

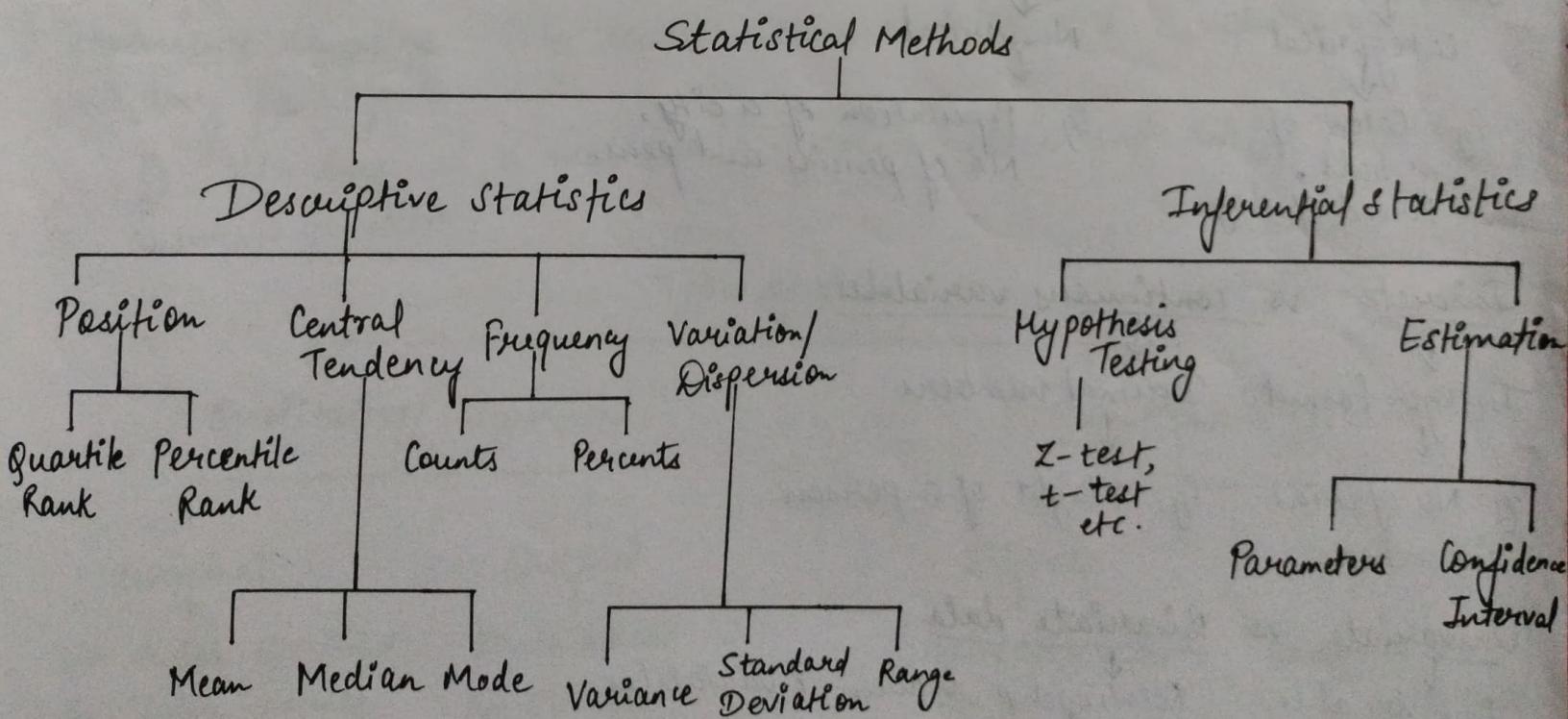
STATISTICS

Statistics is the branch of Mathematics where we collect, organize, analyse and represent the data for better decision making. We apply statistics to different problems.

Statistics refers to a scientific approach used to :

- 1) Collect data
- 2) Interpret and Analyze data
- 3) Assess the reliability of conclusions based on sample data.

For Example: Let us consider that a researcher wants to know which medicine is more effective - A or B , based on the symptoms, then, using the statistical methods, the researcher might conclude with 95% confidence that one medicine was superior to the other.



Descriptive statistics: It is a summary that describes or summarizes the collection of information / data. It summarizes the sample data rather than learning about the population that sample data is representing.

Inferential statistics: It is the process of data analysis where we make the conclusions about population data using sample data.

Population: It is the entire group that we want to draw the conclusions about.

sample (n): It is a small part of population.

Variables in Statistics:

Qualitative vs. Quantitative Variables:

↓
categorical

Eg: Color of
a ball,
Breed of dog.

↓
Numerical

Eg: Population of a city,
No. of pencils and pens.

Discrete vs. continuous variables:

↓
Integer format

Eg: No. of coins
in a box.

↓
Decimal numbers

Eg: Weight of a person.

Univariate vs. Bivariate data:

↓
One Variable

Eg: Average weight of
all students

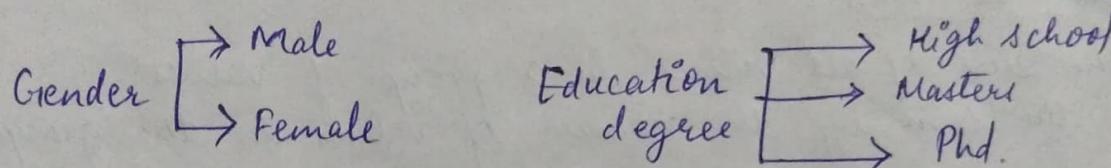
↓
Relationship between 2 variables.

