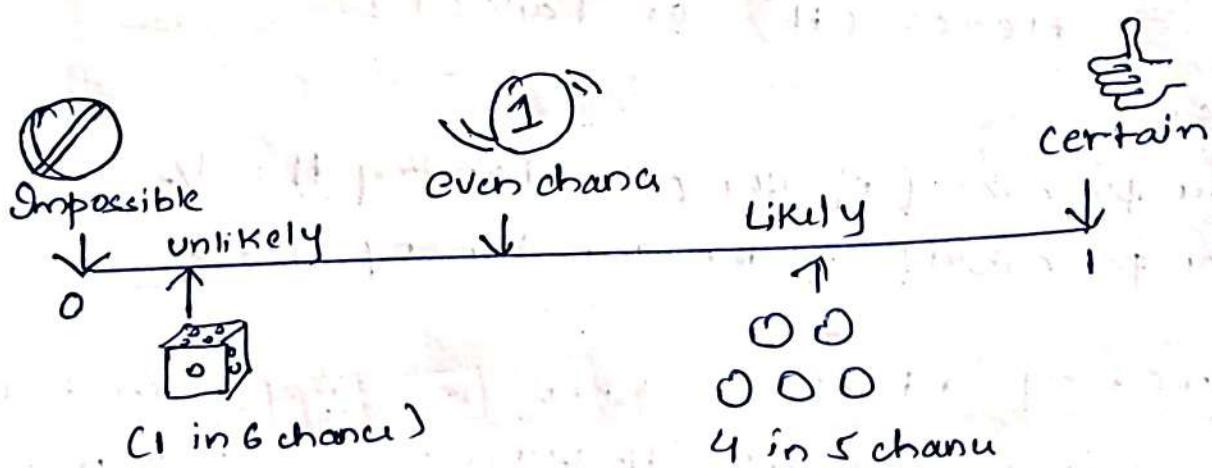
e.g. 5 marbles in a bag \rightarrow 4 blue and 1 red.

↙ P.C. of blue marble gets picked) $\Rightarrow \frac{4}{5}$

$\therefore 4 \rightarrow$ there are 4 blue marbles

$5 \rightarrow$ there are 5 marbles in total

Probability Line : we can show probability on a Probability line



Probability is always b/w 0 and 1

Probability is just a guide

It does not tell us exactly what will happen.
It is just a guide.

e.g. If we toss a coin 100 times, how many heads will come up?

Probability says heads have $\frac{1}{2}$ chance so we can expect 50 heads, but in reality we can get 48 heads or 55 heads etc.

✳ Some words have special meaning in Probability

⇒ Experiment : A repeatable procedure with a set of possible results

e.g. throwing dice

We can throw the dice again and again so it is a repeatable.

→ the set of possible outcomes for any single throw is $\{1, 2, 3, 4, 5, 6\}$

⇒ Outcome : A possible result

e.g. 6 is one of the outcomes of a throw of a dice.

⇒ Total : A single performance of an experiment

e.g. conducted a coin toss experiment & after 4 trials the results received

Outcome	Trial 1	Trial 2	Trial 3	Trial 4
H	✓	✓	✓	✓
T				✓

So the 3 trials had the outcome Head and
One trial had the outcome tail.

⇒ Sample Space : All the possible outcomes
of an experiment.

eg choosing a card from a deck

→ 52 cards in deck

so sample space is all 52 possible cards

the sample space is made up of sample
points

Sample point : Just one of the possible outcomes

eg : Deck of cards

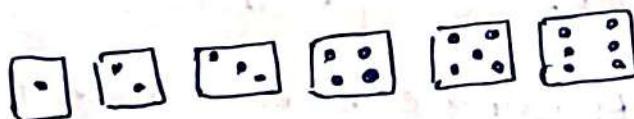
→ 5 of the clubs is a sample point

→ the King of hearts is a sample
point

Note "King" is not a sample point, there
are 4 Kings so 4 different sample points

eg : throwing dice

there are 6 diff sample points in the
sample space



⇒ Event: One or more outcomes of an experiment

Ex of Events

An event can be just one outcome

→ getting a tail when tossing a coin

→ Rolling a "5"

An event can include more than one outcome

→ choosing a king from a deck of cards

(Any of 4 kings)

→ Rolling an "even no." (2, 4 or 6)

① Additional Rule of Probability

The addition rule of probability is a fundamental concept used to find the probability of the union of two events (i.e. the probability that at least one of the events occur). The formula depends on whether the events are mutually exclusive or not.

① When events are mutually exclusive

If two events A and B cannot occur at the same time (no overlap), the probability of A or B occurring is the sum of their individual probabilities.

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

Eg Rolling a dice, let A be rolling a 2 and B be rolling a 5, so we can't roll both a 2 and a 5 in one roll, A and B are mutually exclusive.

$$P(A) = \frac{1}{6}, P(B) = \frac{1}{6} \Rightarrow P(A \cup B) : P(A) + P(B)$$

$$\Rightarrow P(A \cup B) = \frac{1}{6} + \frac{1}{6} = \frac{2}{6} = \frac{1}{3}$$

② When events are not mutually exclusive

→ If two events A and B can occur simultaneously (overlap exists), the formula accounts for the overlap by subtracting the probability of both events happening together.

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

Eg Drawing a card from a deck, A be drawing a heart, B be drawing a face card. Some cards (eg the King of hearts) satisfy both A and B so they aren't mutually exclusive.

$$P(A) = \frac{13}{52} \quad P(B) = \frac{12}{52} \quad P(A \cap B) = \frac{3}{52}$$

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

$$\Rightarrow \left(\frac{13}{52} + \frac{12}{52} - \frac{3}{52} \right) : \frac{22}{52} = \frac{11}{26}$$

Keypoints to remember here

- $P(A \cup B) = P(A) + P(B)$ if events are mutually exclusive
- $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ for non-mutually exclusive events
- $P(A \cap B)$ is critical in the second case to avoid double counting.

Examples

① Mutually Exclusive Events

toss a coin

$$\hookrightarrow A: \text{getting Heads} \Rightarrow P(A) = 0.5$$

$$\hookrightarrow B: \text{getting Tails} \Rightarrow P(B) = 0.5$$

What is the prob of getting heads or tails $P(A \cup B)$

$$\Rightarrow 0.5 + 0.5 = 1$$

② Non Mutually exclusive

In a class of 30 students

$$18 \text{ students like math} \Rightarrow P(A) = 18/30$$

$$12 \text{ students like science} \Rightarrow P(B) = 12/30$$

$$5 \text{ like both} \quad P(A \cap B) = 5/30$$

$$\Rightarrow P(A \cup B) = 18/30 + 12/30 - 5/30 = 25/30 = 5/6$$



Multiplication Rule of Probability

Used to find the probability of two events occurring together (i.e. the intersection of two events)

This formula depends on whether the events are dependent or independent

① Independent Events

Two events A and B are independent, if the occurrence of one does not affect the probability of the other.

→ For independent events, the multiplication rule is $P(A \cap B) = P(A) \cdot P(B)$

e.g. tossing 2 coins → A getting Heads on first coin $P(A) = 0.5$

and B getting heads on second coin $P(B) = 0.5$

∴ the outcome of two coins are independent

$$P(A \cap B) = P(A) \cdot P(B) = (0.5)^2 = 0.25$$

② Dependent Events : 2 events are said to be dependent if the occurrence of one affects the probability of the other

for Dependent events the multiplication rule is

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

Here $P(B|A)$ is the conditional prob of B given that A has occurred.

e.g Drawing cards from a deck without replacement

A \Rightarrow drawing a heart on first draw, $P(A) = 13/52$

B \Rightarrow drawing another heart on 2nd draw

After drawing 1 heart, there are only 12 hearts

left out of 51 cards

$$P(B|A) = 12/51$$

$$\Rightarrow P(A \cap B) = P(A) \cdot P(B|A) = \frac{13}{52} \times \frac{12}{51} = \frac{1}{17}$$

examples

① For independent events

you roll a dice and flip a coin

A : Rolling a 6 $\Rightarrow P(A) = 1/6$

B : Getting heads $\Rightarrow P(B) = 1/2$

② Deck of 52 cards (without replacement)

A : Draws a king $\Rightarrow 4/52$

B : Draws another king $\Rightarrow 3/51$

$$\Rightarrow \frac{4}{52} \times \frac{3}{51}$$

OR \Rightarrow Union

AND \Rightarrow Intersection

$$\Rightarrow \frac{3}{13 \times 51}$$

Permutations And Combinations

Permutations and combinations are fundamental concepts in probability and counting.

they help us calculate the no. of ways to arrange or select objects depending on whether order matters.

① Permutations (Order Matters)

Permutations count the arrangements of objects where the order is important.

Formula : $P(n, r) = \frac{n!}{(n-r)!}$

where

$n \rightarrow$ total no. of objects

$r \rightarrow$ no. of objects to arrange

$! \rightarrow$ Factorial (eg $4! = 4 \times 3 \times 2 \times 1$)

eg How many ways can you arrange 3 letters from A, B, C, D?

here $n=4$, $r=3$

$$P(4, 3) = \frac{4!}{(4-3)!} = \frac{4 \times 3 \times 2 \times 1}{1} = 24$$

so there are 24 arrangements ABC, ACB, BAC, BCA,

② Combinations (order doesn't matter)

Combinations count selections of objects where the order does not matter.

formula for combinations

$$C(n, \gamma) = \frac{n!}{(n-\gamma)! \gamma!}$$

| $n \rightarrow$ total nos of objects
 $\gamma \rightarrow$ nos of objects to choose

eg

How many ways can you choose 3 letters from A, B, C, D where order doesn't matter

$$n = 4, \gamma = 3$$

$$C(4, 3) = \frac{4!}{1! \cdot 3!} = \frac{4 \times 3!}{3!} = 4$$

So total of 4 combinations

A B C, A B A, A C D, B C D

Definitions

* Factorial : there are $n!$ ways of arranging n distinct objects into an ordered sequence, permutations where $n=\gamma$.

* Combination : the number of ways to choose a sample of γ elements from a set of n distinct objects where order does not matter and replacements are not allowed.

(*) Permutation : the nos of ways to choose a sample of σ elements from a set of n distinct objects where order does matter and replacements are not allowed, when $n = \sigma$ this reduces to $n!$ a simple factorial of n .

(*) Combination Replacement : the nos of ways to choose a sample of σ elements from a set of n distinct objects where order does not matter and replacements are allowed.

(*) Permutation Replacement : the nos of ways to choose a sample of σ elements from a set of n distinct objects where order does matter and replacements are allowed.

* $n \rightarrow$ Population $\sigma \rightarrow$ subset of n or sample

Combinations formula :

$$C(n, \sigma) = \frac{n!}{(n-\sigma)! \sigma!}$$

The formula shows us the number of ways a sample of " σ " elements can be obtained from a larger set of " n ".

Combination Problem

* Choose 2 prizes from a set of 6 prizes

Q You won first place in a contest and are allowed to choose 2 prizes from a table that has 6 prizes. How many diff combination of 2 prizes can you possibly choose.

$$\Rightarrow r = 2, n = 6$$

$$6C_2 = \frac{6!}{4!2!} = \frac{6 \times 5}{2} = 15 \text{ ways}$$

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

* Choose 3 students from a class of 25
Q She wants to find out how many unique teams of 3 can be created from her class of 25

$$25C_3 = \frac{25!}{23! \times 2!} = \frac{25 \times 24 \times 23}{12} = 2300$$

* Choose 4 menu items from a menu of 18 items
(These orders is not important just the combination matters)

$$18C_4 = \frac{18!}{14! \times 4!} = \frac{18 \times 17 \times 16 \times 15}{4 \times 3 \times 2} = \frac{108}{18 \times 17 \times 10} = 3060$$

Handshake Problem

In a group of n people, how many diff handshakes are possible?

→ Let's find the total handshakes possible, that is to say if each person shook hands once with every other person in group, what is the total nos of handshakes that occur?

⇒ each person in the group will make a total of $n-1$ handshakes, so there are n people, that would be n times $(n-1)$ total handshakes.
In other words total nos of people multiplied by nos of handshakes that each can make will be the total handshakes.

A grp of 3 would make a total of $3(3-1) = 6$.
each person registers 2 handshakes with the others 2 people in grp. ∴ total handshakes ⇒ $n(n-1)$

↳ this includes each handshake twice
1 with 2, 2 with 1, 1 with 3,
3 with 1, 2 with 3, 3 with 2

∴ Original Ques wants to find diff handshakes

$$\Rightarrow \frac{6}{2} = 3$$

Handshake Problem as a combination prob

We can solve Handshake Problem as a combination Problem as $C(n, 2)$

$n \rightarrow$ nos of people in grp
 \Rightarrow 2 (the nos of people involved in a diff handshake)

the order of the items chosen in the subset does not matter so for a group of 3 it will count 1 with 2, 1 with 3, 2 with 3

$$n_{C_2} = \frac{n!}{(n-2)! 2!} = \frac{n(n-1)}{2}$$

Sandwich Combination Problem

It's a classic math problem and asks something like how many sandwich combinations are possible? calculate the possible sandwich combinations if we can choose 1 item from each of the 4 categories

- 1 bread from 8 options
- 1 meat from 5 options
- 1 cheese from 5 options
- 1 topping from 3 options

As the multiplication of number possible options in each category $\rightarrow 8 \times 5 \times 5 \times 3 = 600$ combination

In terms of combinations even
 the nos of possible options for each category
 is equal to the nos of possible combinations
 for each category, so we are only making
 1 selection. for eg $C(8,1) = 8$, $C(5,1) = 5$ and

$$C(3,1) = 3$$

$$\Rightarrow 8C_1 \times 5C_1 \times 5C_1 \times 3C_1$$

$$8 \times 5 \times 5 \times 3 = 40 \times 15 = 600$$

Random Variables

Random Variable (X) converts outcomes of experiments to something measurable.

for eg : As a data analyst at a bank, we are trying to find out which of the applicants are more likely to default on their loan.

cust-no	Income	Loan-amt	Dependents	Is-default
1	10lakh	75lakh	3	Yes
2	15lakh	50lakh	2	No
3	20lakh	40lakh	1	No.

