

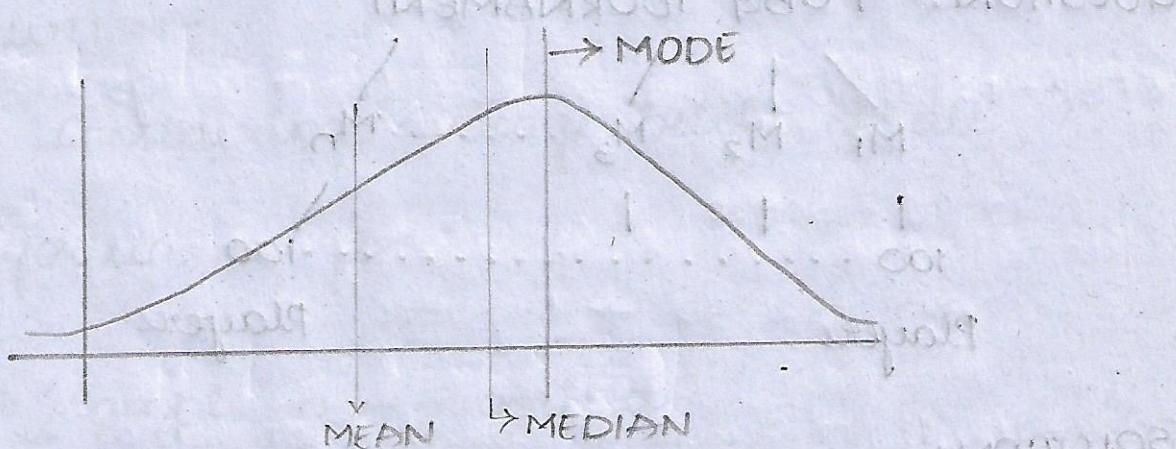
MEAN, MEDIAN AND MODE GETS SHIFTED IN A DISTRIBUTION:

If we have large values, the mean gets attracted to those large values.

→ the outliers has impact on median.

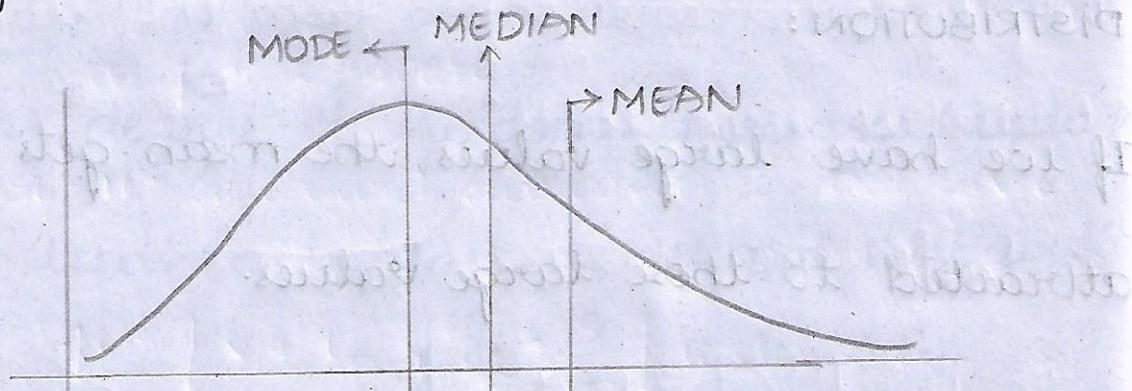
→ so, the mean will be shifted towards right because of the outliers.

\* If the distribution is negatively skewed, the mean will be shifted towards the left.



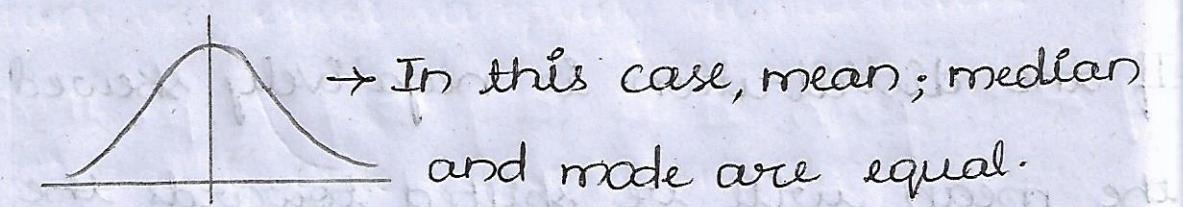
(156)

\* If the distribution is positively skewed, the mean will be shifted towards the right.