Eg: Relationship between
height and weight of students.

Sampling Techniques:

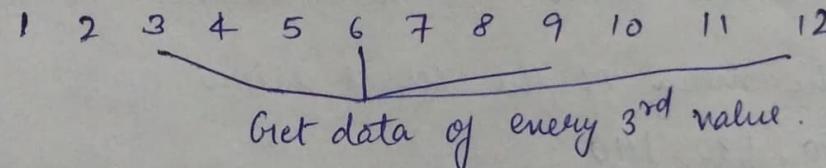
- 1] Simple Random Sampling: Every member of the population (N) has an equal chance of being selected for your sample (n).
- 2] Stratified Sampling: Population (N) is split into non-overlapping groups.

Eg:



- 3] Systematic Sampling: It is a probability sampling method where researchers select members from population at n^{th} interval.

Eg:



- 4] Convenience Sampling: Only those who are interested in the survey will only participate.

Eg: Standing at a mall or a grocery store and asking people to answer questions.

Types of data

Qualitative / Categorical

Nominal

No order specified

Eg: Color, names of people.

Ordinal

Rank based

Eg: Medals received, score between 1-10

Quantitative / Numerical

Discrete

Integer

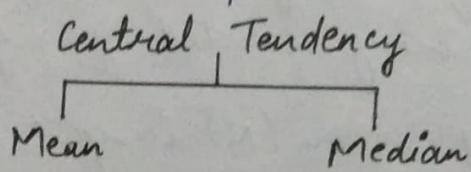
Eg: No. of coins

Continuous

Decimal

Eg: Height, weight

Measures of Central Tendency: Central Tendency is the statistical measure that identifies a single value as a representative of an entire distribution.



To compute the mean:

1. Find the sum of all values in a group of values.
2. Divide the sum by the number of values in the group.

$$\text{Mean} = \frac{\text{Sum of all values}}{\text{Number of values}}$$

μ - Population mean \bar{x} - Sample mean.

- Q. Compute the mean of the scores below

100, 170, 130, 140

$$\mu = \frac{540}{4} = \underline{135}$$

To compute the median:

1. List scores from smallest to largest.
2. With an odd number of scores, the median is the middle score.
3. With an even number of scores, the median is the sum of the middle two scores divided by 2.

$$\text{Median} = \text{Sum of middle two scores} / 2$$

- Q. Find median score

$$\begin{aligned} 1. \quad & 100, 170, 130, 140, 160 \xrightarrow{\text{Sorting}} 100, 130, 140, 160, 170 \xrightarrow{\text{Median}} \\ 2. \quad & 100, 170, 130, 140 \xrightarrow{\text{Sorting}} 100, \underbrace{130, 140, 170}_{\frac{130+140}{2}} = \underline{135} \rightarrow \text{Median.} \end{aligned}$$

As measures of central Tendency, the mean and the median each have advantages and disadvantages.

- The median may be a better indicator of the most typical value if a set of scores has an outlier. An outlier is an extreme value that differs greatly from other values.
- However, when the sample size is large and does not include outliers, the mean score usually provides a better measure of central tendency.

Eg: If we consider annual income of 4 people to be ₹ 7,00,000, ₹ 8,00,000, ₹ 9,00,000 & ₹ 1,00,00,000, then the 4th value will be an outlier. If we choose a measure to estimate the annual income of a ~~typical~~ person in general, the mean will greatly overestimate the income (because of the outlier); while the median will not.

Mode: It is the most frequent score in the dataset. Normally, the mode is used for categorical data where we wish to know which is the most common category.

A problem with mode is that it will not provide us with a very good measure of central tendency when the most common mark is far away from the rest of the data in the dataset.

Measures of dispersion: Statisticians use summary measures to describe the amount of variability or spread in a set of data. The most common measures of variability are the range, the Interquartile Range (IQR), variance and standard deviation.

1. Range: It is the difference between the largest and smallest values in a set of values.

Eg: 2 4 9 5 7 3 → Range = 9 - 2 = 7

It is easy to calculate but one drawback is, it ignores the middle values.

$$\text{Eg: } 1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7 \ 8 \ 9 \rightarrow \text{Range} = 9 - 1 = 8$$

$$1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 9 \rightarrow \text{Range} = 9 - 1 = 8$$

Range is same in both cases, but there is a huge difference between the distribution of the data points.

2. Interquartile Range (IQR): It is a measure of variability, based on dividing a dataset into quartiles.

What is a quartile?

→ Before understanding quartiles, we need to understand percentiles.

Percentile: A percentile is a value below which a certain percentage of observations lie.

Eg: 99 percentile - It means a person has performed better than 99% of the people.

$$\text{percentile} = \frac{\text{No. of values below a given no.}}{\text{Total no. of values.}}$$

$$\text{Eg: } X = \{1, 2, 3, 4, 5, 5, 7, 20, 11, 10\}$$

$$y = 10$$

$$\text{percentile} = \frac{7}{10}$$

Quartile: Quartiles divide a rank-ordered data set into four equal parts. The values that divide each part are called the first, second and third quartiles, and they are denoted by Q_1 , Q_2 and Q_3 respectively.

Relationship between quartiles and percentiles.

$Q_1 \rightarrow P_{25}$, $Q_2 \rightarrow P_{50}$, $Q_3 \rightarrow P_{75}$. Q_2 is the median value in a dataset.

The InterQuartile Range (IQR) is a measure of variability, based on dividing a dataset into quartiles. The advantage of IQR over Range is that it takes all the values of dataset into consideration.