Now instead of processing Yes/No , It's easier to define a random variable X indicating whether the customer is predicted default or not.

the values will be assigned according to below rule

$X = 1 \Rightarrow$ customer defaults

$X = 0 \Rightarrow$ customer doesn't default

So data will be like

cust-no	Income	Loan-amt	Dependants	X (random var)
1	10 lakh	75	3	1
2	15 lakh	50	2	0
3	20 lakh	40	1	0

Now the table is more quantified, ie converted to nos now the data is entirely in quantitative terms

Probability Distribution

A probability distribution is a statistical function that describes all the possible values and likelihoods that likelyhood that a random variable can take within a given range. The range is bounded by the variable's minimum and maximum possible values but the specific shape of distribution depends on the type of variable and the context.

Key Concepts of Probability Distribution

- ① Random Variable : A variable that can take on different values each with an associated Probability.
- Discrete random variable : taken on a finite or countable set of values eq (rolling a dice, number of heads in coin tossed).
- Continuous Random Variable : Can take any value within a range eq height, weight, temperature

② Probability Distribution Function (PDF)

2.1 Probability Mass Function (PMF) :
for discrete variables, It is the probability mass function (PMF) for continuous variables It gives the probability density.

2.2 Cumulative Distribution Function (CDF)

the cumulative probability that a random variable is less than or equal to a particular value.

3. Types of Probability Distributions

→ Discrete Distributions

→ Binomial Distribution : Models the nos of successes in a fixed number of trials

e.g flipping a coin

→ Poisson Distribution : Models the nos of occurrences of an event in a fixed interval or time or space.

→ Continuous Distributions

→ Normal distribution : the classic bell shaped curve, many natural phenomena follow this distribution

→ Uniform distribution : Every outcome in a range is equally likely.

Note : The probabilities of all outcomes sum up to 1.

For continuous distributions, the area under the curve (total probability) is

e.g. for Discrete Probability Distribution :-

↳ Tossing a fair dice

Random variable x outcomes $(1, 2, 3, 4, 5, 6)$

↳ For continuous Probability distribution

Heights of students in class, Modelled by Normal dist with Mean μ and σ^2

Expected Value

Definition: Random variable X can take values $x_1, x_2, x_3 \dots x_n$ such that

$$\text{Expected Value } E(X) = x_1 * P(X=x_1) + x_2 * P(X=x_2) + x_3 * P(X=x_3) + \dots + x_n * P(X=x_n)$$

Note: Here $P(X=x_i)$ denotes the probability that X is equal to x_i

$$\text{eg } E(X) = 0 * P(X=0) + 1 * P(X=1) + 2 * P(X=2) \dots$$

the expected value should be interpreted as the "avg value" you get after the experiment has been conducted an infinite number of times

Q A Casino game contains numbers 0 to 36 written in irregular sequence, the player can bet on any no from 0 to 36, for eg Kaiti bets £100 on nos 5, now a ball would be dropped in a wheel which is then given a spin, If the ball lands on the pocket marked 5, Kaiti would win $(\pm 100) \times 36 = \pm 3600$ resulting in net winnings of $\pm 3600 - \pm 100 = \pm 3500$, however if the ball lands in any other pocket \rightarrow loss of £100

let's see what is the expected value for Kriti's net winnings if she plays this game and bet £100 on Number 5.

$$E(x) = \sum p_i x_i$$

Here

p_i \Rightarrow Probability of outcome i

x_i \Rightarrow Net winning for outcome i

Step 1 : Identify outcomes and probabilities

① there are 37 pockets on the wheel (nos 0 to 36)

② the probability of the ball landing on any specific no

$$P(\text{Landing on } 5) = \frac{1}{37}, P(\text{Not Landing on } 5) = \frac{36}{37}$$

③ the net winning for each case

if ball lands on 5 $\Rightarrow 3500 \Rightarrow (3600 - 100)$

if ball lands on another no $\Rightarrow -100 \Rightarrow (0 - 100)$

Step 2 : Apply formula

$$EV(x) = P(\text{Landing on } 5) * x(\text{Landing on } 5) + P(\text{Landing on another no}) * x(\text{Landing on another no})$$

$$EV(x) = \frac{1}{37} (3500) + \left(\frac{36}{37}\right) (-100)$$

$$\Rightarrow 94.59 \rightarrow 94.30 - 2.71$$

So the expected value of Kriti's net winning is -2.71

Probability without Experiment

↳ for this we will use Binomial Distribution
the probability of an event without conducting any experiment is determined through theoretical Probability, which is based on reasoning and calculations rather than actual outcomes.

formula for theoretical Probability :-

$$P(E) = \frac{\text{Number of favourable outcomes}}{\text{total nos of Possible outcomes}}$$

Key points :-

favourable outcomes : the outcomes that meet the condition of the event you are studying.

total outcomes : All the outcomes that can possibly occur in the given situation.

e.g. ① Rolling a fair die

$$\rightarrow \text{total possible outcomes} = 6$$

$$\rightarrow \text{favourable outcome for rolling a 4} \Rightarrow 1$$

$$\rightarrow \text{Probability} = \frac{1}{6}$$

② Drawing a card from a standard deck

$$\rightarrow \text{total possible outcomes} = 52$$

$$\rightarrow \text{favourable outcome for drawing heart} = 13$$

$$\rightarrow \text{Prob} = \frac{13}{52} = \frac{1}{4}$$

Q A bag contains 5 red balls and 3 blue balls
 If a ball is selected at random, what is the probability of

- ① Picking a red ball?
- ② Picking a blue ball?

Ans ① Identify the total nos of balls

$$\text{total balls} = 5 + 3 = \underline{\underline{8}}$$

② Calculating Probabilities

$$P(\text{Red}) = \frac{5}{8} \quad | \quad P(\text{Blue}) = \frac{3}{8}$$

Q If we randomly pick two balls from the bag without replacement what is the probability that

- ① Both Balls are red?
- ② One is red & another is blue?

Step ① Use combination to calculate probability

total ways to select 2 balls from the bag:

$${}^8C_2 = \frac{8!}{6!2!} = \frac{8 \times 7}{2} = \underline{\underline{28}}$$

Prob of selecting 2 red balls

ways to choose 2 reds out of 5 red

$${}^5C_2 = \frac{5!}{3!2!} = \frac{5 \times 4}{2} = 10$$

$$\Rightarrow \text{Prob}(2 \text{ red}) = \frac{10}{28} = \frac{5}{14}$$

Prob of 1 red ball and 1 blue ball

ways to choose 1 red ball from 5 red and
1 blue ball from 3

$${}^5C_1 \times {}^3C_1 = \frac{5!}{4! \times 1!} \times \frac{3!}{2! \times 1!} = 15$$

$$\Rightarrow P(1 \text{ red } 1 \text{ blue}) = \frac{15}{28}$$

Q Bag contains \rightarrow 5 red balls
 \downarrow 3 blue balls } Balls are drawn
with replacement
If we draw 4 balls what's the probability
of getting exactly 3 red balls and 1 blue ball?

(Step) Understanding Keypoints

→ Since the balls are drawn with replacement,
the probability of drawing a red or blue ball
remains constant for each draw.

$$P(\text{Red}) = 5/8 \quad P(\text{Blue}) = 3/8$$

→ Order doesn't matter, and sequence with
3 red ball and 1 blue ball counts

\Rightarrow for this problem we will use "Binomial
formula for exact probability"

The Binomial formula :

It is used to calculate the probability of achieving a specific number of success (K) in a fixed no of independent trials (n), when the probability of success ($P(\text{success})$) remains constant for each trial.

→ formula :

$$P(c_k) = n_{c_k} * P^k \cdot (1-P)^{n-k}$$

where :

n : total nos of trials (eg 4 draws in problem)

K : Desired numbers of success (eg 3 red balls in problem)

P : Probability of success on a single trial

$$\text{eg } P(Red) = 5/8$$

$1-P$: Probability of failure on a single trial

$$\text{eg } (P(B)) = 3/8)$$

n_{c_k} : Binomial Coefficient, which calculates the nos of ways to arrange k success in n trials

Understanding Components

- ① Binomial Coefficient (n_{C_k}) represents the no. of ways to choose K successes from n trials and it's denoted by

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

It's important because in binomial distribution the order of success doesn't matter for eg → If you are drawing 4 balls and want 3 red ball, the arrangement RRRB, RBRR etc are all valid and counted in 4C_3

- ② Probability of Success: P^K

this is the prob of getting K success in n trials for eg

→ Prob of drawing red ball is $P = 5/8$

→ Prob of drawing blue ball is $(5/8)^3$ → this term becomes $(5/8)^3$

- ③ Prob of failure $(1-P)^{n-k}$

this account for trials that don't result in success

→ Prob of drawing blue ball = $3/8$

→ if there is 1 blue ball ($n-k=1$) $\Rightarrow (3/8)^1$

Applying on our ques ($n=4$, $K=3$ (red)), $P(R) = 5/8$

$$1-P = 3/8$$

$$P(K=3) = {}^4C_3 \left(\frac{5}{8}\right)^3 \left(\frac{3}{8}\right)^1 \Rightarrow \frac{4!}{1!3!} \left(\frac{5}{8}\right)^3 \left(\frac{3}{8}\right)^1$$

$$\Rightarrow \frac{4 \times 375}{4096} = \frac{375}{1024} \approx 0.366$$

$P \Rightarrow$ Prob of choosing red ball

$1-P \Rightarrow$ Prob of choosing blue ball

⇒ Probability Distribution for General Probability P

Now lets extend this case to the more generic one

X	$P(X=x)$
0	$(1-P)^4 \rightarrow$ OR 4 Blue
1	$4P(1-P)^3 \rightarrow$ 1R 3B
2	$6P^2(1-P)^2 \rightarrow$ 2R 2B
3	$4P^3(1-P) \rightarrow$ 3R 1B
4	$P^4 \rightarrow$ 4R 0B

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

X	$P(X=x)$
0	${}^n C_0 (P)^0 (1-P)^n$
1	${}^n C_1 (P)^1 (1-P)^{n-1}$
2	${}^n C_2 (P)^2 (1-P)^{n-2}$
:	⋮
n	${}^n C_n (P)^n (1-P)^0$

Conditions for Binomial Probability Distribution

- ① total no. of trials is fixed at n
- ② Each trial is binary i.e. has only two possible outcomes (success or failure)
- ③ Prob of success is same in all trials
- ④ $P(X=x) = P(\text{getting } x \text{ successes in } n \text{ trials})$

$$= {}^n C_x (P)^x (1-P)^{n-x}$$

Cumulative Probability

Cumulative probability is the probability that a random variable takes a value less than or equal to a given value, in other words it is the sum of probabilities for all possible outcomes upto and including a specific outcome.

For binomial distribution, cumulative probability answers Ques like
→ What is the probability of getting atmost 3 red balls
↳ this would mean summing the prob of getting 0, 1, 2, and 3 red balls.

* Cumulative Probability in a Binomial Distribution
↳ for this we sum the probabilities for each k (number of success) from 0 upto a target value K .

$$P(X \leq K) = \sum_{i=0}^K \binom{n}{i} p^i (1-p)^{n-i}$$

where $n \rightarrow$ total trials

$p \rightarrow$ probability of success

$K \rightarrow$ upper limit for success

$\binom{n}{i} \rightarrow$ Binomial Coefficient for i successes in n trials.

Example

Cumulative prob. for Atmost 3 red balls

Problem: Using the same scenario of 5 red balls and 3 blue balls, if we draw 4 balls with replacement, what is the probability of getting at most 3 red balls.

$$\text{Given: } n=4 \quad (\text{4 draws})$$

$$P = 5/8 \quad (\text{Prob of red})$$

$$1 - P = 3/8 \quad (\text{Prob of blue})$$

What we want

$$P(X \leq 3) = P(X=0) + P(X=1) + P(X=2) + P(X=3)$$

We calculate the prob for each K using binomial formula.

$$1) \text{ Prob}(X=0) = {}^4C_0 \left(\frac{5}{8}\right)^0 \left(\frac{3}{8}\right)^{4-0} \\ 1 \times 1 \times \frac{3 \times 3 \times 3 \times 3}{8 \times 8 \times 8 \times 8} = \underline{\underline{\frac{81}{4096}}}$$

$$2) \text{ Prob}(X=1) = {}^4C_1 \left(\frac{5}{8}\right)^1 \left(\frac{3}{8}\right)^{4-1} = \underline{\underline{\frac{540}{4096}}}$$

$$3) \text{ Prob}(X=2) = \underline{\underline{\frac{1350}{4096}}} \quad 4) \text{ Prob}(X=3) = \underline{\underline{\frac{1500}{4096}}}$$

$$\Rightarrow P(X \leq 3) = \frac{81 + 540 + 1350 + 1500}{4096} = \underline{\underline{\frac{3471}{4096}}} = 0.84$$



Understanding PDF and CDF

When dealing with random variables, Probability Density Functions (PDFs) and Cumulative Distribution Functions are essential tools to describe the behaviour of continuous random variables.

① What is a PDF (Probability Density Func)?

→ Imagine you are studying the daily commute times of people, the commute time vary continuously like someone's commute could be 15.4 minutes, 20.8 minutes, or even 30.1 minutes → the PDF is a func that shows how likely it's for a random variable (like commute time) to be in a small range of values.

→ For eg the PDF tells us how likely it is for a commute to be in between 20 and 25 minutes

Key Idea

the area under the curve of the PDF b/w two values (say 20 and 25) gives the prob of the random variable falling within that range.

