

Statistics

Thursday, November 7, 2024 7:42 PM

Definition:- Statistics is the science of collecting, organizing and analyzing data.

Data : "facts or pieces of information"

Ex- heights of students, IQ of students

Types of Statistics

Descriptive Statistics

Def:- It consists of organizing and summarizing data

- 1) Measures of central tendency [Mean, Median, Mode]
- 2) Measure of dispersion [Variance, standard deviation]
- 3) Different type of distribution of data



Inferential Statistics

Def:- It consists of using data you have measured to form conclusion.



- 1) Z test
 - 2) T test
 - 3) CHI square test
- } *hypothesis testing*

Ex- Let say there are 20 statistics classes at your college. And you have collected the heights of students in the class. Heights are recorded 175cm , 180cm , 140cm, 140cm, 135cm, 160cm, 135cm, 190cm.

Descriptive question : What is the average height of the entire class *which means mean*

$$175 + 180 + 140 + 140 + 135 + 160 + 135 + 190 / 8$$

Inferential question : Are the height of the students in class similar to what you expect in the entire college

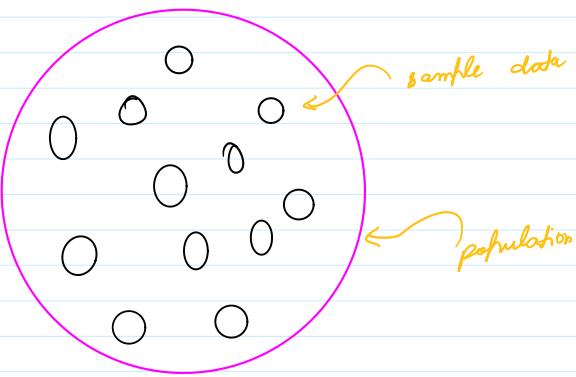
sample data

population data

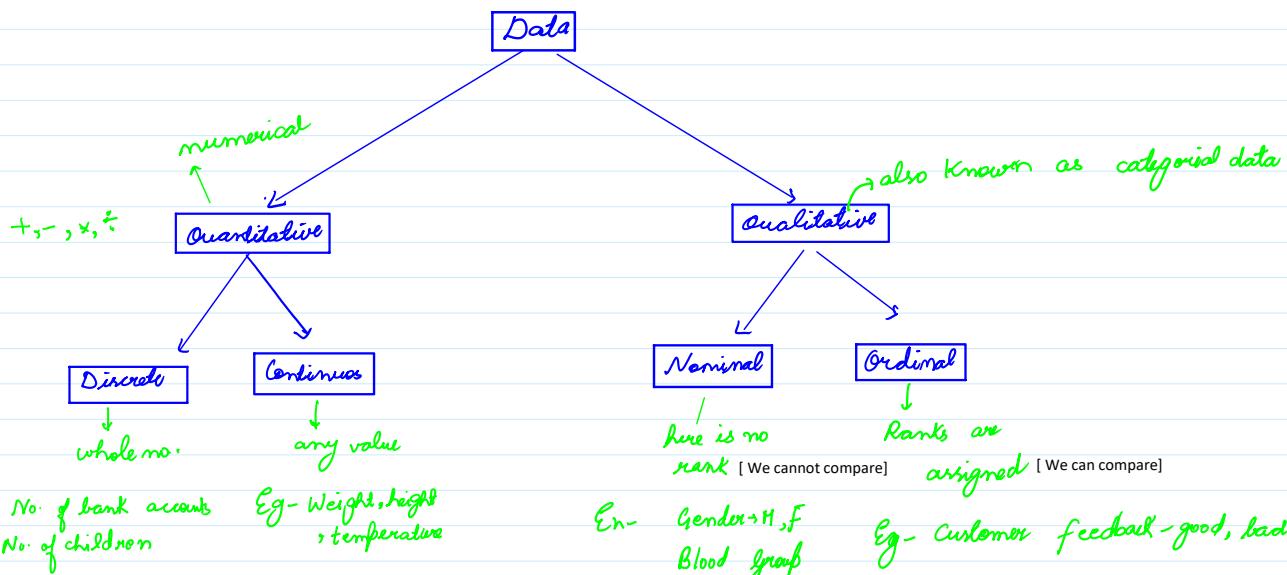
Population Data: The group you are interested in studying