3. Variance: It is the average squared deviation from the population mean.

$$\sigma^2 = \sum (x_i - \mu)^2 / N$$

where, σ^2 = population variance

μ = population mean

x_i = i th element from population

N = no. of elements in population.

$$\text{For sample variance, } s^2 = \sum (x_i - \bar{x})^2 / (n-1)$$

Why sample variance is divided by $n-1$?

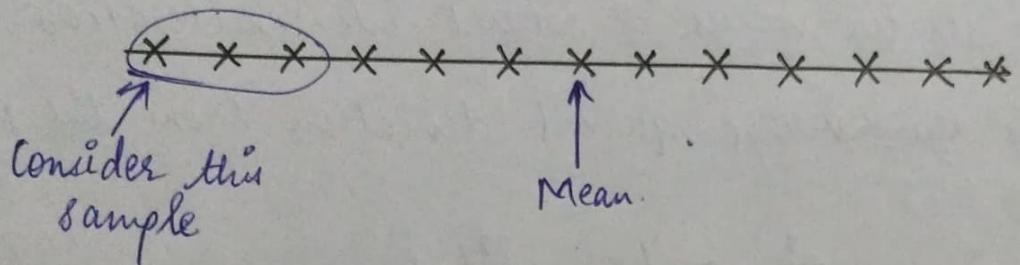
\Rightarrow	Population (N) $\mu = \frac{\sum_{i=1}^N x_i}{N}$ $\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu_x)^2}{N}$	Sample (n) $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$ $s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$
---------------	--	---

When we take out any sample, we need to infer some information from it so that we can draw some conclusions about the population data. This is possible only when we have

Population mean \approx Sample mean
and

Population variance \approx Sample variance.

Consider the below example,



Now, when we calculate, we get ~~pop~~ huge difference between the population mean and sample mean and the same with population variance and sample variance.

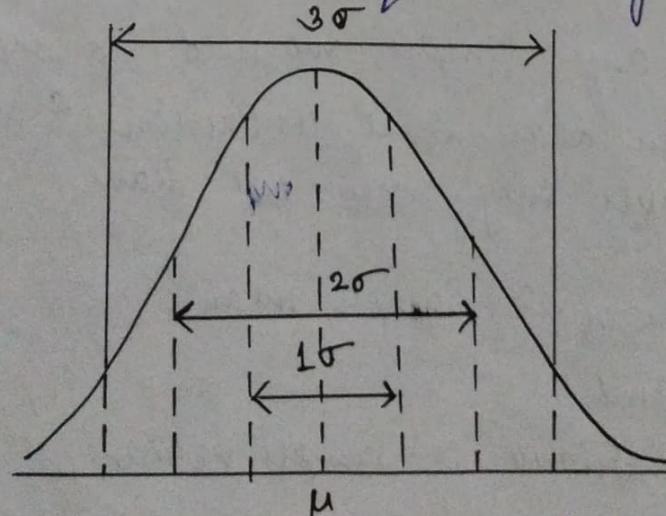
So, after a lot of research, ~~that~~ it was found out that instead of having ' n ' in the denominator of sample variance, if we have $n-1, n-2, n-3, \dots$ then

$$\text{sample } \frac{\text{variance}}{n-1} \approx \text{Population variance}$$

And hence, it was found that " $n-1$ " gave the closest answer ~~than~~ as compared with $n-2, n-3, \dots$

Hence, sample variance is divided by $n-1$.

4. Standard Deviation: It is the square root of the variance.



This is also used to identify whether a data point is an outlier or not.

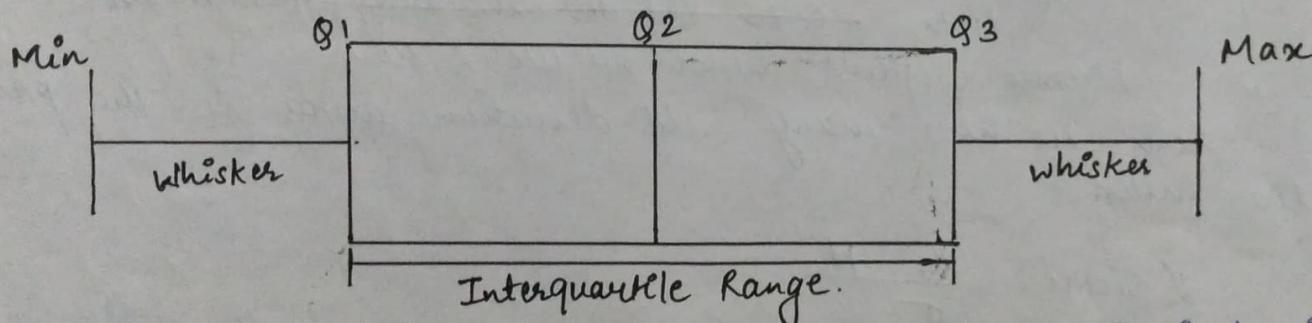
Five Number Summary (In order to remove outliers):