Why can't we find the probability at an exact value

for continuous data the probability of a single exact value (eg $P(X=20)$) is 0 ∵ there are infinitely many possible values for a continuous var

and the probability gets spread out over a range instead we will focus on intervals like $P(20 \leq x \leq 25)$

PDF Properties

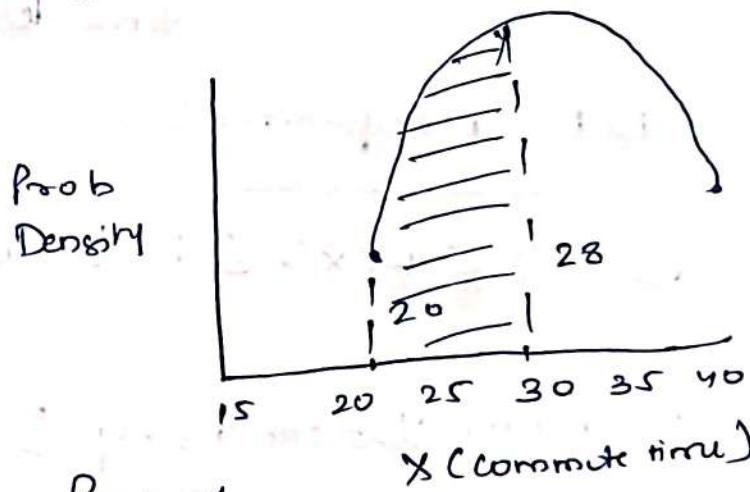
- ① The PDF is non-negative : $f(x) \geq 0$ for all x .
this means the graph of PDF is always above horizontal axis.

Probability Density func is a function in which area under the curve gives us the cumulative Probability,

for eg. the Area
under the curve

b/w 20 the smallest

Possible value of x and
28 gives the cumulative Prob of
 $x = 28$.



Using PDF we can calculate the probability like
 $P(25 \leq x \leq 35)$; the prob that a commute
is b/w 25 and 35 minutes

What is CDF (Cumulative Distribution Function)

The CDF is the cumulative probability upto a certain value of x , It tells us the prob that the random variable X is less than or equal to x .

$$F(x) = P(X \leq x)$$

How does CDF work?

The CDF is calculated by adding up all the Probabilities from the smallest possible value upto x .

$$F(x) = \int_{-\infty}^x f(t) dt$$

It is essentially the area under the PDF curve from the leftmost side usually $-\infty$ to x .

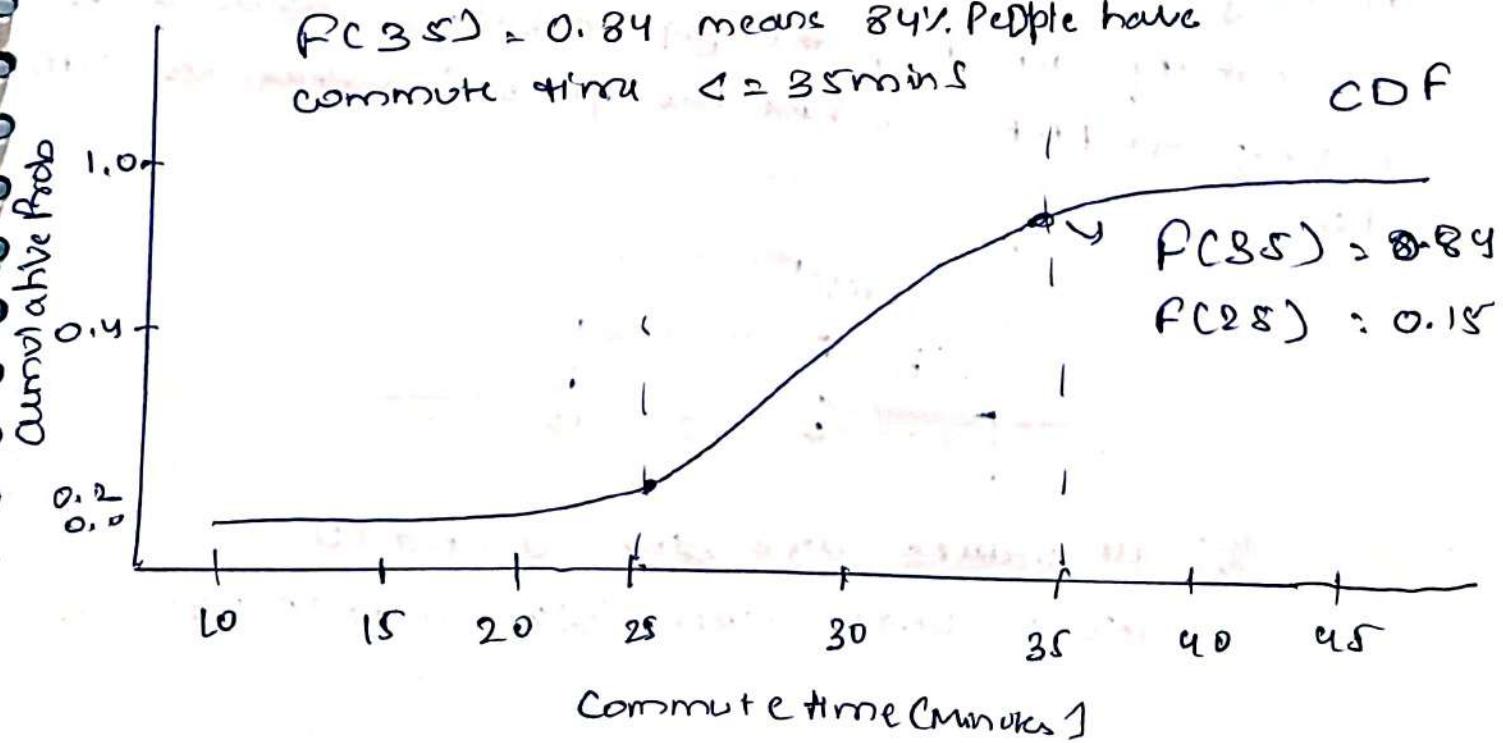
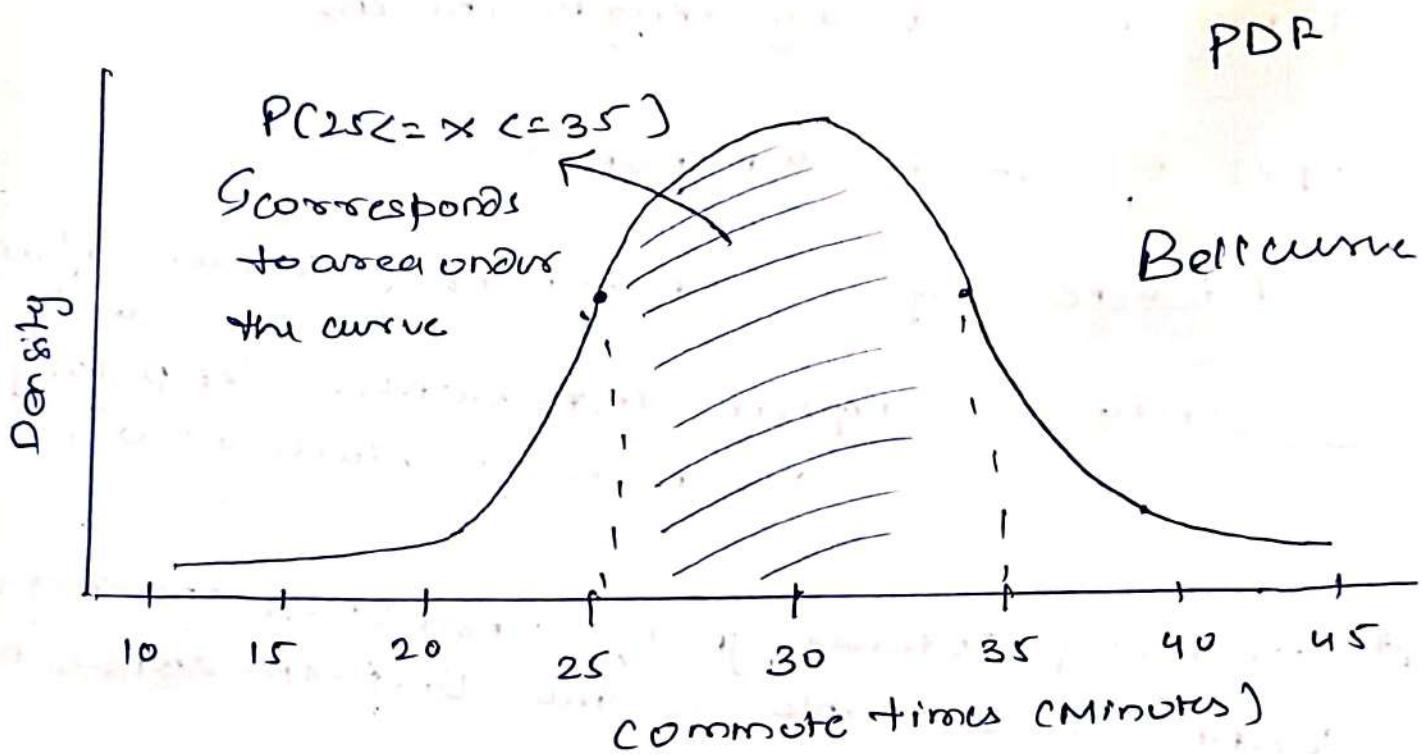
e.g. $F(35) = .85 \Rightarrow$ Means 85% of people have commute time less than or equal to 35 min.

CDF Properties

- ① $0 \leq F(x) \leq 1$: Cumulative prob grows from 0 to 1
- ② Non Decreasing : The CDF either stays constant or increases as x increases. It can't decrease because probability only accumulates.
- ③ Boundary Values
 - $\rightarrow F(x) \rightarrow 0$ as $x \rightarrow -\infty$: The prob of very small values is close to 0.
 - $\rightarrow F(x) \rightarrow 1$ as $x \rightarrow \infty$: The cumulative prob eventually reaches 1 as we include all possible values.

Visualizing PDF and CDF

- ① PDF curve showing the likelihood of diff commute times.
- ② CDF curve showing cumulative probability



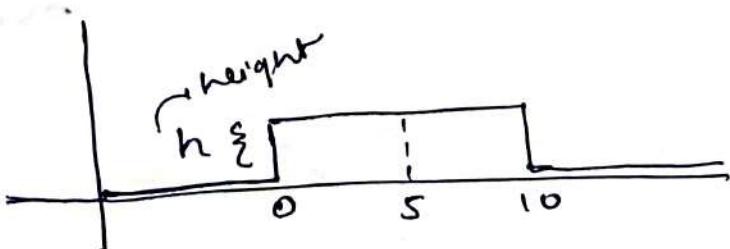
Uniform Distribution

In uniform distribution, every outcome has same chance of happening, you can think of it like drawing a number from a hat where all nos are equally likely to chosen.

Types of Uniform distribution

- ① Discrete Uniform distribution (like rolling a fair die)
- ② Continuous Uniform distribution (like picking a random nos b/w 0 to 10)
- ③ Commonly observed type of Distribution among continuous variable is the Uniform distribution

In uniform PDF, all possible values have same Probability density, the figure below shows such a uniform PDF, where the possible values are 0 to 10.



∴ all values are b/w 0 and 10
area under curve b/w 0 and 10 is 1

Clearly this area is the area of rectangle with length 10 and unknown height h

$$\Rightarrow 10 \times h = 1 \Rightarrow h = \frac{1}{10} = \underline{\underline{0.1}}$$

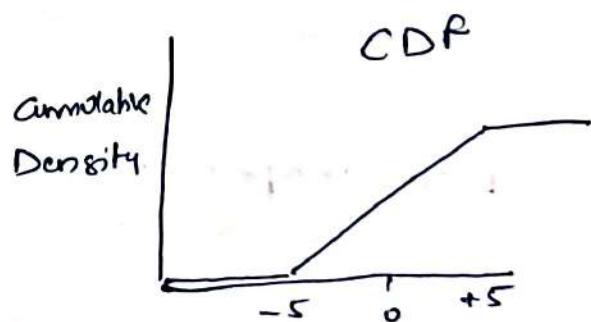
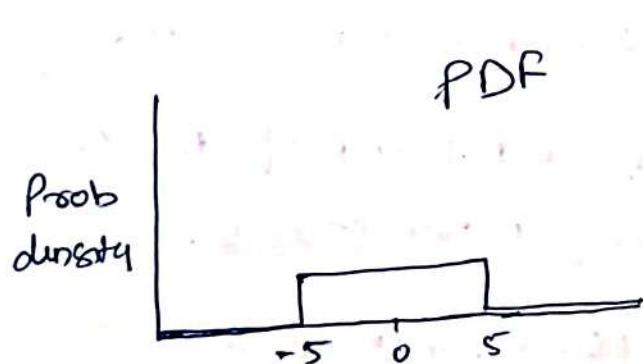
$\therefore 1 \times b = \text{Area}$ \therefore PDF for all values between 0 to 10 is $\underline{\underline{0.1}}$

Now cumulative prob for $x=0.5$

$$0.1 \times 0.5 = 0.05$$

\hookrightarrow $\begin{matrix} h \\ \text{length} \end{matrix}$

Note: PDF's are more commonly used in real life, the reason is that it is much easier to see patterns in PDFs as compared to CDFs



PDF shows uniformity as the probability density value remains constant for all possible values. whose CDF doesn't show any trend that shows a variable is uniformly distributed.