Sample Data : A subset of population

Eg-



Yha hamr raye lete hai sbki or uske basis pr hamr conclusion nikalte hai



DMC	DC	ISI	BUI	FWI	Classes	Region
2.4	4.6	1.3	3.4	0.5	Fire	0
4.1	3.4	1	3.9	0.4	Not fire	1
2.6	6.9	0.3	2.7	0.4	Fire	0

Scale of Measurement:

- 1) Nominal scale data
- 2) Ordinal scale data
- 3) Interval scale data
- 4) Ratio scale data

1. Nominal Scale Data

- Includes qualitative / categorical data
- Ex- Gender, Colors
- Order does not matter

Survey of favourite color

Red → 5 Person
Blue → 3 Person
Green → 2 Person
Yellow → 1 Person

0 0 0

Red \rightarrow 5 Person

Blue \rightarrow 3 Person

Orange \rightarrow 2 Person

Order means there is no point to compare red is greater than blue and orange but we can conclude this 50% red color person opt, 30% blue color person opt and remaining 20% color orange opt

2. Ordinal Scale Data

- Ranking is important
- Order matters
- Difference cannot be measured.

Ex- App with review Best, good, Better

- Ranking \rightarrow
- 1 \rightarrow Best
 - 2 \rightarrow Good
 - 3 \rightarrow Better

Ex- In a Race

- Difference
cannot
be measured
- Ranking \rightarrow
- 1st \rightarrow Rahul
 - 2nd \rightarrow Shivam
 - 3rd \rightarrow Loveria

3. Interval Scale Data

- Order matters
- Difference can be measured
- Ratio cannot be measured
- No "0" starting point

Ex- Temperature variable

30 F°
60 F°
90 F°

$$30 : 90 = 1 : 3$$

ya hamara year wali kah skte 90f mai 3 times
garbi laggi nahi hogi 30f comparison
mai

No 0 starting point \rightarrow yeah negative se thi start hot skta ha

4. Ratio Scale Data

- Order matters
- Difference are measurable including ratio
- Contains "0" starting point

Ex- Student marks of a class

Ex- students marks of a class

Random Quiz?

- 1) Length of different rivers in the world → Ratio scale Data
- 2) favorite food based on gender → Nominal scale data

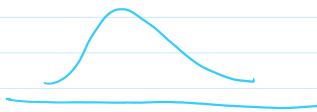
Measure of Central Tendency

- 1) Mean / Average
- 2) Median
- 3) Mode

1. Mean

Population (N)

$$X = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6\}$$



Sample (n)

$$\text{Sample} = \frac{\sum_{i=1}^n x_i}{n}$$

$$\text{Population Mean} (\mu) = \frac{\sum_{i=1}^n x_i}{N}$$

$$= \frac{1+1+2+2+3+3+4+5+5+6}{10} = 3.2$$

2. Median

$$X = \{4, 5, 2, 3, 2, 1\}$$

Step 1: Sort the Random variable $\{1, 2, 3, 4, 5, 6\}$

Step 2: No. of elements count = 6

Step 3: if count == even

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

$$\text{Median} = \frac{2+3}{2} = 2.5$$

Step 4: if count == odd

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

median

Why we use median

$$X = \{1, 2, 3, 4, 5\}$$

$$\bar{x} = \frac{1+2+3+4+5+6}{5}$$

$$X = \{1, 2, 3, 4, 5, 100\}$$

→ This large no. is called outlier

$$\bar{x} = \frac{1+2+3+4+5+100}{6}$$

$$\bar{x} = \frac{1+2+3+4+5+6}{5}$$

$$= 3$$

Important

Yeah kha large no. 100 add Kone par Kihna zyada shift hoh zya mean mai skdum isliye study median