1. Minimum
2. First Quartile (25 percentile - Q1)
3. Median (50 percentile - Q2)
4. Third Quartile (75 percentile - Q3)
5. Maximum

$$IQR = Q3 - Q1$$

$$\text{Lower fence} = Q1 - 1.5 * (IQR)$$

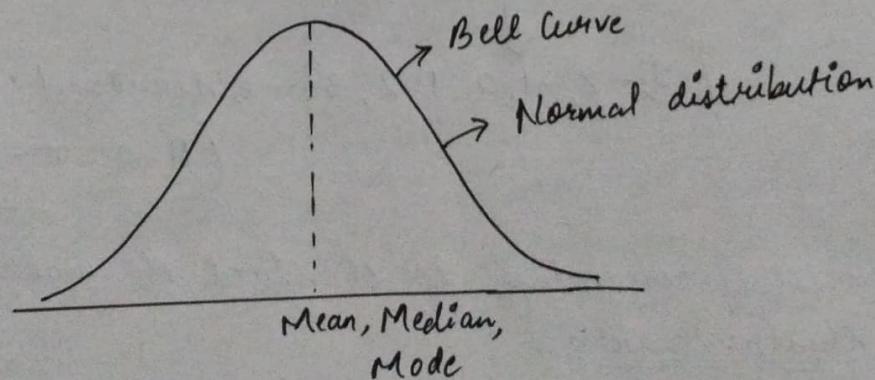
$$\text{Higher fence} = Q3 + 1.5 * (IQR)$$



The plot shown above is called the Box plot. This is basically used to get the outliers. Whatever points lie beyond this box plot are treated as the outliers.

Distributions:

① Gaussian / Normal Distribution:



Empirical formula - (68-95-99.71. rule)

68.1% of data lies in 1st std. deviation.

95% of data lies in 2nd std. deviation.

99.71% of data lies in 3rd std. deviation.

This rule is applied on Gaussian/Normal distribution.

Z-score:

Let us consider, $\mu=4$ & $\sigma=1$

When we say 4.5, it means +0.5 std. deviation away from mean.

4.75, ~~+0.75 std. deviation away from mean~~ calculation becomes difficult, hence we use Z-score.

Z-score tells us how many std. deviations away is the point from the mean.

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

$$= \frac{4.75 - 4}{1} = 0.75 \text{ std. deviation to the right.}$$

Dataset - $\{1, 2, 3, 4, 5, 6, 7\}$

↓
Z-score

↓
 $\{-3, -2, -1, 0, 1, 2, 3\} \rightarrow$ Standard Normal distribution.
($\mu=0, \sigma=1$)

The process of converting to standard Normal distribution is called standardization.

Normalization: $\{\mu=0, \sigma=1\}$

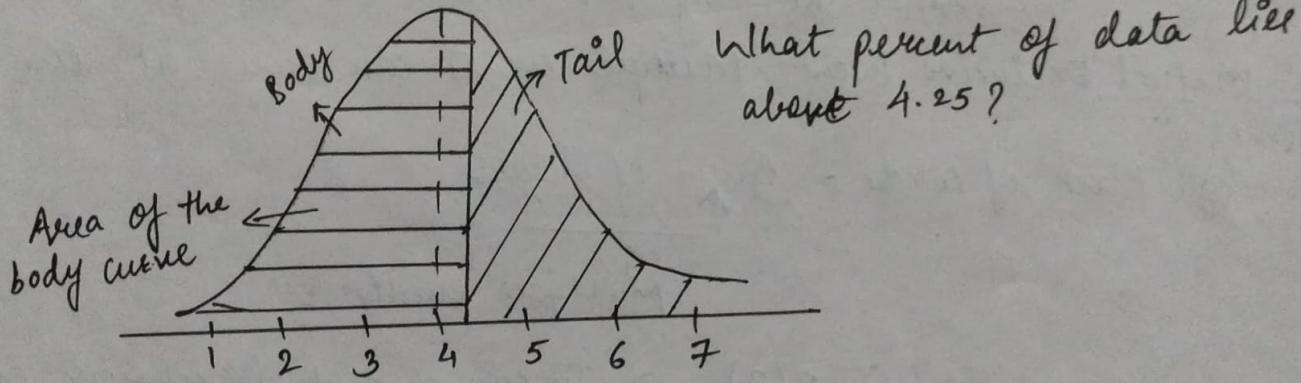
If we want to shift the values between a particular range then we use normalization. For example, between 0 to 1.

Min Max Scales: (0 to 1)

Where to use Normalization?

→ In deep learning, we have images which have pixels which range between 0 (0-255) where we do normalization and bring it in between 0 to 1.

Q.



$$\Rightarrow Z = \frac{x_i - \mu}{\sigma} = \frac{4.25 - 4}{1} = 0.25$$

Z-score: For area in a left tail, we need to look at the left tail z-table. or $1 - \text{right Area}$

$$\therefore \text{Area of tail} = 1 - 0.5987 = \underline{\underline{0.4013}}$$

i.e. 40% of data lies above 4.25

Probability: It is a measure of the likelihood of an event.

Eg: Flipping a coin - {H, T}

Probability of head or tail = $\frac{1}{2}$

Addition Rule (Probability, "or")

Mutual Exclusive Event: Two events are mutually exclusive if they cannot occur at the same time.

Eg: Rolling a dice {1, 2, 3, 4, 5, 6}

We cannot get any two numbers at the same time.

Non Mutual Exclusive Event: Multiple events can occur at the same time.

Eg: Deck of cards - Queen of hearts
↓ ↑
 ↓
Multiple events.

$$P(A \text{ or } B) = P(A) + P(B) \rightarrow \text{In case of mutual exclusive event.}$$

$$P(A \text{ or } B) = P(A) + P(B) - P(A \cap B) \rightarrow \text{In case of non-mutual exclusive event.}$$

Multiplication Rule:

Independent Events:

Eg: {1, 2, 3, 4, 5, 6}

First time I can get 1, next time 3 and so on.

Each and every event is independent of each other.