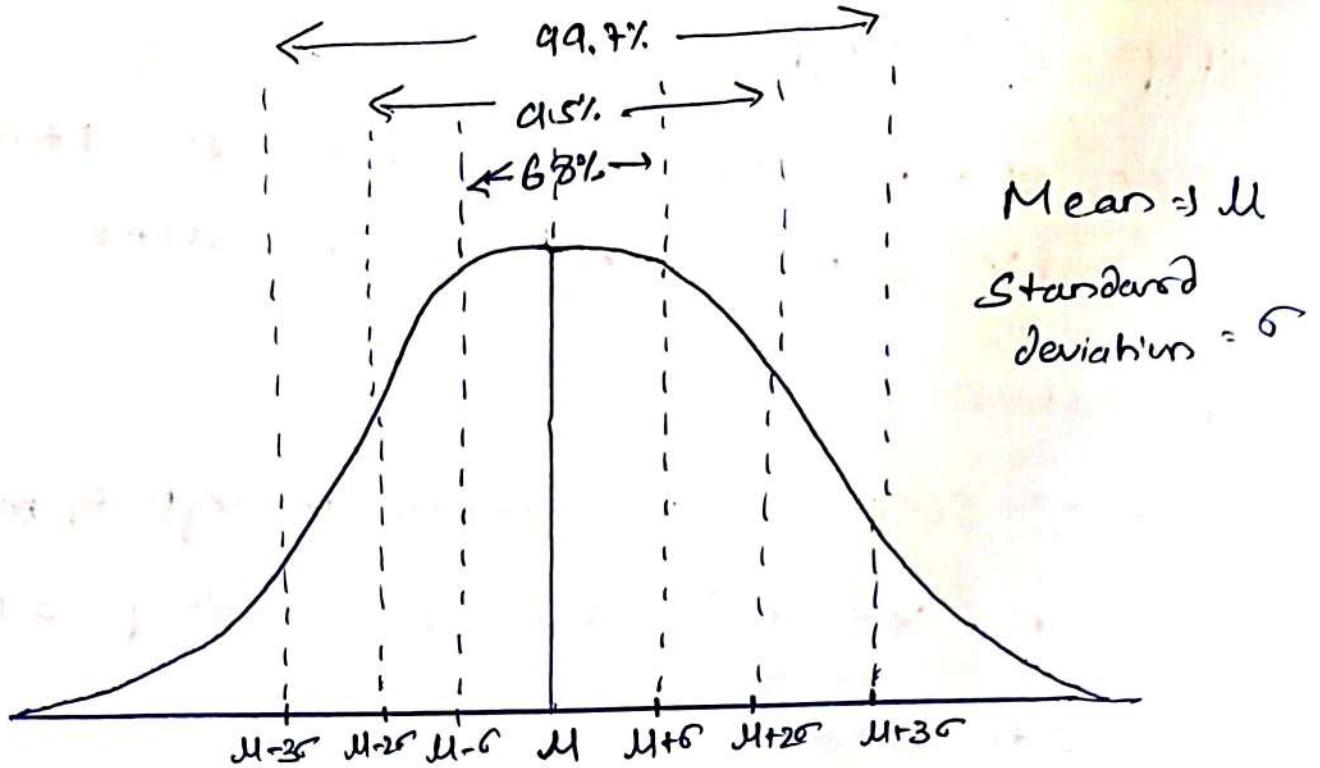


Normal Distribution

Normal distribution also called as bell curve or Gaussian distribution is one of the most common distribution. It represents a situation where data is symmetrically distributed around the mean, meaning most of the values cluster around the avg and fewer values are found as you move away from it.

Key features of Normal distribution

- Symmetry : the curve is perfectly symmetric around the mean (μ)
- Mean = Median = Mode : these 3 measures of central tendency are same
- Bell shaped curve : highest point of the curve occurs at the mean.
- Empirical Rule (68-95-99.7 Rule)
 - ① 68% of data lies within 1st standard deviation (σ) of the mean.
 - ② 95% of data lies within 2 standard deviations
 - ③ 99.7% of data lies within 3 standard deviations



eq of normal distribution

Scenario

↳ Let's say the heights of adult males in a country are normally distributed with mean (μ) of 170 cm and standard deviation (σ) of 10 cm.

↳ Most men are around 170 cm tall.
Fewer men are shorter than 140 cm ($\mu - 3\sigma$) and longer than 200 cm ($\mu + 3\sigma$).

Q Let's say we need to find the cumulative prob for a random variable X which is normally distributed, we do not know what the value of X is or, for that matter, what is the value of μ and σ is we only know that $x = \mu + \sigma$
Can you find cumulative prob ie prob of var being less than $\mu + \sigma$?

Given: random var X is normally distributed

we know $x = \mu + \sigma$ we need to find $P(x \leq \mu + \sigma)$

Step by Step Soln

- 68% of Data lies b/w $\mu - \sigma$ to $\mu + \sigma$
→ 34% Data lies b/w μ and $\mu + \sigma$

Cumulative prob

↳ 50% 50% of Data lies to the left of mean

$P(x \leq \mu) = 50\%$ ∴ adding 34% from

μ to $\mu + \sigma$

→ 84%.

∴ $P(x \leq \mu + \sigma)$ is 84%.

✳ Standard Normal Variable (Z score)

A Z score (also called as standard score) tells us how many standard deviations a data point is from the mean of the dataset in a normal distribution. It allows us to standardize different datasets so we can compare them or calculate probabilities.

Key Concepts:

① Position Relative to mean

→ A Z score measures the position of a value x relative to the mean μ of a dataset

- A +ve Z score means X is above mean.
- A -ve Z score means X is below mean
- A Z score of 0 means X is equal to mean

② Scaling by Standard deviation

The Z score considers how far a value is from the mean, scaled by standard deviation (σ) this ensures that the comparison is not biased by spread of data.

③ Standard Normal Distribution

- By converting values to Z score any normal distribution can be transformed into standard normal distribution. (mean=0, standard dev=1)

How Z scores work?

- Let's say we have students dataset

Mean height $M = 170\text{cm}$

standard deviation (σ) = 10cm

If student height is 180cm , Z score tells us how many standard deviations 180 is above the mean

Interpreting Z scores:

$Z = 0$ (exactly at mean)	$ $	$Z = 1$ (1 standard dev above mean)
$Z = -1$ (1 deviation below mean)	$ $	$Z = 2$ (2 dev above mean)

Q Comparing heights

Dataset: heights of men $\rightarrow \mu = 170, \sigma = 10$

A man's height is $x = 180$

$$Z \text{ score} = \frac{\text{Value} - \text{Mean}}{\text{Standard Dev}} = \frac{180 - 170}{10} = 1$$

$\Rightarrow x = 180$ is 1 standard deviation above the mean

if height were $x = 160$

$$Z = \frac{160 - 170}{10} = -1$$

We can refer to Z table for Probability
using Z score

e.g. 1.65 \rightarrow row will be 1.60

and col will be 0.05

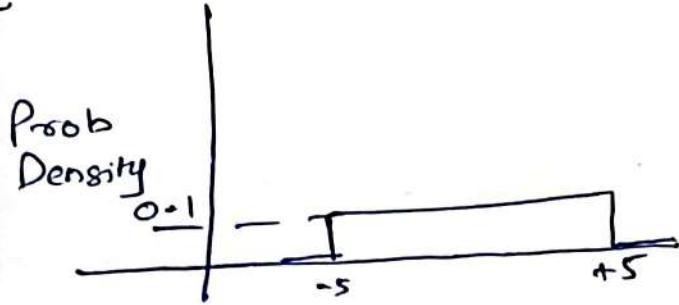
Q regulatory authority selects a random tablet from a batch; based on previous knowledge you know that batch Z2 has mean $\rightarrow 510\text{mg}$ std dev = 20mg

Prob that tablet select by authority has Paracetamol level $< 550\text{mg}$

$$\frac{x - \mu}{\sigma}, \frac{550 - 510}{20}, \Rightarrow Z = 2 \quad (2)$$

from Z table $P(Z < 2) = \underline{0.977}$

Q



this graph represents the PDF of a uniformly distributed random variable X .

As we can see the prob density is equal for all the possible values of x ($-5 \leq x \leq 5$)

What is the prob of the random variable X lying b/w -1.5 to $+2.5$ i.e. $P(-1.5 < x < 2.5)$

Ans: the prob of the variable lying b/w -1.5 and 2.5 would be equal to the area under the PDF

$$\begin{aligned} \rightarrow \text{Calculating the area of rect from } (-2.5 \text{ to } 2.5) \text{ and} \\ \text{of height } (0.1) &\Rightarrow (0.1)(2.5 - (-1.5)) \\ &= 0.1(2.5 + 1.5) = 0.4 \\ &\Rightarrow \underline{\underline{0.4}} \end{aligned}$$

Q

An astronomer's error in estimating the dist b/w earth and Uranus follows a normal distribution with $\mu = 0 \text{ km}$ and standard deviation $\sigma = 1000 \text{ km}$

What's the prob that the error is

- ① overestimating by 2330 km or more
- ② within $\pm 500 \text{ km}$

Soln ① we need to find $P(X \geq 2330)$

$$Z^2 = 2330 - 0 / 1000 = 2.33$$

Now using Z-table

$$P(Z \geq 2.33) = 1 - P(Z \leq 2.33)$$

$$\Rightarrow 1 - (0.9901) = 0.0099 = 0.99\%$$

② $P(-500 \leq X \leq +500)$

$$\text{for } x = -500 \Rightarrow Z = \frac{-500 - 0}{1000} = -0.5$$

$$\text{for } x = +500 \Rightarrow Z = +0.5$$

$$P(Z \leq 0.5) = 0.6915$$

$$P(Z \leq -0.5) = 0.3085$$

$$P(-0.5 \leq Z \leq 0.5) = P(Z < 0.5) - P(Z < -0.5)$$

$$0.6915 - 0.3085$$

$$0.3830 \Rightarrow 38\%$$

* Central limit theorem

Sample : In statistics, a sample is a subset of individuals, items or observations taken from a larger population to study and make conclusion about the population. Studying an entire population is often impractical.

Keypoints : A sample should represent the population well to ensure that inferences made from the samples are valid.

→ Bias in sampling can lead to inaccurate results

(*) Notations and formulae related to population and samples

Population/Sample	Term	Notation	Formula
($x_1, x_2, x_3, \dots, x_n$)	Population Size	N	Nos of items/ele in population
	Population Mean	μ	$\sum_{i=1}^N \frac{x_i}{N}$
	Population Variance	σ^2	$\frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$
Sample ($x_1, x_2, x_3, \dots, x_n$) sample of Population	Sample Size	n	Nos of items/ele in the sample
	Sample mean	\bar{x}	$\frac{\sum_{i=1}^n x_i}{n}$
	Sample Variance	s^2	$\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$

Now suppose we conducted a survey and got that 50.5% of sample are preferring Option 2 over option 1, is it fair to conclude that 50.5% of the people in population will also prefer option 2 over option 1?

6

Write now can't say...

Sample Distribution

A Sampling distribution is the probability distribution of a statistic (like the sample mean, proportion, variance etc) calculated from multiple random samples of the same size drawn from a population.

Eqs and Visualizations

Imagine we have a population of 10 students with the following test scores

Population : 50, 55, 60, 65, 70, 75, 80, 85, 90, 95

$$\hookrightarrow \text{Population Mean}(\bar{M}) \Rightarrow 72.5$$

Step 1 : Take random samples

Suppose you take random samples of size 2 ($n=2$) repeatedly from this population and calculate the Sample mean (\bar{x}) for each sample

→ Here are some possible samples and their means

$$\text{Sample 1 } \{50, 55\} \rightarrow \text{Mean} \Rightarrow (50+55)/2 = 52.5$$

$$\text{Sample 2 } \{65, 70\} \rightarrow \text{Mean} = 67.5$$

$$\text{Sample 3 } \{85, 95\} \rightarrow \text{Mean} = 90$$

$$\text{Sample 4 } \{60, 80\} \rightarrow \text{Mean} = 70$$

Similarly many samples can be withdrawn

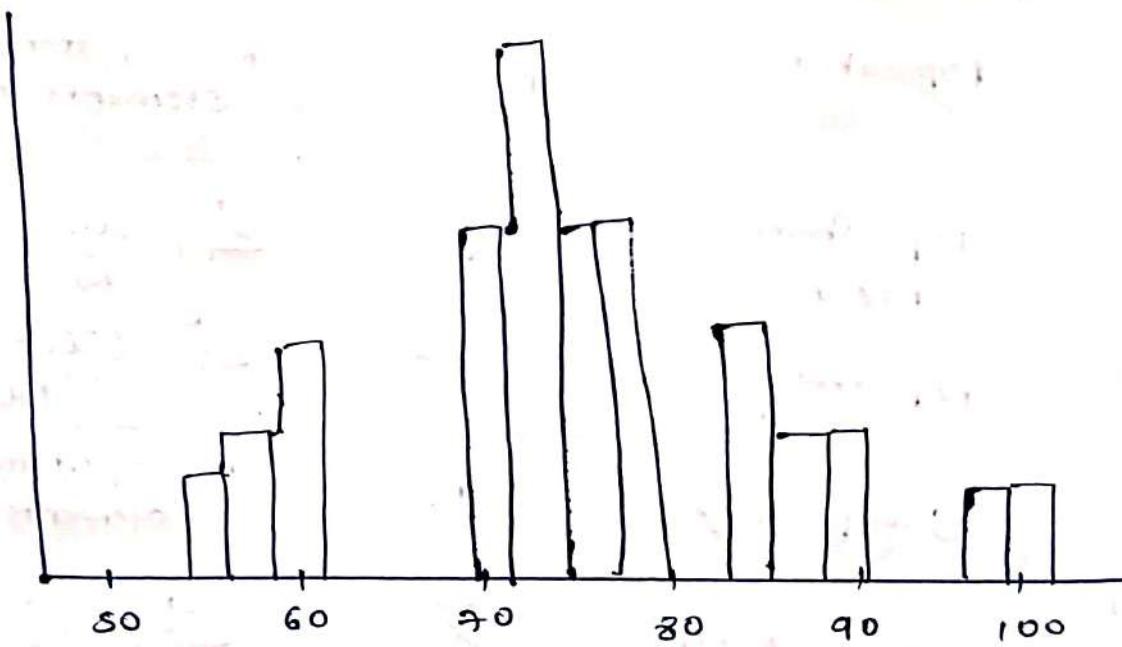
Step2: Create a sample distribution

If we list the mean's from all possible samples of size 2 we get the sampling distribution of the sample mean.