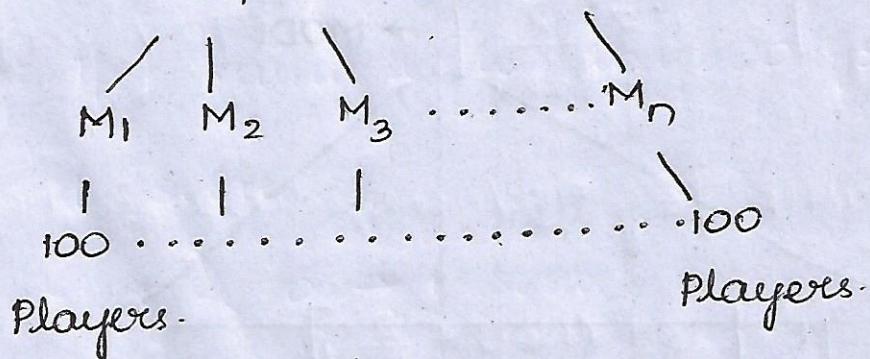


Then the relationship is  $\text{mean} < \text{median} <$

mode



QUESTION: PUBG TOURNAMENT



SOLUTION:

(157)  $\Rightarrow$  In average, the damage is 800, which is mean damage.

$$\Rightarrow \text{Variance}(\sigma^2) = 400$$

$$\Rightarrow \sigma = 20$$

$$\Rightarrow Z = 1.20$$

$$\Rightarrow \frac{D - \mu}{\sigma} = 1.20$$

$$\Rightarrow \frac{D - 800}{20} = 1.20$$

$$\Rightarrow D = 1.20 \times 20 + 800$$

$$\Rightarrow D = 824$$

QUESTION: AVERAGE HEIGHT OF ALL GORILLAS.

SOLUTION :

Considering average height of all gorillas as  $\mu$ .

1. Sample  $n \Rightarrow 5$  gorillas.

2.  $\bar{x}$

3.

CASE-I: Standard deviation is given.

$$\mu = \left[ \bar{x} \pm z^* \frac{\sigma}{\sqrt{n}} \right] \text{ with } \% \text{ confidence}$$

as per  $z^*$  score

CASE-II: Standard deviation is not given

$$\mu = \left[ \bar{x} \pm t_{n-1, \alpha/2} \cdot \frac{s}{\sqrt{n}} \right] \text{ with } \%$$

confidence as per t score

where  $s/\sqrt{n}$  is the standard error of  
standard error sampling distribution.

→  $z^*$  and t-score depends on sample size.

For Sample Size  $n$ :

CASE-I: If 'n' value is less than 30, go with  
t-score

$$\Rightarrow n < 30 \rightarrow t\text{-score.}$$

CASE-II: If 'n' value is greater or equal to 30  
go with  $z^*$ -score.

159

$$\Rightarrow n > 30 \text{ or } n \geq 30 \rightarrow z^* - \text{score}$$

↳ LOOKS LIKE A NORMAL

check  $\sigma$  value:

DISTRIBUTION.

$\sigma$  value

NOT GIVEN

GIVEN

CHECK FOR SAMPLE SIZE 'n'