Dependent Events:

Eg: Bag of 3 red marbles and 2 green marbles

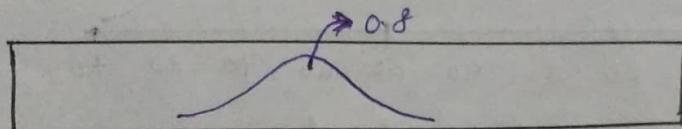
So picking out any marble will have impact on ~~on~~ probability of other marble.

Permutation and Combination:

$$\text{Permutation: } {}^n P_r = \frac{n!}{(n-r)!}$$

$$\text{Combination: } {}^n C_r = \frac{n!}{r!(n-r)!}$$

p Value:



Suppose, this is a space bar of a laptop. The area marked is the area where a person touches the most. Suppose out of 100 times, the person touches 80 times in the place shown above.

p-value is the probability that a particular statistical measure, such as the mean or standard deviation, of an assumed probability distribution will be greater than or equal to observed results.

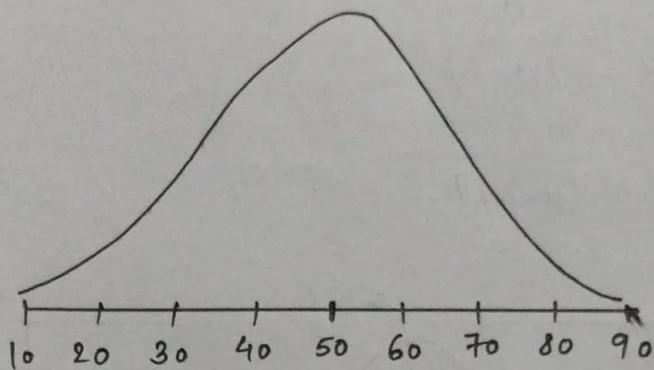
Hypothesis Testing:

Consider an example of tossing a coin 100 times. If we get head for 50 times then we say that the coin is fair.

Null Hypothesis: Coin is fair

Alternate Hypothesis: Coin is unfair.

Let, the standard deviation be 10.



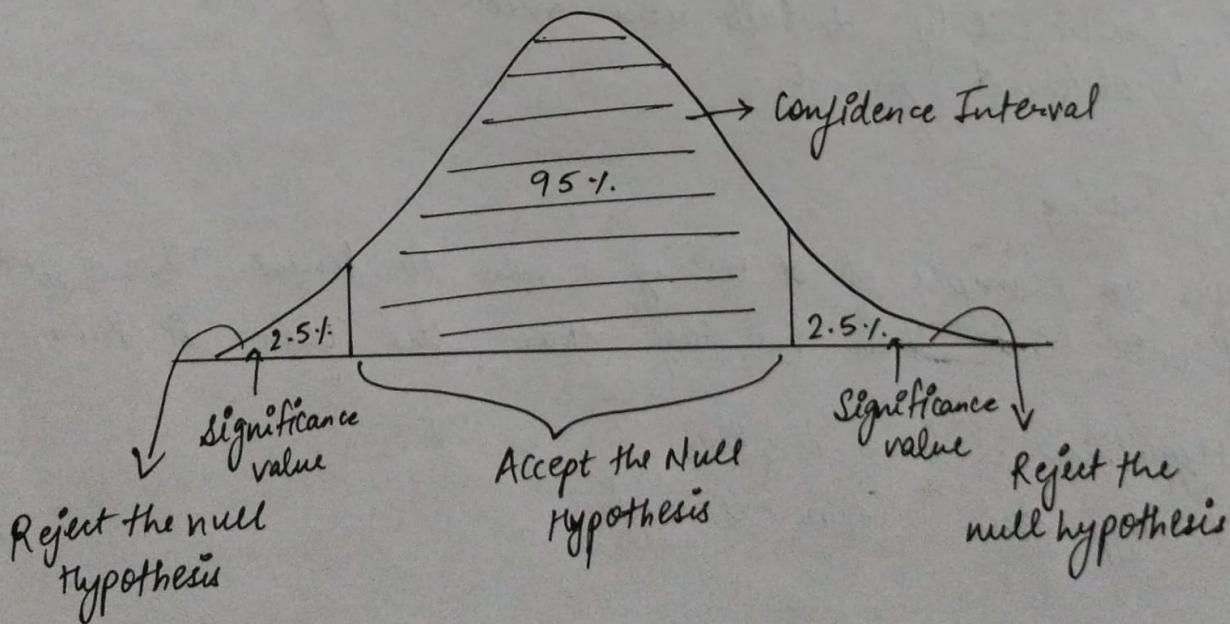
i) Suppose we got 30 times head

In order to prove coin is fair, we need to have the value nearer to the mean.

Now, how much away can the value be from mean is given by Significance value.

Let $\alpha = 0.05 \rightarrow$ Defined by domain expert.

$1 - 0.05 = 0.95 \rightarrow$ Confidence Interval.



Significance value \neq P-value

Type 1 and Type 2 error:

Null Hypothesis (H_0)

Alternate Hypothesis (H_1)

Outcome 1: We reject the Null Hypothesis, when in reality it is false \rightarrow Yes, good decision.

Outcome 2: We reject the Null Hypothesis, when in reality it is true \rightarrow No, bad decision.

Outcome 3: We retain the Null Hypothesis, when in reality it is False \rightarrow No, Bad Decision

Type 2 error

Outcome 4: We retain the Null Hypothesis, when in reality it is True \rightarrow Yes, Good decision.