The sampling dist will look like a histogram showing freq of each sample mean.

Visualizations

Sampling Distribution of Sample Mean
(Sample Size = 2)



Observations:

- the sample means are centered around the Population mean ($\mu = 72.5$)
- the distribution is symmetric because population itself is symmetric
- As the sample size increases the spread (variability) of sampling distribution decreases making it more concentrated around mean.

* the standard deviation of the distribution of the sampling means is given by

$$\text{Sampling distribution's Standard Deviation} = \frac{\sigma}{\sqrt{n}}$$

Population SD
Sample size

the higher the value of the standard deviation the fatter the curve will be, hence the values will be farther apart.

✳ Notations and formulae related to Sampling Distribution

Population/ Sample	Term	Notation	Formula
Population ($x_1, x_2, x_3, \dots, x_n$)	Population Size	N	Nos of Items/ elements in population
Sample (x_1, x_2, \dots, x_n) (sample of population)	Population Mean	μ	$\sum_{i=1}^n \frac{x_i}{N}$
	Population Variance	σ^2	$\sum_{i=1}^n \frac{(x_i - \mu)^2}{N}$
Sampling distribution of sample mean ($\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n$) K sample means	Sample Size	n	Nos of Items/ elements in sample
	Sample Mean	\bar{x}	$\sum_{i=1}^n \frac{x_i}{n}$
	Sample Variance	s^2	$\sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1}$
Sampling distribution of sample mean ($\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n$) K sample means	Sampling distribution mean	$\mu_{\bar{x}}$	$\mu_{\bar{x}} = \mu$
	Sampling dist standard dev	$s_{\bar{x}}$ (standard error)	σ / \sqrt{n}

Central Limit theorem (CLT)

It's a fundamental concept in statistics that explains how the sampling distribution of the sample mean approaches a normal distribution as the sample size increases, regardless of the original population distribution.

Keypoints of CLT

- ① Sample Mean approximates population Mean :
the mean of sampling distribution \bar{X} is equivalent to the population mean (μ)
- ② Standard error ($\sigma_{\bar{X}}$) : the standard deviation of the sampling distribution (called as standard error) is given by $\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$

where
 $\sigma \rightarrow$ Standard Deviation of Population
 $n \rightarrow$ sample size
- ③ Normality of large sample sizes : for $n > 30$ the sampling distribution of sample mean becomes approximately normal, even if the population distribution is not normal
- CLT helps us make predictions and performs statistical tests, even if the original data is not normally distributed

Every day examples

If we measure the avg. height of 30 students in different classrooms (multiple samples)

→ their avg heights will

center around population mean (true avg height of all students)

→ form a normal distribution even if individual heights don't follow normal distribution

So for an office of 30000 employees, we wanted to find the avg commute time. So instead of asking all, we can ask only 100. So for sample size of 100 → Mean(\bar{x}) \Rightarrow 36.6 Standard Dev. (S) = 10

However it would not be fair to conclude the same for total population because the flaws of sampling process must have lead to some error. Hence sample mean's value has to be reported with some margin of error.

$$36.6 \pm \text{marginal error}$$

Practice Ques

→ Exit Poll for MCD, we asked 100 randomly selected voters to name the party they voted for

Data Collected is

Contesting Party	No of voters
BJP	58
INC	42

from this sample we have to estimate the % of voters that might have voted for BJP
So we defined X as the proportion of people that voted for BJP

X	freqv	frequency Distribution
1	58	of X
0	42	

Mean of X which is equal to $(0+0+0\dots+42+58)$
and $(1+1+1\dots 58 \text{ times})$

$$\text{Mean} = \frac{(0 \times 42) + (1 \times 58)}{100} = 0.58$$

Means in terms of X (Voted for BJP)

Standard Deviation (C_x) $\sqrt{\frac{(x_i - \bar{x})^2}{n-1}}$

$$\Rightarrow \frac{((0 - 0.58)^2 \times 42) + ((1 - 0.58)^2 \times 58)}{100-1}$$

$$\Rightarrow 0.2461$$

Q We define Y (Proportion of people voted for Congress)

$$\text{Mean} = 0.42$$

	Y	Votes
Freq	0	58
Distribution	1	42

$$\sigma^2 = \frac{(0 - 0.42)^2 \times 58 + (1 - 0.42)^2 \times 42}{100-1}$$

$$\Rightarrow \frac{(-0.42)^2 \times 58 + (0.58)^2 \times 42}{99}$$

$$\Rightarrow \frac{0.1864 \times 58 + 0.3364 \times 42}{99}$$

$$\Rightarrow \frac{10.281 + 14.128}{99} = 0.2460$$

$$\Rightarrow SD = \sqrt{\sigma^2}, \sqrt{0.2460} = 0.4959$$

$$\Rightarrow 49.59\% = 49.6\%$$

Q3 Voter Sample

Let's say we have a sample distribution of \bar{Y} the proportion of people that voted for 'inc'. Mean of sample distribution is $M\bar{X} = 0.50$

$SE(\sigma/\sqrt{n})$ is 0.048

To find σ for population

$$0.048 = \frac{\sigma}{\sqrt{n}} \quad n = 100$$

$$\sigma = 0.048 \times \sqrt{100} = 0.48 \text{ or } 48\%$$

* Estimating Mean Using CLT

Earlier we tried to estimate the mean commute time of 30,000 employees of an office by taking a small sample of 100 employees and finding their mean commute time.

The sample mean was $\bar{X} = 36.6$ minutes & standard deviation was $S = 10$ minutes.

Recall that we also said about the population mean ie daily commute time of all 30,000 employees $M = 36.6$ (sample mean) \pm some margin of error

We can find the margin of error using CLT.

→ To find the margin of error using the central limit theorem (CLT) let's go step by step

Given Data:

- ① Sample mean (\bar{x}) = 36.6 minutes
- ② Sample standard deviation (s) = 10 minutes
- ③ Sample Size (n) = 100
- ④ The population mean (μ), is approximately

$$\bar{x} \pm \text{margin of error}$$

(Step 1): Standard error of Mean

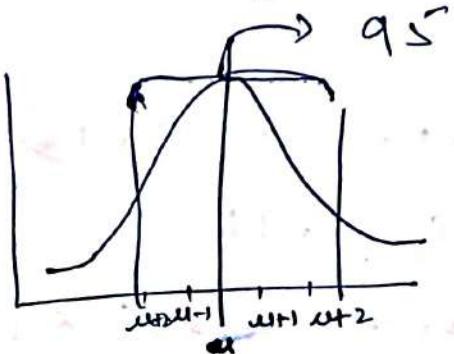
By Central limit theorem

$$SE = \frac{s}{\sqrt{n}} = \frac{10}{\sqrt{100}} = 1$$

∴ Sample size $100 > 30$ so the sampling distribution is a normal distribution.

But the Mean is unknown.

Let's see $P(\mu - 2 < \bar{x} < \mu + 2) = ?$



$\Rightarrow 0.954$ (by normal distribution)
Since 68, 95, 99.7

$$P(36.6 - 2 < \mu < 36.6 + 2) = P(\mu - 2 < \bar{x} < \mu + 2) \\ = 95\%$$

$$P(36.6 - 2 < \mu < 36.6 + 2) = 95.4\%$$

Prob. associated with the claim is called confidence level (Here it's 95.4%).

Maximum error made in sample means is called Marginal error. (Here it is 2 min.)

→ Final interval of values is called confidence interval (34.6, 38.6)

Step by step approach using Z score

Given Data

① Sample mean (\bar{x}) = 36.6 minutes

② Sample standard deviation (s) = 10 minutes

③ Sample size (n) = 100

→ SE (standard error) = $\frac{s}{\sqrt{n}} = \frac{10}{\sqrt{100}} = 1$

$$\Rightarrow SE = 1$$

Marginal error is calculated as

$$\text{Margin of Error} = Z \cdot SE$$

∴ Sample data is normally distributed

Based on CLT

$$P(\mu - 2 < \mu < \mu + 2) = 0.95$$

Z for 95% value is 1.96

$$\Rightarrow \text{Margin of error} = 1.96 \times 1 = \underline{1.96}$$

$$\vec{M} = \bar{x} \pm \text{margin of error}$$

$$\Rightarrow 36.6 \pm 1.96 \Rightarrow \underline{\underline{(34.64 - 38.56)}}$$

* Now we can say that for $\% \text{ confidence level}$
the confidence interval is given by

$$\left(\bar{x} - z \times \frac{s}{\sqrt{n}}, \bar{x} + z \times \frac{s}{\sqrt{n}} \right)$$

$\hookrightarrow z * (SE)$
 $\hookrightarrow \sigma_{\bar{x}} \text{ or } \frac{s}{\sqrt{n}}$ for sample

Q Maggie Noodles example

↳ Max Permissible lead in products 2.5 PPM

↳ Sample Size = 100

Sample Mean (\bar{x})
for lead content = 2.3 PPM

Standard Dev (S) = 0.3 PPM

Now to find the confidence Interval

↳ Confidence Level \Rightarrow It's a sensitive task

let's take it as $99\% \quad z = 2.576$

\Rightarrow Confidence Interval $\Rightarrow 2.3 \pm (2.576)(0.3) \sqrt{100}$

$$\Rightarrow 2.3 \pm \frac{2.576(0.3)}{10}$$

$$(2.223, 2.387) \checkmark$$

Q Amt of Paracetamol specified by drug regulatory authority \rightarrow 500mg
 \rightarrow allowed errors $\Rightarrow 10\%$.

Below 450mg \rightarrow Quality Issue & above 550mg would be serious regulatory issue

- \rightarrow Given sample of 100 tablets
- \rightarrow Mean Paracetamol $\Rightarrow 530\text{mg}$
- \rightarrow Standard dev: 100mg
- \rightarrow CI = 95%.

$$\text{CI} : \mu \pm Z \cdot SE$$

Z score for 0.95 (95%)

Note: 95% of confidence level implies that 95% of the data lies within the confidence interval

This leaves 2.5% on each tail split equally
2.5% on the lower side and 2.5% on the upper side.

To determine the upper bound Z score, we need the cumulative prob until that point
(ie: the central tendency (95%)
+ the lower tail (2.5%).)

Hence the cumulative probability becomes

$$95\% + 2.5\% = 0.975$$

Hence Z score for 0.975 = 1.96

So margin of error for 95% confidence level

$$Z \times SE \Rightarrow (1.96) \times \frac{\sigma}{\sqrt{n}} : \frac{1.96 \times 100}{\sqrt{100}}$$

$$\frac{1.96 \times 100}{10} = 19.6$$

→ Confidence Interval for 95% confidence level

$$530 \pm 19.6 \Rightarrow (510.4, 549.6)$$

✳️ Confidence Interval for 90% confidence level

90% confidence level means 5% left on each tail

→ Z for 90% + 5% → for 0.95 is 1.64

$$\therefore Z \times SE = (1.64) \left(\frac{100}{\sqrt{100}} \right) : \frac{1.64 \times 100}{10} \\ 16.4$$

$$530 \pm 16.4 \quad (513.6, 546.4)$$

Hypothesis Testing

Let's understand the basic difference b/w inferential statistics and hypothesis testing.

Inferential Statistics

Purpose : to make generalizations about a population based on data collected from a sample.

Key concept : The population mean (or another parameter) is unknown, so we estimate it using :

→ Sample statistics : such as sample mean (\bar{x})

→ Confidence Intervals : These provide a range within which the population parameter is likely to fall, along with confidence level (e.g. 95% confidence interval)

Use case : When you don't have prior knowledge about the population and need to draw conclusions from it

Hypothesis Testing

→ Purpose : To test an assumption or claim (hypothesis) about a population parameter using sample data. Hypothesis testing is used to confirm your conclusions (or hypothesis) about the population parameter (which you know from EDA or your intuition). Through hypothesis testing, you can determine whether there is enough evidence to conclude if the hypothesis about the population parameter is true or not.

To understand about hypothesis, let's take an ex " A criminal trial → A person has been charged with crime, Jury has to decide whether the defendant is innocent or guilty.

Hypothesis 1 : Defendant is Innocent

Hypothesis 2 : Defendant is not Innocent (Guilty)

These two opposing hypothesis are called as the "null hypothesis" and the "alternate hypothesis"

eg of food product containing excess lead

⇒ Null Hypothesis

Avg lead content is within the allowed limit
of 2.5 PPM

⇒ Alternate hypothesis

Avg lead content is more than 2.5 PPM

In our eg of Criminal trial

Null hypothesis (H_0) \rightarrow Defendant is ^{Notation}innocent

Alternate hypothesis (H_1) \rightarrow Defendant is not innocent

If the Defendant is found guilty \Rightarrow there is a "Rejection of Null Hypothesis"

Key concepts:

→ you start with null hypothesis (H_0) \rightarrow a statement of no effect or no difference.

→ the alternate hypothesis (H_1) is what you aim to support

through hypothesis testing, you can check whether the sample data provides enough evidence to reject H_0 in favour of H_1

In Maggie Noodles eq

↳ Null Hypothesis (H_0) \Rightarrow avg lead content is less than or equal to 2.5 PPM.

If we fail to reject null hypothesis, we can conclude that Maggie Noodles Doesn't contain excess lead.

Note: In hypothesis testing we never "accept" the null hypothesis we can only fail to reject it because

\rightarrow the null hypothesis (H_0) represents a default assumption such as there is no effect or the means are equal.