$n < 30$

$n \geq 30$

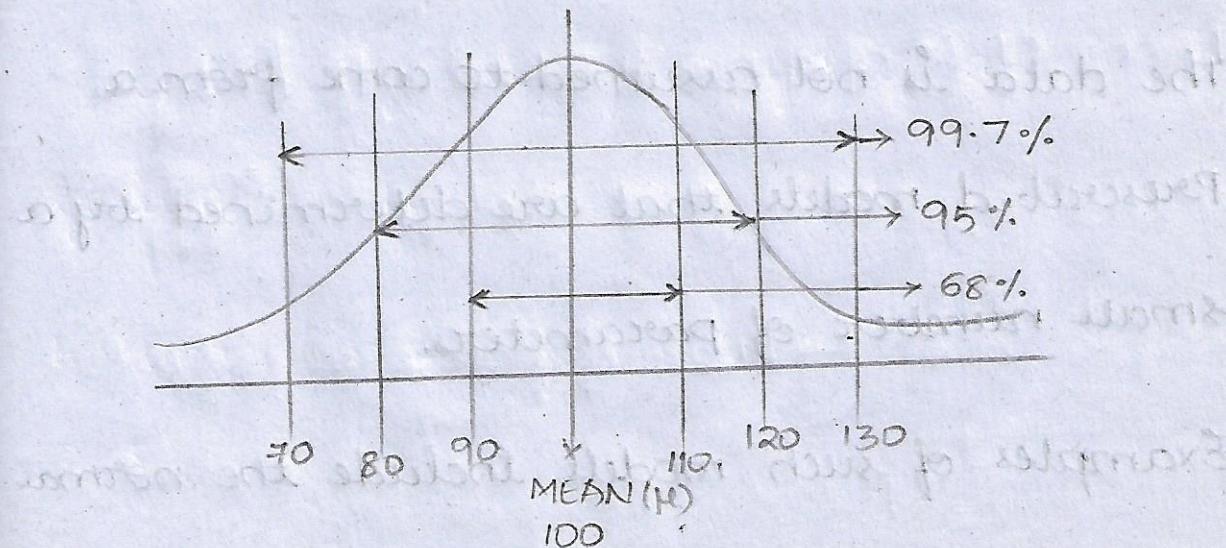
t-score

$z^*$ -score

QUESTION:  $X \sim N(100, 10)$  what is the probability

$$\text{of } X > 130 \rightarrow P(X > 130)$$

SOLUTION:



$\therefore$  The probability of  $x$  160 is 99.7%.

$$\Rightarrow P(X > 160) \rightarrow 99.7\%$$

The rest probability is  $100\% - 99.7\% = 0.03\%$

(or)

$$\frac{100 - 99.7}{2}$$

### PARAMETRIC STATISTICS:

It assumes that the sample data comes from a population that can be adequately modeled by a probability distribution that has a fixed set of parameters.

$$N(\mu, \sigma^2)$$

### NON-PARAMETRIC STATISTICS:

The data is not assumed to come from a prescribed models that are determined by a small number of parameters.

Examples of such models include the normal

(161)

distribution model and the linear regression model.

- \* The Non-parametric statistics are very powerful and invented after the modern computers.

### BOOTSTRAPPING:

A technique for producing a self-compiling compiler - that is, a compiler (or) assembler written in the source programming language that it intends to compile.

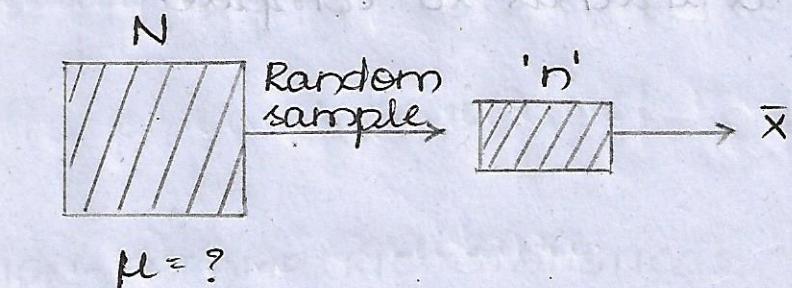
(162)