Confusion Matrix :

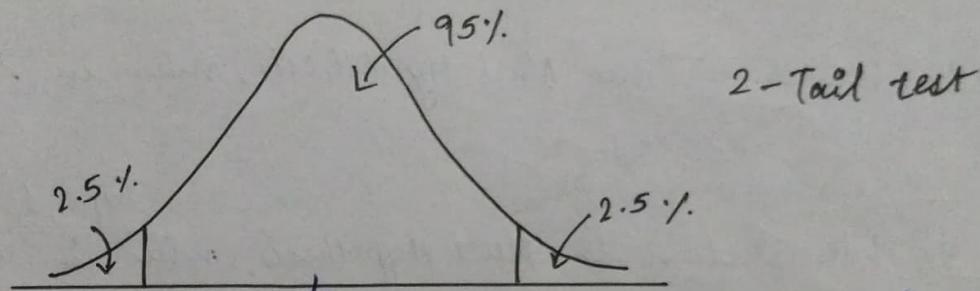
		P	N	
		TP	TN	Type 2
T	T			
	F	FP	FN	

Type 1

1 Tail and 2 Tail test:

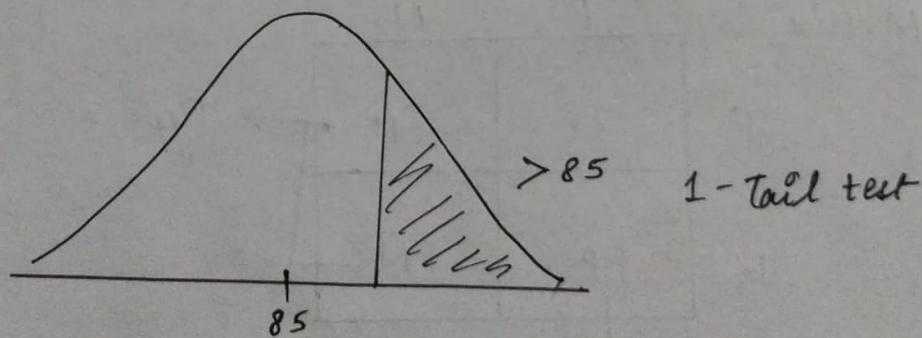
Q. Colleges in Karnataka have 85% placement rate. A new college was recently opened and it was found that a sample of 150 students had a placement rate of 88%, with a standard deviation of 4%. Does this college have a different placement rate than other colleges?
 $\alpha = 0.05$

\Rightarrow 2-Tailed test



From the given question, we can have the placement rate either greater than 85 or less than 85, that is why it is a 2-tailed test.

Suppose, we have a question as greater than 85%, then it will be a 1-tail test.



Point Estimate: The value of any statistic that estimates the value of a parameter.

$$\begin{array}{ccc} \overline{x} & \longrightarrow & \mu \\ \text{sample mean} & & \text{Population mean} \end{array}$$

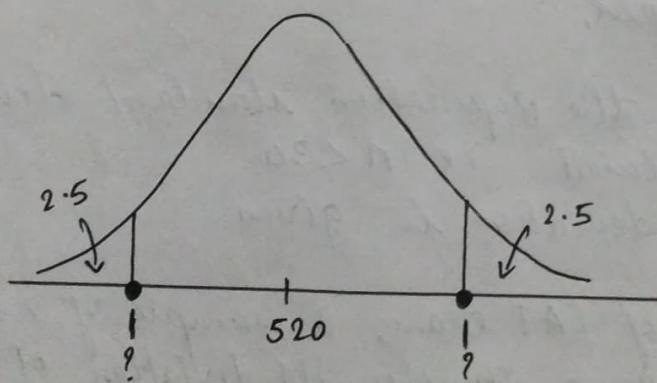
Through \bar{x} , we are estimating μ , this is called point estimate.

Confidence Interval:

Point Estimate \pm Margin of Error.

- Q. On the quant test of CAT exam, the standard deviation is known to be 100. A sample of 25 test takers has a mean of 520 score. Construct a 95% CI about the mean.

$$\Rightarrow \sigma = 100, n = 25, \alpha = 0.05, \bar{x} = 520$$



- ① When population standard deviation is given, we apply z-test
Point Estimate \pm Margin of error

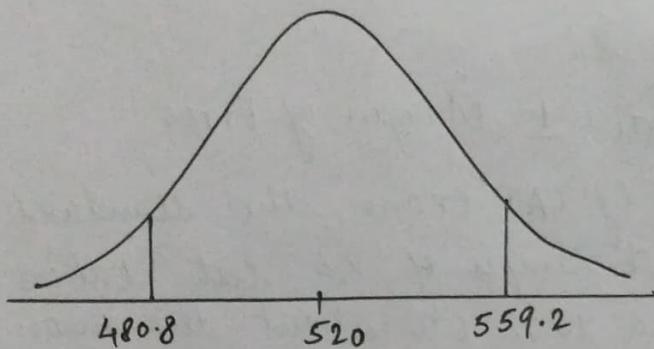
$$\bar{x} \pm Z_{\alpha/2} \left[\frac{\sigma}{\sqrt{n}} \right] \rightarrow \text{Standard error.}$$

For a Z-test to happen,

- ① We know the population standard deviation or
- ② We do not know our population standard deviation but our sample size ≥ 30 .

$$\text{Upper bound} - \bar{x} + Z_{0.05} \frac{\sigma}{\sqrt{n}} = 1 - 0.025 = 0.975 \xrightarrow{Z_{0.05} \text{ tabl.}} 1.96 \\ - 520 + (1.96)(20) = 559.2 \xrightarrow{\text{C.I.}}$$

$$\text{Lower bound} - \bar{x} - Z_{0.05} \frac{\sigma}{\sqrt{n}} = 520 - (1.96)(20) = 480.8 \xrightarrow{\text{Interv.}}$$



For a t-test to happen,

- ① We do not know the population standard deviation.
- ② Sample size is small i.e. $n < 30$.
- ③ Sample standard deviation is given.