\rightarrow failing to reject H_0 doesn't prove it is true \rightarrow It simply means the sample data does not provide proper evidence against H_0

\rightarrow Lack of evidence is not equivalent to evidence of absence.

Just like the court of law, failing to reject doesn't mean Innocence (acceptance) It simply means there wasn't enough evidence to prove guilty (reject H_0)

Both Null and alternate hypothesis can't be true at the same time, Only one of them is true

Key Rule for formulating Hypothesis

Null Hypothesis (H_0) :

- Represents a statement of no difference, no effect or a baseline assumption.
- Always include the signs :
 $=, \leq$ or \geq
 - ↳ Because hypothesis testing seeks evidence to reject H_0 and these signs allows us to test for deviations.

Alternative hypothesis (H_1) :

- Represents what we aim to support or find the evidence for.
- Always include the signs :
 $\neq, >$ or $<$
 - ↳ Because H_1 is formulated to oppose H_0 .

Guidelines for the Claims

- Claims with "at least"
eg. the mean income is at least \$50,000.
 - $H_0 : \mu \geq 50,000$
 - $H_1 : \mu < 50,000$

→ claim with "atmost"
eq defect rate is atmost 5%.

$$H_0 : P \leq 0.05$$

$$H_1 : P > 0.05$$

→ claim with "greater than"
eq the new drug is more effective than old one

$$H_0 : M_{\text{new}} \leq M_{\text{old}}$$

$$H_1 : M_{\text{new}} > M_{\text{old}}$$

→ claim with equals
eq Mean life span of product is 500 hrs

$$H_0 : M = 500$$

$$H_1 : M \neq 500$$

For eg Flipkart claimed that its total valuation in December 2016 was atleast \$ 14 billion. Here the claim contains \geq sign, so the null hypothesis is the original claim.

Total valuation $\geq \$14 \text{ Billion} \rightarrow \text{Null Hypothesis}$

Total valuation $< \$14 \text{ Billion} \rightarrow \text{Alternate Hypothesis}$

Making a decision

e.g.: Apurva's claim : Avg score in archery
is equal to 70.

Over 5 games of Archery :

- ↳ Avg score of Apurva = 20 (Less likely to believe her claim)
- ↳ Avg score = 65 → More likely to believe her claim.

So we reject her claim when her avg score was 20 but fail to reject if her avg score is 65 → so what score b/w 20 and 65 is the boundary where we decide to change our opinion.

↳ this point is called as critical point.

④ Critical Region : It's like a danger zone for the null hypothesis (H_0) In hypothesis testing, it's a range of values where, If your test statistic lands → we decide to reject the null hypothesis.

think of it in this way :

- Imagine throwing a dart at the target
- null hypothesis says that the dart will land on the safe zone (outside the critical region)

- the critical region is like a danger zone on the target.
- if your dart (test statistic) hits the red zone, it's an unexpected (based on null hypothesis) that you decide, "the null hypothesis doesn't make any sense"

④ Critical Point : It's the boundary value that separates the critical region from the non-critical region. It is determined based on the significance level (α) and the type of test (e.g. one-tailed or two-tailed)

How it works:

→ for a significance level of 5% ($\alpha = 0.05$)

↳ One-tailed test: the critical point is the value of the test statistic that leaves 5% in one tail of distribution.

↳ Two-tailed test: the critical points leave 2.5% in each tail

④ Tailed Tests

- Determines how the rejection region is distributed
- One-tailed test: tests whether the sample statistic is either significantly greater than or significantly less than the hypothesized value
- Two-tailed test: tests whether the sample statistic is significantly different (either greater or less) than the hypothesized value.

Eg's with visualizations

Eg1 One-tailed test (Right tail)

Scenario: A factory claims that its machine produces bolts with an avg strength of at least 50Kg. We want to test if the avg strength is greater than 50Kg.

1. Hypothesis

$$\Rightarrow H_0 : \mu \leq 50$$

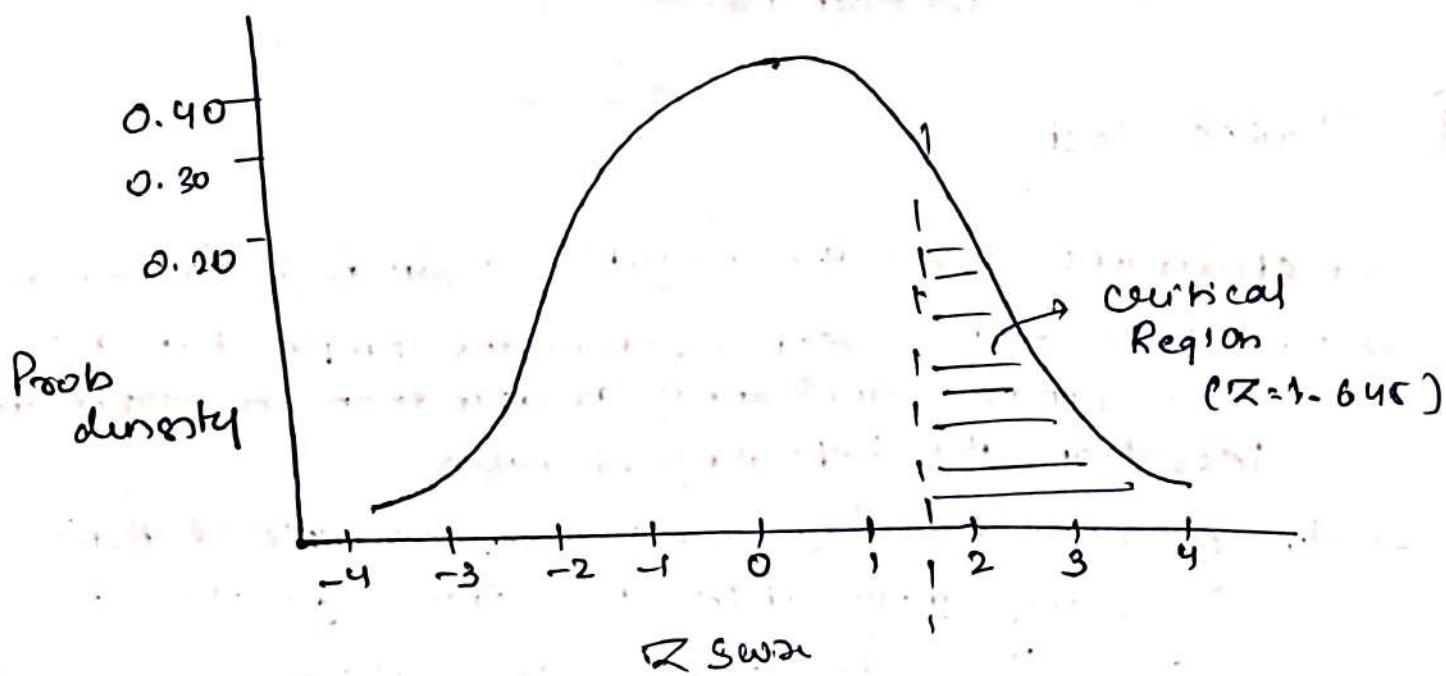
$$H_1 : \mu > 50 \text{ (Right tailed test)}$$

2. Significance level (α) = 0.05 \rightarrow 5%, the critical region lies in the right tail of distribution

3. Critical point: For $\alpha = 0.05$, the critical Z value is 1.645

If the test statistic $Z > 1.645 \rightarrow \text{reject } H_0$

4. Visualization



- Here the shaded area represents the critical region for a right-tailed test ($H_1: \mu > 50$)
- the critical point $Z = 1.645$ separates the critical region from the rest of the distribution.
- If the test statistic falls in shaded region \rightarrow we reject H_0 .

(eq2) two-tailed test

Scenario: A pharma company claims that the avg effectiveness of drug is 75%. We want to test if the actual effectiveness is diff from 75%.

1. Hypothesis:

$$H_0: \mu = 75\%$$

$$H_1: \mu \neq 75\% \text{ (two tailed test)}$$

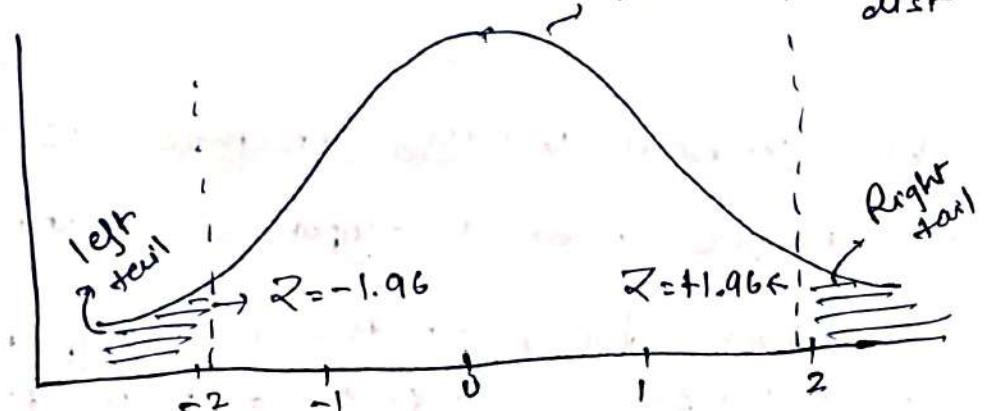
2. Significance Level (α) = 5% \rightarrow Critical region is split b/w two tails of the distribution.

3. Critical points: For $\alpha = 0.05$ $Z = -1.96, 1.96$

↳ If the test statistic

$$Z < -1.96 \text{ or } Z > 1.96 \Rightarrow \text{reject } H_0$$

Visualization



* Quick rule to identify the type of test

① Two-tailed test (\neq)

- Use when you are testing for a difference (each direction).
- Rejection regions are on the both tails of the distribution
- eg testing if mean is not equal to 50 $\Rightarrow H_1: \bar{M} \neq 50$

② Lower-tailed test ($<$)

- Use when you are testing if the population parameter is less than a certain value.
- Rejection Region is on the left side of the distribution

eg testing if mean is < than 50 $\Rightarrow H_1: \bar{M} < 50$

③ Upper-tailed test ($>$)

- Use when you are testing if the population parameter is greater than a certain value
- Rejection Region is on the right side of distribution

$$H_0: \bar{M} > 50$$

* Critical Value Method

The critical value method is a statistical technique used in hypothesis testing to decide whether to reject the null hypothesis (H_0) based on a comparison of test statistic to critical values.

Steps in the Critical Value method:

1. formulate the hypothesis

① $H_0 \Rightarrow$ Null Hypothesis \Rightarrow Assumes no effect or no difference (eq $\mu = \mu_0$)

② $H_1 \Rightarrow$ Alternative Hypothesis \Rightarrow indicates the presence of an effect or difference eq ($\mu \neq \mu_0$, $\mu > \mu_0$ or $\mu < \mu_0$)

2. choose significance level α

(commonly used α values 0.05 (5%) or 0.01 (1%)

3. Determine the critical value:

① Critical values (Ucv and Lcv) are the boundaries of the rejection region for H_0

② use the Z table (for normal distributions) to find these values based on α and whether the test is two tailed or one tailed.

4. Compute test statistic

$$\text{for a Z test} \Rightarrow Z = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}}$$

$\bar{x} \rightarrow$ Sample mean

$\mu_0 \rightarrow$ Hypothesized population mean

$\sigma \rightarrow$ Population standard deviation

$n \rightarrow$ Sample size

5. Compare the test statistic with critical value

→ For a two tailed test reject H_0 if
 $Z < LCV$ (lower critical value) or $Z > UCV$
(upper critical value)

→ For One tailed test, reject H_0 if Z is beyond the critical value in appropriate tail.

6. Make a decision

→ If the test statistic falls in the rejection region reject H_0
→ Otherwise fail to reject H_0

Ex Company claims that avg weight of product is 500gm, you suspect that it's incorrect and test this with a sample of 50 products, the sample mean weight $\Rightarrow 495$ gms, standard dev $\Rightarrow 10$ gms
 $\alpha = 0.05$

$$H_0 \Rightarrow \mu = 500, H_1 \Rightarrow \mu \neq 500$$

$$\alpha = 0.05 \quad (\text{two tailed test}), \quad \alpha = 0.025 \text{ in each tail}$$

Using Z table Z value for (0.025)

$$LCV = -1.96, UCV = +1.96$$

Compute test statistic $\Rightarrow \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} = \frac{495 - 500}{10/\sqrt{50}}$
 $\Rightarrow -3.55$

Z value (-3.55) < -1.96 It falls in rejection region
∴ Reject H_0 → Avg weight is not 500g.

Q What would be the value of cumulative prob of UCV of the significance level (α) for an upper tailed test is 3%?

In an upper-tailed test with a significance level (α) of 3%, the cumulative prob up to the Upper critical value is (UCV).

This is because the rejection region (critical region) is the area on the upper tail of the distribution corresponding to significance level.

Calculation: cumulative prob? $1 - \alpha \Rightarrow 1 - 0.03 = 0.97$

→ the cumulative probability of upper critical value is 0.97 → 97%.

Critical Value Method for Hypothesis testing

The Critical value method involve comparing a test statistic to a critical values to decide whether to reject or fail to reject the null hypothesis , let's break the process into steps :

Step 1 : formulate the hypothesis

① Null Hypothesis (H_0) : this is the default assumption we are testing, for eg : H_0 ; the mean salary of employee is \$50000

② Alternative Hypothesis (H_1) : this is the opposing statement we aim to support

for eg : the mean salary of employees is not \$50,000 (two tailed test)

Step 2 : Determine the Significance Level (α)

→ Significance level (α) : this is the prob of rejecting the null hypothesis

if not specified assume $\alpha = 0.05 \Rightarrow 5\%$

for $\alpha = 0.05$ in a two tailed test
 $\alpha/2 = 0.025$ in each

Step 3 : Calculate the Z value (Z_c)

→ the Z value is the critical value corresponding to $\alpha/2$, it can be found using Z table

for $\alpha = 0.05$

$Z_c = 1.96$ (from the Z table for 0.025 in each tail for a standard normal distribution).