## HYPOTHESIS TESTING:

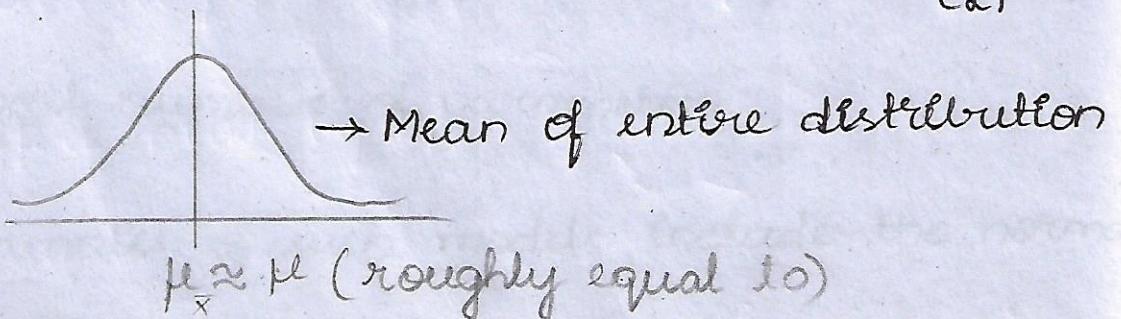
A testing where an analyst tests an assumption regarding a population parameter.

→ It is used to assess the plausibility of a hypothesis by using sample data.

→ Such data may come from a larger population, or from a data-generating process.



From CI (Confidence Interval),  $\mu = [\bar{x} \pm \frac{s}{\sqrt{n}}]$   
 comes from CLT



It defeats cdt  $\rightarrow$  when  $\sigma$  is given.

CLAIM: It is about the entire population  
looking at a very small sample.

Is the claim correct or not?

$\rightarrow$  We are going to test the claims of the population is called as hypothesis testing.

The Hypothesis is classified into 2 types.

They are

1. NULL HYPOTHESIS : STATUS QUO.

The statement is assumed to be true.

It suggests that no statistical relationship

and significance exists in a set of given

single observed variable, between two sets

of observed data and measured phenomena.

$\rightarrow$  Represented by  $H_0$

$\rightarrow$  It is often an initial claim that is

(164)

based on previous analyses or specialized knowledge.

Also  
→ States that a population parameter such as the mean, standard deviation, etc is equal to hypothesized value.

## 2. ALTERNATIVE HYPOTHESIS: BOLD CLAIM

The Contradictory statement.  
It is one in which a difference (or an effect) between two or more variables is anticipated by the researchers; that is, the observed pattern of the data is not due to a chance occurrence.

→ The concept of alternative hypothesis is a central part of formal hypothesis testing.

→ Represented by  $H_1$

- The alternative hypothesis is what we might believe to be true or hope to prove true.
- Also states that the population parameter is smaller, greater or different than the hypothesized value in the null hypothesis.

the Hypothesis Tests are classified according to the form of the alternative hypothesis in the following way.

- If  $H_1$  has the form of  $\mu \neq \mu_0$ , then the test is called as two-tailed test.
- If  $H_1$  has the form of  $\mu < \mu_0$ , then the test is called as a left-tailed test.
- If  $H_1$  has the form of  $\mu > \mu_0$ , then the test is called as a right-tailed test.

## ONE-SIDED HYPOTHESIS:

- Also called as Directional Hypothesis.
- used to determine whether the population parameter differs from the hypothesized value in a specified direction.
- It can specify the direction to be either greater than or less than the hypothesized value.
- A one-sided test has greater power than a two-sided test, but it can't detect whether the population parameter differs in the opposite direction.

## TWO-SIDED HYPOTHESIS:

- Also called as Non-Directional Hypothesis.
- used to determine whether the population parameter is greater than or less than the

hypothesized value.

→ A two-sided test can detect when the population parameter differs in either direction, but has less power than a one-sided hypothesis test.

Z-SCORE:

A numerical measurement that describes a value's relationship to the mean of a group of values.

→ Measured in terms of standard deviations from the mean.

→ If a z-score is '0', it indicates that the data point's score is identical to the mean score.

\*\*

→ Z-score is calculated when the standard deviation is given.

## WHY DO WE USE Z-SCORE?

- The Z-score (or) standard score is a very useful statistic because it
  - ↳ Allows us to calculate the probability of a score occurring within our normal distribution.
  - ↳ Enables us to compare two scores that are from different normal distributions.

### t-SCORE :

A t-score (or) t-statistic is the result of applying a T-test and it represents a point from a STUDENT'S T-DISTRIBUTION with  $n-k$  degrees of freedom where ' $n$ ' is the total sample size and ' $k$ ' is the number of test groups in a A/B test.

(169)

WHY DO WE USE t-SCORE?