Q. On the quant test of CAT exam, a sample of 25 test takers has a mean of 520 with a sample std. deviation of 80. Construct 95% confidence interval about the mean.

$\Rightarrow n = 25, \bar{x} = 520, s = 80, \sigma = 0.05$
Population std. dev is not given \rightarrow t-test

$$\bar{x} \pm t_{\alpha/2} \left[\left(\frac{s}{\sqrt{n}} \right) \right] \rightarrow \text{Standard Error}$$

$$\text{Degree of freedom} = n - 1 = 25 - 1 = 24$$

$$\text{Upper bound} = 520 + 2.064 \left(\frac{80}{\sqrt{25}} \right) = 553.024$$

$$\text{Lower Bound} = 520 - 2.064 \left(\frac{80}{\sqrt{25}} \right) = 486.97$$

$$[486.97 \rightarrow 553.024]$$

Standard Error: It indicates how different the population mean is likely to be from a sample mean.

Chi square Test: It claims about population proportions.

It is a Non-parametric test that is performed on categorical (Nominal or Ordinal) data.

Parametric vs Non-Parametric:

Parametric statistics are based on assumptions about the distribution of population from which the sample was taken.

Non-Parametric statistics are not based on assumptions, that is, the data can be collected from a sample that does not follow a specific distribution.

- Q. In the 2000 Indian Census, the age of the individuals in a small town were found to be the following:

Less than 18	18-35	>35
20%	30%	50%

In 2010, age of $n=500$ individuals were sampled. Below are the results.

<18	18-35	>35
121	288	91

Using $\alpha = 0.05$, would you conclude the population distribution of ages has changed in the last 10 yrs?

⇒

<18	18-35	>35	
121	288	91	→ Observed.] Huge difference
500×0.2 $= 100$	500×0.3 $= 150$	500×0.5 $= 250$	→ Expected

H_0 = The data meets the distribution of 2000 census.

H_1 = The data does not meet the distribution of 2000 census.

$$\alpha = 0.05 \quad (95\% \text{ CI})$$

Degree of Freedom = $n - 1 = 3 - 1 = 2$
 \downarrow
No. of categories
(This is a 2-tail test)

Chi-square table.

If χ^2 is greater than 5.99, then reject H_0 .

$$\begin{aligned}\chi^2 &= \sum \frac{(f_o - f_e)^2}{f_e} \\ &= \frac{(124 - 100)^2}{100} + \frac{(288 - 150)^2}{150} + \frac{(91 - 250)^2}{250}\end{aligned}$$

$$\underline{\chi^2 = 232.494} > 5.99 \rightarrow \text{Reject } H_0.$$

\therefore Conclusion is, the data does not meet the distribution of 2000 census.

If P-value < significance value

~~Accept~~ Reject the Null Hypothesis

Covariance :

X	Y
Weight	Height
50	160
60	170
70	180
75	181

$x \uparrow \quad y \uparrow$

$x \downarrow \quad y \downarrow$

Quantify relationship between X & $Y \rightarrow$ Covariance.

Covariance

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

In case of sample this is $\frac{1}{(n-1)}$

$x \uparrow y \uparrow$
 $x \downarrow y \downarrow$

\rightarrow +ve
correlation

$x \uparrow y \downarrow$
 $x \downarrow y \uparrow$

-ve
correlation

$0 \rightarrow$ No correlation

Covariance: No particular magnitude. We have direction but magnitude wise, it is not a good choice. So, we move to Pearson's correlation.

Pearson Correlation Coefficient: It restricts the value between $(-1, 1)$. The more towards $+1$, the more it is positively correlated and the more towards -1 , the more it is negatively correlated.

$$r_{(X, Y)} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

This captures the linear properties well.

Spearman Rank correlation:

$$\text{Spear}(X, Y) = \frac{\text{Cov}(R(X), R(Y))}{\sigma_{R(X)} \sigma_{R(Y)}}$$

Height (X)	Weight (Y)	R(X)	R(Y)
170	75	2	2
160	62	3	3
150	60	4	4
145	55	5	5
180	85	1	1

Let's use this correlation because it ~~too~~ supports Non-linear properties.

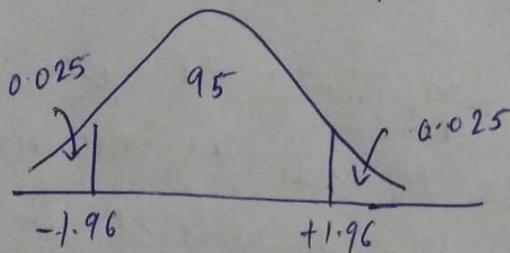
- Q. The average weight of all residents in Bangalore city is 168 pounds with a standard deviation of 3.9. We take a sample of 36 individuals and the mean is 169.5 pounds. C-I = 95%.