Step 4 : Determine the critical values (UCV and LCV)

① Critical values

→ upper critical value (UCV) : $M + Z_c \cdot \frac{\sigma}{\sqrt{n}}$

→ lower critical value (LCV) : $M - Z_c \cdot \frac{\sigma}{\sqrt{n}}$

example:

→ Assume :

$M = 50,000$ (hypothesized population mean)

$\sigma = 4000$ (population standard deviation)

$n = 25$ (sample size)

→ standard error (CSE) = $\frac{4000}{\sqrt{25}} = \frac{4000}{5} = 800$

→ Critical values:

UCV → $M + Z(CSE) = 50000 + 1.96(800)$

→ LCV : $M - Z(CSE) = 50000 - 1.96(800)$

⇒ (48,432, 51,568)

Step 5: Compute test statistic

→ the test statistic measures how far the sample mean (\bar{x}) is from the hypothesized mean (μ)

$$\text{formula} = Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

Step 6: Make a decision

① Decision Rule:

→ Reject H_0 if \bar{x} or (Z) lies outside the critical region defined by LCV and UCV

→ fail to reject H_0 if \bar{x} or (Z) lies inside the critical region

e.g. Sample mean (\bar{x}), Assume $\bar{x} = \$1000$

$$\text{test statistic } (Z) = \frac{\$1000 - \$1000}{\$000} \\ \Rightarrow 1.25$$

∴ 1.25 is b/w -1.96 to +1.96

↳ Fail to reject H_0

Q Problem Statement:

Company claims that the avg time taken to resolve the customer query is 15 min, A random sample of 36 queries had a mean resolution time of 16 min with a standard dev of 3 min, $\alpha = 0.05$

Step 1 Formulating Hypothesis

H_0 (Avg resolution time is) $15\text{min} \Rightarrow \bar{M} = 15$

H_1 (Alternate hypothesis) $\bar{M} \neq 15$ (two tailed)

Given $\alpha = 5\%$.

→ ②nd step ∵ it's two tailed $\Rightarrow \alpha/2 = \frac{0.05}{2} = 0.025$

3rd step ~~Graph~~

Z critical value depends on the significance level (α) and the type of test (one tailed or two tailed)

Below is the guide for common significance level

for One tailed test, the critical value is found by

$$1 - \alpha$$

α (Significance Level)	Cumulative Prob $(1 - \alpha)$	Z critical value (Z_c)
0.10 (10%)	0.90	1.28
0.05 (5%)	0.95	1.645
0.01 (1%)	0.99	2.33

For two tailed test, divide the significance level by 2 ($\alpha/2$) and then find the critical value for $1 - \alpha/2$

α (Significance Level)	Cumulative Prob $(1 - \alpha/2)$	Z critical value
0.10 (10%)	0.95	± 1.645
0.05 (5%)	0.975	± 1.96
0.01 (1%)	0.995	± 2.575

Hence $\alpha = 0.05$ \Rightarrow It's a two tailed test

$$\Rightarrow \alpha/2 = 0.025$$

$$\text{for } Z_{\text{critical}} \text{ or } Z_C = 1 - \alpha/2 = 1 - 0.025 \\ = 0.975$$

$$\Rightarrow Z(0.975) = \pm 1.96$$

$$\Rightarrow UCV = +1.96, LCV = -1.96$$

(Step 4) To compute test statistic

$$\bar{x} = 16, \mu = 15 \quad | \quad SE = \sigma/\sqrt{n} = 3/\sqrt{36} = \frac{1}{2} = 0.5$$
$$\sigma = 3, n = 36$$

$$Z \text{ test statistic: } Z = \frac{\bar{x} - \mu}{SE} = \frac{16 - 15}{0.5} = \frac{1}{0.5} = 2$$

(Step 5) Decision Making

$$Z \text{ test} = 2 \quad ; \quad Z \text{ test} > Z_{\text{critical}}$$

$$Z_{\text{critical}} \Rightarrow \pm 1.96 \quad \text{By Reject } H_0$$

Q A manufacturer claims that avg. life of product is 36 months. Auditor selects a sample of 49 items \rightarrow Avg. life for these 34.5 months. Population standard dev = 2 months, $\alpha = 3\%$.

$$H_0: \mu = 36$$

$$H_1: \mu \neq 36 \quad (\text{two tailed})$$

$$\text{Given: } \alpha = 0.03 \Rightarrow Z_C = 1 - \alpha/2 = 1 - 0.015 \\ = 0.985$$

$$\Rightarrow Z \text{ score (0.985)} = \pm 2.14$$

Step 2 : To find UCV, LCV from Z_c

$$\text{Critical Value} = \bar{U} \pm (Z_c \cdot SE)$$

$$36 \pm (2.77 \times \frac{4}{\sqrt{49}}) = 36 \pm (2.77 \times \frac{4}{7})$$

$$\Rightarrow 36 \pm 1.24 \Rightarrow (34.76, 37.24)$$

$$\Rightarrow UCV \Rightarrow 37.24 \quad \bar{x} = 34.5 \\ LCV \Rightarrow 34.76$$

$$\bar{x} < LCV \Rightarrow \text{Reject } H_0$$

Q Govt regulatory body have specified the max permissible amount of lead in any food product is 2.5 PPM

Suppose we take 100 random samples, mean lead content in samples $\Rightarrow 2.6 \text{ PPM}$ with standard dev of 0.6

$$\Rightarrow SE = 0.6$$

One thing to notice here is that the standard dev of sample is given as 0.6 instead of population standard dev \rightarrow so we can approximate the population standard dev to sample standard dev

$$\text{Given } \alpha = 0.03 \Rightarrow 3\%$$

H_0 : Avg lead content $\leq 2.5 \text{ PPM}$

H_1 : Avg lead content $> 2.5 \text{ PPM}$ (Right tailed test)

$$\alpha = 0.03 \Rightarrow Z_c \text{ for } 1 - 0.03 = 0.97 \Rightarrow 1.89$$

Q Cadbury states that avg weight of one of its chocolate Product 'Dairy Milk Silk' is 60g, $\alpha = 2\% = 0.02$

Sample size = 100, $\bar{x} = 62.69$

SE = 10.79

$$Z_c = 1 - \alpha/2 = 1 - \frac{0.02}{2} \cdot 1 - 0.01 = 0.99$$

Z_c value of 0.99 is 2.33

$$\mu \pm Z SE \Rightarrow 60 \pm (2.33) \frac{10.79}{\sqrt{100}}$$

$$60 \pm 2.33(10.79) \Rightarrow +62.49 \\ +60.00 - (2.33)(10.79) \\ (57.51, 62.49) \Rightarrow 57.51$$

Sample mean lies outside the range \rightarrow Failed test

(*) Pvalue Method

What is Pvalue?

The p-value is the probability of observing the data (or something more extreme) assuming the null hypothesis (H_0) is true.

\rightarrow Null hypothesis (H_0): this is the default assumption
eg there is no difference, no effect
or no relationship

→ Alternative Hypothesis (H_1) : the hypothesis you want to test e.g there is a difference, an effect or a relationship.

A low p-value (usually $< \alpha$, where $\alpha = 0.05$ or some significance level) suggests evidence against H_0 meaning you reject the null hypothesis

A high P-value suggests that the data is consistent with H_0 , meaning you fail to reject the null hypothesis

Simplified steps for the Pvalue test method

① formulate the Hypothesis

e.g H_0 : the mean weight of apples is 150gms
 H_1 : the mean weight of apples is not 150gms
(two tailed test)

② Collect Data and Compute the test statistic

formula of Z score $\Rightarrow Z = \frac{\text{Sample mean} - \text{Pop mean}}{\text{Standard error}}$
for test statistic

③ find the Pvalue

use the Z score and a Z table to find the cumulative prob.

④ Compare the pvalue with α

→ If $P < \alpha$, reject H_0

→ If $P \geq \alpha$, fail to reject H_0

E9 One tailed test

Scenario : A teacher claims that the avg score of her students is greater than 75

$$H_0 : \text{Avg score} \leq 75$$

$$H_1 : \text{Avg score} > 75$$

Given Data : \rightarrow Sample mean = 78

$$\rightarrow \text{Population mean} = 75$$

$$\rightarrow \text{Standard Dev} = 10$$

$$\rightarrow \text{Sample size} = 30$$

$$\text{Step 1 : Calculate Z score } \rightarrow Z = \frac{78 - 75}{10/\sqrt{30}} = 1.64$$

Step 2 : find the P value

from the Z table the cumulative prob for $Z = 1.64$

$$\text{is } 0.9495$$

$$\boxed{\text{for one tailed test } \Rightarrow P = 1 - 0.9495 = 0.0505}$$

Step 3 : Decision

$$\text{if } \alpha = 0.05$$

$$P = 0.0505 > 0.05 \Rightarrow \text{fail to reject } H_0$$

Eg two tailed test

Scenario : A company claims that avg battery life of Product is 10 hours

$$H_0 : \text{Avg battery life} = 10 \text{ hrs.}$$

$$H_1 : \text{Avg battery life} \neq 10 \text{ hrs}$$

Given Data:

Sample mean = 9.5
Population mean = 10
Standard Dev = 1.5
sample size = 40
 $\alpha = 0.05$

Calculate Z score

$$Z = \frac{9.5 - 10}{1.5/\sqrt{40}} = \frac{-0.5}{0.237} = -2.11$$

finding the P value based on Z value

So the cumulative prob for $Z = -2.11$ is 0.0174

for two tailed test $\Rightarrow P = 2 \times 0.0174 = 0.0348$

Now taking decision

$P = 0.0348 < 0.05 \Rightarrow P < \alpha \rightarrow$ rejecting H_0

So In Short

After formulating the null and alternate hypothesis the steps to follow in order to make a decision using the p-value method:

- ① calculate the value of Z score for the sample mean point on the distribution
- ② calculate the P value from the cumulative probability for the given Z score using Z table
- # → ③ Make a decision on the basis of P value (multiply it by 2 for two tailed test) with respect to the given value of α (significance level)

Q A manufacturer claims that the avg life of its product is 36 months, we select a sample of 49 units of the product and calculates the avg life to be 34.5 months, the population standard dev is 4 months to test the manufacturer claim at 3% significance level using p-value test

$$\Rightarrow H_0 \Rightarrow \mu = 36 \text{ months}, H_1 \Rightarrow \mu \neq 36 \text{ months}$$

Step1: To calculate value of Z score for sample mean point on the distribution

Calculating Z-score for sample mean (\bar{x}) = 34.5

$$Z = \frac{\text{Sample mean} - \text{pop. mean}}{\sqrt{\frac{4}{49}}}$$

$$\Rightarrow Z = \frac{34.5 - 36}{\sqrt{\frac{4}{49}}} = \frac{-1.5 \times 7}{40} = \frac{-21}{40} = -2.625$$

Step2: Calculating P value for $Z = -2.625$

Z value for $-2.625 \rightarrow 0.0044$

\therefore Its two tailed $\rightarrow 2 \times 0.0044 = 0.0088$

$P = 1 - 0.0088 \Rightarrow$ we don't have to use this here

Note : the calculation process for p-values can depend on which side of distribution we are dealing with, and whether it's one tailed or two tailed, let's see why we didn't use $1 - P$ here.

for a Z score of $\boxed{-2.62}$ the cumulative prob from the Z table ($P = 0.0044$) already represents the area to the left of Z score on the standard normal curve

- Since the Z score is -ve, the area we need is already represented directly as 0.0044 from Z table
- For a two tailed test, we simply double the prob to account for both tails

When to use $1 - P$

We use $1 - P$ when the observed Z-score lies on the right hand side (Positive side of mean), the cumulative prob from Z-table will give the area to the left and we subtract it from 1 to find the area to right.

Hence

① If Z score is -ve (-2.62)

Use cumulative prob directly from

Z-table ($P = 0.0044$) & for 2 tailed $\Rightarrow 2 \times 0.0044$

② If Z score +ve ; find cumulative prob from Z-table $P(0.9956) \Rightarrow$ for 1 tailed $\Rightarrow 1 - P = 1 - 0.9956 = 0.0044$
For two tailed $\Rightarrow 2(1 - 0.9956) \Rightarrow 0.0088$

(+) Hypothesis testing for Paracetamol Content

Scenario: A pharma company produces tablet with 500mg of paracetamol as the ideal content.

- Quality Issues of Paracetamol $< 500\text{mg}$
- Regulatory Issues of Paracetamol $> 500\text{mg}$

Objective to test whether the manufacturing process is running successfully using a sample of 900 tablets

Given Data
→ Sample size (n) = 900
→ Sample mean (\bar{x}) = 510mg

→ Sample standard dev = 110mg (s)
→ $\alpha = 5\%$.

Steps:
 $H_0 \Rightarrow \mu = 500\text{mg}$
 $H_1 \Rightarrow \mu \neq 500\text{mg}$

$$Z = \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{510 - 500}{110/\sqrt{900}} = \frac{10}{3.67} = 2.72$$

for $Z = 2.72 \Rightarrow \text{Prob} = 0.9967$

∴ two failed $\Rightarrow P = 2 \times (1 - 0.9967) = 0.0066$

? $P < \alpha$

∴ Reject H_0

Types of Errors

In hypothesis testing, two types of errors can occur based on whether we reject or fail to reject the null hypothesis.

These are :

① Type 1 Error (false Positive) :

→ Definition : Rejecting the null hypothesis H_0 when it is actually true.

→ Cause : Happens due to random sampling variation or setting a very low significance level (α)

→ Probability : The significance level α represents the likelihood of making a type I error.