→ The t-distribution is a probability distribution that is used to estimate population parameters where the sample size is small and/or when the population variance is unknown.

P-VALUE:

The p-value or calculated probability of finding the observed, or more extreme results when the null hypothesis ( $H_0$ ) of a study question is true.

→ The definition of 'extreme' depends on how the hypothesis is being tested.

CONFIDENCE INTERVAL:

It refers to the probability that a population parameter will fall between a

set of values for a certain proportion of times.

### SIGNIFICANCE LEVEL ( $\alpha$ ):

The probability of rejecting the null hypothesis when it is True.

For Example:

→  $\alpha$  of 0.05 indicates 5% risk of concluding that a difference exists when there is no actual difference.

### CRITICAL VALUE:

It is a point on the test distribution that is compared to the test statistic to determine whether to reject the Null hypothesis.

→ If the absolute value of your test statistic is greater than the critical value,

(171)

we can declare statistical significance and reject the null hypothesis.

#### ONE-TAILED TEST:

A statistical test in which the critical area of a distribution is one-sided so that the it is either greater than or less than a certain value, <sup>but not</sup> both.

→ If the sample being tested falls into the one-sided critical area, the alternative hypothesis will be accepted instead of the null hypothesis.

#### TWO-TAILED TEST:

A statistical test method in which the critical area of a distribution is two-sided and tests whether a sample is greater than or less than a certain range of values.

→ It is used in null-hypothesis testing and testing for statistical significance.

### RIGHT-TAILED TEST:

A test where the rejection region is located to the extreme right of the distribution.

→ It is conducted when the alternative hypothesis contains the condition of greater than a given quantity ( $x$ )

$$\text{e.g., } H_1 > x$$

→ The upper tail is stated as "right tail" since the highest numbers will appear on the right. It is because, the highest values on a number line are to the right.

### LEFT-TAILED TEST:

A test where the rejection region is located to the extreme left of the distribution.

- It is conducted when the alternative hypothesis contains the condition of less than a given quantity ( $x$ ).  
i.e.,  $H_1 < x$ .

### TYPE II ERROR:

The mistake of failing to reject the null hypothesis when it is false.

### TYPE I ERROR:

The mistake of rejecting the null hypothesis when it is true.

174

QUESTION: HYPOTHESIS TEST FOR A COIN.

SOLUTION:

COIN  $\begin{cases} \text{UNBIASED} \rightarrow P(H) = \frac{1}{2} - 0.5 \\ \text{BIASED TOWARDS HEAD} \rightarrow P(H) > 0.5 \end{cases}$

STEP-1: Figure out the test statistics

e.g., flip the coin 5 times.

$\rightarrow$  As we can count number of Heads,

denote as  $x$ .

STEP-2: Identify the  $H_0$  &  $H_1$ .

$\Rightarrow H_0 = \text{COIN IS UNBIASED}$

$\Rightarrow H_1 = \text{COIN IS BIASED TOWARDS HEAD}$

STEP-3: Compute the p-value.

p-value is written as  $P(x=5 | H_0)$

$|$   $\rightarrow$  denotes CONDITIONAL PROBABILITY

Given that, the coin is unbiased, find the probability of getting 5 heads in 5 flips is

$$\rightarrow \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2}$$

Exactly 5 heads in 5 flips, then

$$\rightarrow \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} = \frac{1}{32}$$

$$\Rightarrow P(X=5 | H_0) = \frac{1}{32}$$

$$\Rightarrow P(X=5 | H_0) = 0.03 = 3\% \text{ chance.}$$

$\therefore$  There is 3% of chance of getting 5 heads in 5 flips for the coin which is unbiased.

$\rightarrow$  Basically we go with 5% p-value.

$\rightarrow$  If anything is less than 5%, then we reject the Null Hypothesis ( $H_0$ )

$\rightarrow$  In the above case,  $H_0$  is rejected as the coin is unbiased.

$\rightarrow$  Accept the Alternative Hypothesis ( $H_1$ ).

\* \*  $\rightarrow$  5% is the threshold value and used in medical field around 2.5%

→ If the p-value is greater than the 5%, then we say that failed to reject the Null hypothesis ( $H_0$ )