⇒

$$H_0: \mu = 168$$

$$H_1: \mu \neq 168$$

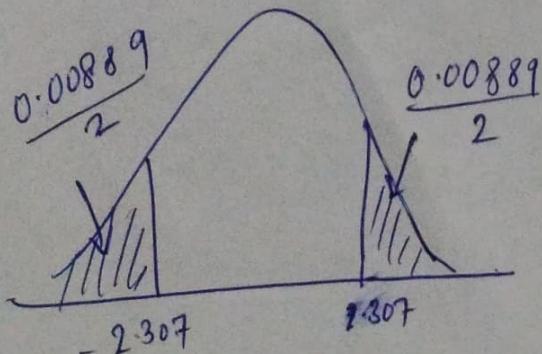
$$\alpha = 0.05$$



We use Z-test here.

$$Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{169.5 - 168}{\frac{3.9}{\sqrt{36}}} = \underline{\underline{2.307}}$$

$2.307 > 1.96 \rightarrow$ We reject the null hypothesis.



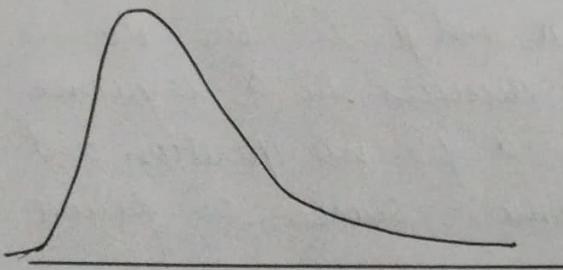
Z-table for 2.307 → 0.99111

$$1 - 0.99111 = 0.00889$$

$$\begin{aligned} P\text{-value} &= 0.0044 + 0.0044 \\ &= \underline{\underline{0.0088}} \end{aligned}$$

$0.0088 < 0.05 \rightarrow$ We reject null hypothesis.

Log Normal Distribution:



Eg: 1) Wealth Distribution.

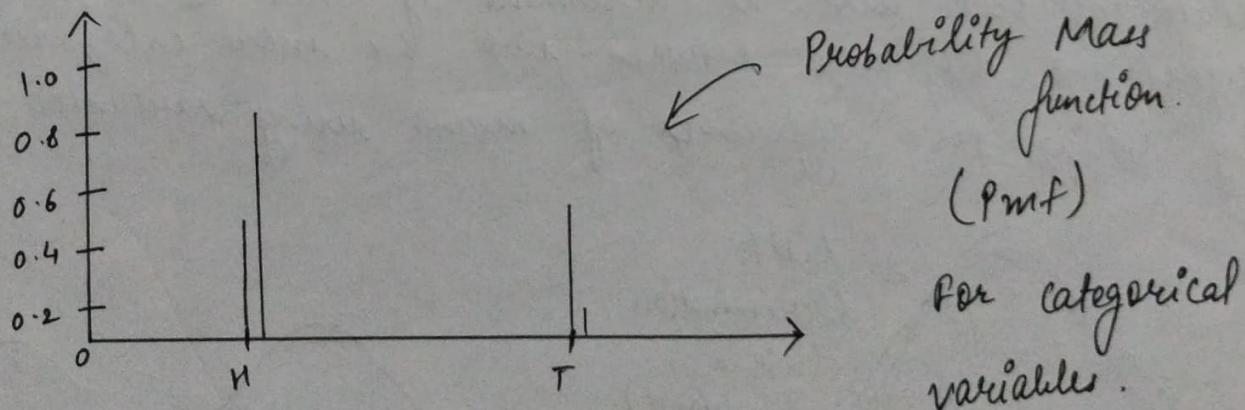
2) People writing big comments.

If $X \approx \text{Log Normal distribution}$
then

$y \approx \ln(x) \rightarrow \text{Normal distribution}$

If a random variable X is log-normally distributed, then $Y = \ln(X)$ has a normal distribution. Equivalently, if Y has a normal distribution, then the exponential function of Y , $X = \exp(Y)$, has a log-normal distribution.

Bernoulli's distribution: It is the discrete probability distribution of a random variable which takes the value 1 with probability p and the value 0 with probability $q = 1 - p$.
It has only 2 outcomes \rightarrow either 0 or 1.



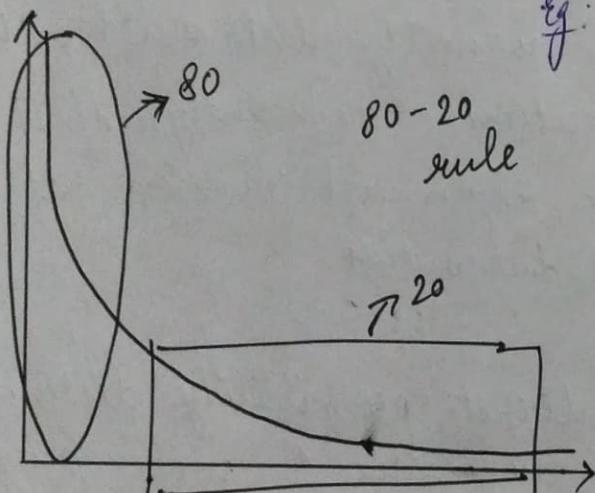
$$\text{PMF} \quad q = 1 - p \quad \text{if } k=0$$

$$\frac{p}{\left(p^k (1-p)^{1-k}\right)} \quad \text{if } k=1$$

Binomial Distribution

The binomial distribution with parameters n and p is the discrete probability distribution of the number of successes in a sequence of n independent experiments, each asking a yes-no question and each with its own Boolean-valued outcome: success or failure.

Pareto Distribution: → Also called as power law.



- Eg:
- 1) 80% of wealth is distributed with 20% of people.
 - 2) 80% of sales is done by 20% of most famous product.

Anova Test (Analysis of Variance): An ANOVA test is a type of statistical test used to determine if there is a statistically significant difference between two or more categorical groups by testing for differences of means using variance.