For e.g. $\alpha = 0.05$ means there is a 5% chance of rejecting H_0 incorrectly.

e.g. : We conclude that a new drug works better than the existing ones, but in reality it doesn't

② Type 2 Error (false negative)

→ Definition : Failing to reject the null hypothesis (H_0) when it is actually false

→ Cause : Happens due to small sample size, large variability in data or setting a high significance level (α)

→ Probability : Represented by β (Power of test = $1 - \beta$)

Eg of type 2 error : we conclude that a new drug is no better than existing one, but in reality it is better

Comparison of Errors :

	Type I error	Type II error
Null Hypothesis (H_0)	Rejected Incorrectly	Not rejected when it should be rejected
Outcome	False Positive	False negative
Risk Control	Controlled by setting significance level (α)	Reduced by increasing sample size or test power

→ type I error represented by α occurs when you reject a true null hypothesis.

→ type II error represented by β occurs when you fail to reject null hypothesis when it's actually false

Q P-value method

Suppose we conduct a hypothesis test and observe that the value of sample mean and sample standard deviation when $n=25$ do not lead to the rejection of the null hypothesis, we calculate the p value as 0.0667, what would happen to P value

If you observe the same sample mean and standard deviation for a larger sample size say 750

Ans When the sample size (n) increases while keeping the sample mean and standard deviation constant, the P value will decrease because:

Impact of Larger Sample size on p value

① Formula for test statistic (Z)

$$Z = \frac{\bar{x} - \mu}{S/\sqrt{n}}$$

(S) → Sample Standard Deviation
(n) → Sample Size

as $n \uparrow \Rightarrow S/\sqrt{n} \downarrow \Rightarrow Z \uparrow$ (Moving further away from 0)

So the P-value Relationship

↳ the p-value is the probability of observing a value as extreme as the test statistic, assuming the null hypothesis is true.

A larger Z score corresponds to a smaller area under the tail of curve, hence a smaller p value.

Q In a scenario, the Null Hypothesis (H_0) is that the person does not have the disease, and the alternate hypothesis $H_1 \rightarrow$ Person have a disease, if an treatment doesn't have a serious side effect it is better to increase the probability of making a type I error to avoid missing a critical diagnosis.

T Distribution

The T test and t-distribution test are the statistical tools used to make conclusions about the population parameters, particularly when the sample size is small or the population variance is unknown.

① What is a t-test?

A t-test is a hypothesis test that helps to determine whether there is a significant difference between the means of two groups or between a sample mean and a known value. It uses the t-distribution to calculate the test statistic when:

- the sample size is small ($n < 30$)
- the population standard deviation (σ) is unknown

Types of t-tests:

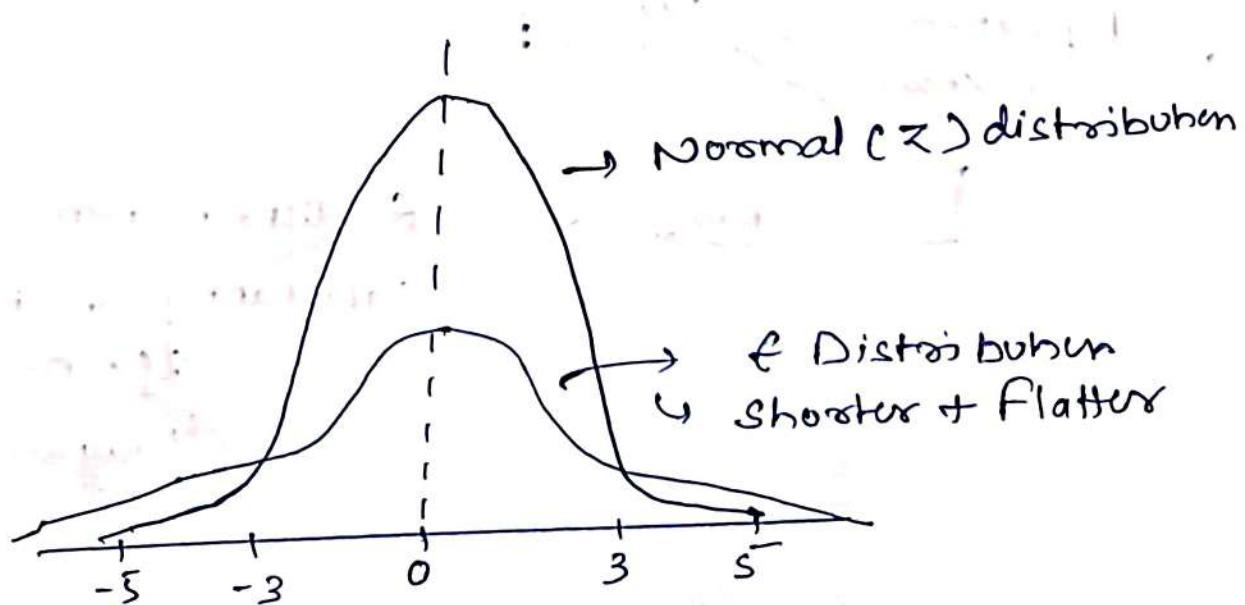
- ① One-sample t-test: Compares the sample mean to a known population mean.
- ② Independent two-sample t-test: Compares the means of two independent groups.
- ③ Paired t-test: Compares means from the same group at different times e.g. (before and after)

2. What is the t -Distribution?

The t -distribution is a probability distribution similar to the normal distribution (bell shaped and symmetric) but;

- It has heavier tails, meaning it accounts for more variability in small samples.
- the shape depends on the degrees of freedom ($df = n - 1$) as the sample size (n) increases, the t distribution approaches the standard normal distribution

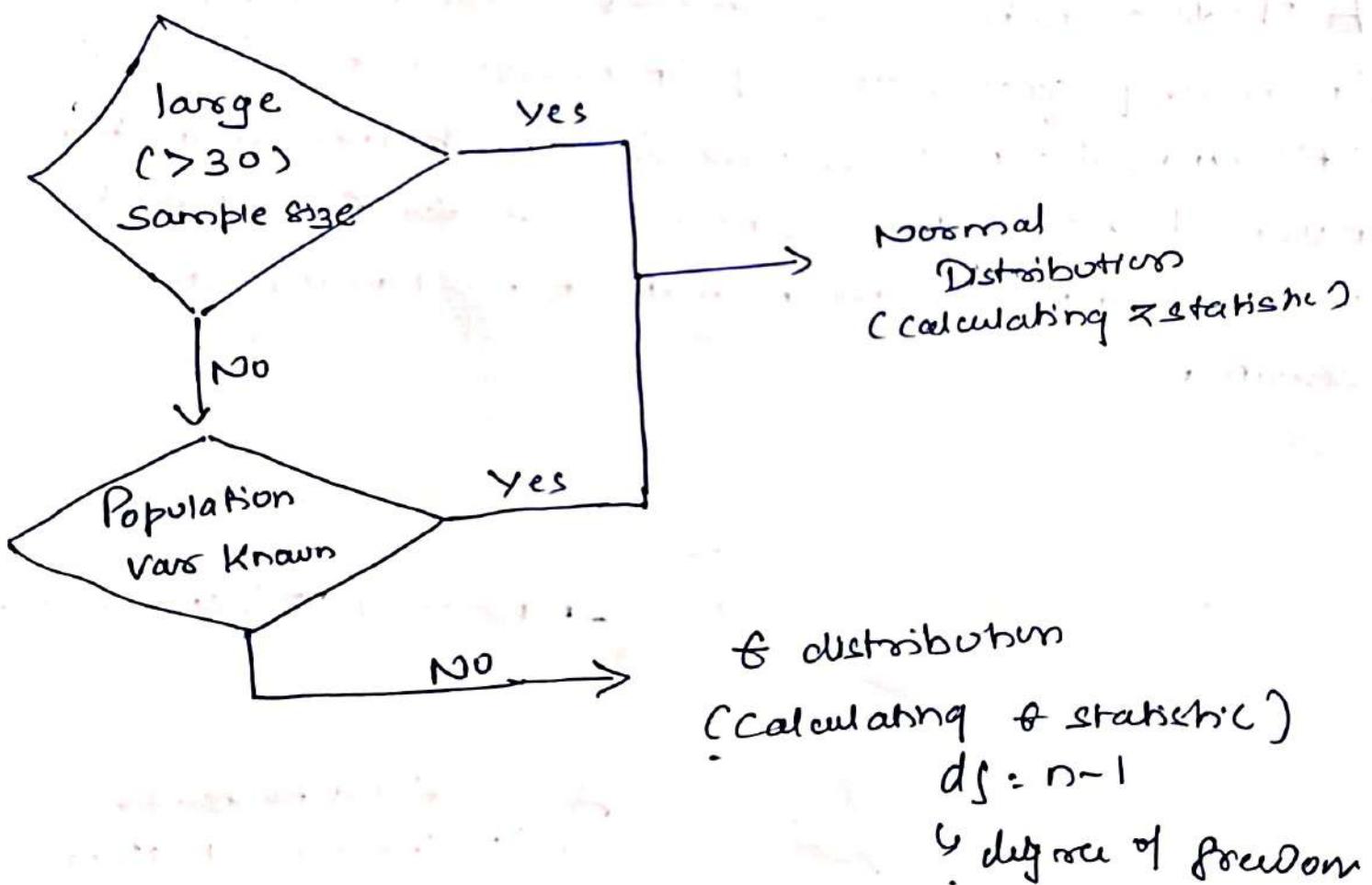
A T distribution is similar to the normal distribution in many cases, for eg it is symmetrical about its central tendency, however it is shorter than the normal distribution and has a flatter tail, which would eventually mean that it has larger standard deviation.



At a sample size beyond 30, the t distribution approximates to the normal distribution.

The most important use of t distribution is that you can approximate the value of the standard deviation of the population (σ) from the sample standard deviation (s). However as the sample size increases more than 30, the t value tends to be equal to the Z value.

Decision Making on flowchart



Let's look at how the method of making a decision changes, if we are using sample's Standard Deviation instead of populations.

Step 1 was to calculate the value of Z_c from the given value of α (significance level), take it as 5%. If not specified.

So for finding Z_c we would use t-table instead of Z-table. The T-table contains the value of Z_c for a given degree of freedom and value of α , here Z_c can also be called as T statistic.

Q you are given the standard deviation of a sample size 25 for a two-tailed hypothesis test of a significance level of 5%.

Given Info :

- Sample size (n) = 25
- Significance level (α) = 5% = 0.05
- Test type : two-tailed hypothesis test

Steps: find degree of freedom (df)

$$df = n - 1 = 25 - 1 = 24$$

Step 2: Significance level per tail $\alpha = 0.05$

$$\Rightarrow \alpha = 0.025$$

Step 3: find t critical from t-table (row corresponding to $df = 24$ & col corresponding to $\alpha = 0.025$)

$$\Rightarrow t_{\text{critical}} = 2.064$$

Ans = 2.064

Q because for sample size = 25, degree of freedom will be $25-1=24$, so look at the value in the t-table corresponding to $df = 24$ and $\alpha = 0.05$ for two-tailed test $\rightarrow 2.064$.

Q You are given the standard deviation of a sample of size 32 for 2-tailed test of $\alpha = 5\%$.

$$n = 32$$

$$df = 32-1 = 31$$

$$\alpha = 0.05 \quad (\because \text{two-tailed} \Rightarrow \alpha = 0.025) \quad \text{(for each tail)}$$

$$t_{\text{critical}} = +1.960 \quad \text{from t-table.}$$

✳ Two Sample Mean Test

Note : To enable Data Analysis tools in excel

File \Rightarrow Options \Rightarrow Addins \Rightarrow Manage \downarrow

Excel Addins \downarrow

OK \leftarrow "Analysis toolpak" \leftarrow No

In Data tab \Rightarrow "Data Analysis Sub tab" gets added here

→ Paired Sample Mean test

- ① Used when data is collected from the same group at two different conditions or at two different times.

e.g.s

- Measuring blood pressure before and after taking medicine.
- Performance score of students before and after training.

→ Unpaired Sample mean test (Independent Sample test)

- ② Used when data is collected from two different groups.

- e.g.s → comparing test scores of GroupA and GroupB under different teaching methods
→ comparing weights of two separate groups of individuals following two diff diets.

Two Sample Mean test

two sample mean test is a statistical method used to compare the means of two independent groups to determine if there is a statistically significant difference between them.

It is commonly used in hypothesis testing, when you have two sets of data and want to know if their population means are different.

Types of Two Sample Mean test

- ① Independent Two Sample Mean Test : Used when the two groups are independent.
- ② Paired Two Sample T test : Used when two groups are related or paired (e.g pre-test and post-test) on same objects

Key Steps in Two Sample Mean test :

① Formulate hypotheses :

- Null Hypothesis (H_0) : Means of the two groups are equal.

$$H_0 : \mu_1 = \mu_2$$

- Alternative Hypothesis (H_1) : The means of the two groups are not equal
 $\rightarrow H_1 : \mu_1 \neq \mu_2$

2. Choose the test

→ Use the t-Test if the sample size are small ($n < 30$) and the population standard deviations are unknown

→ Use the z-test if the sample size are large ($n \geq 30$) and population standard dev are known.

3. Calculate the test statistic : the test statistic depends on whether the variance of the two groups are assumed to be equal.

→ equal variances :

$$t = \frac{(\bar{x}_1 - \bar{x}_2)}{\sqrt{s_p^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} \quad s_p^2 \rightarrow \text{Pooled var}$$

→ unequal variance :

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

4. Determine the degree of freedom (df)

→ For equal variance $\Rightarrow df = n_1 + n_2 - 2$

→ For unequal variance $\Rightarrow df = \left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2} \right)^2$

$$\frac{\left(\frac{s_1^2}{n_1} \right)^2}{n_1 - 1} + \frac{\left(\frac{s_2^2}{n_2} \right)^2}{n_2 - 1}$$

→ Calculate the Pvalue
and test statistic