$$\bar{x} = \frac{1+2+3+4+5+100}{6}$$

$$\approx 19$$

$$X = \{1, 2, 3, 4, 5\}$$

$$\text{Median} = 3$$

$$X_2 = \{1, 2, 3, 4, 5, 100\}$$

$$\text{Median} = \frac{3+4}{2} = 3.5$$

$$3 \rightarrow 3.5$$

→ not large shift in median after adding outliers

➤ Median is used to find the central tendency when outliers is present

3. Mode → Maximum frequency

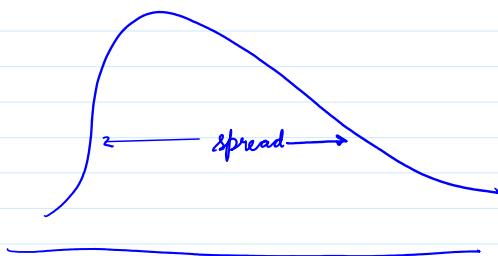
$$X = \{1, 1, 1, 7, 2, 3\}$$

$$\text{Mode} = 1$$

➤ We are learning mean, median and mode toh yeah kha prr use hogia industry main EDA and feature Engineering

Go through the Measure of Central Tendency Code In google Colab

Measure of Dispersion [Spread of the data]



1. Variance

Population variance

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

$x_i \rightarrow$ Data points

sample Variance

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

Ques : Why we divide Sample variance by $n-1$?
Ans. So that we can create an unbiased estimator of the population variance.

\downarrow
this scenario called
 $n-1$ degrees of freedom

$x_i \rightarrow$ Data points
 $\bar{x} \rightarrow$ Sample mean

$x_i \rightarrow$ Data points
 $\mu \rightarrow$ Population mean
 $N \rightarrow$ Population size

Variable
 this scenario called
Bessel correction
 $x_i \rightarrow$ Data points
 $\bar{x} \rightarrow$ Sample mean
 $n \rightarrow$ Sample size

Ex - $\{1, 2, 3, 4, 5\}$

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

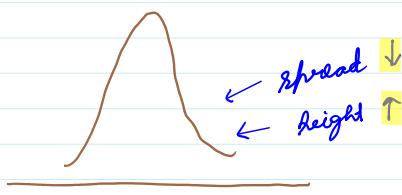
$$\bar{x} = \frac{1+2+3+4+5}{5} = 3$$

x_i	\bar{x}	$(x_i - \bar{x})^2$
1	3	4
2	3	1
3	3	0
4	3	1
5	3	4
		10

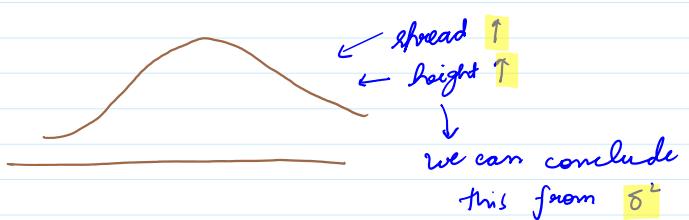
$$\delta^2 = \frac{10}{5-1} = \frac{10}{4} = 2.5$$

what does means let us understand

$X = \{ \}$
 $\delta^2 = 2.5$



$Y = \{ \}$
 $\delta^2 = 7.5$



2. Standard Deviation

Population standard deviation

$$\sigma = \sqrt{\text{Variance}}$$

Sample standard deviation

$$std = \sqrt{s^2}$$

$$S^2 = \text{sample variance}$$

Explore python codes based on variance, standard deviation,

Random Variable

Random Variable is a process of mapping the output of a random process or experiments to a number

Ex - Tossing a coin
 yeah result has head/tail outcome ka
 $\vee \quad | \quad n \quad : 1 \quad H \quad T \quad 1$

Ex- Tossing a coin
 yeah result has Head/Tail outcome now

$$X = \begin{cases} 0 & \text{if Head} \\ 1 & \text{if Tail} \end{cases}$$

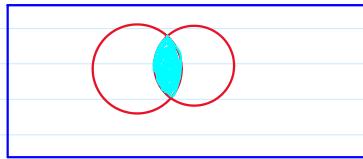
Sets

$$A = \{1, 2, 3, 4, 5, 6, 7, 8\}$$

$$B = \{3, 4, 5, 6, 7\}$$

1. Intersection

$$A \cap B = \{3, 4, 5, 6, 7\}$$

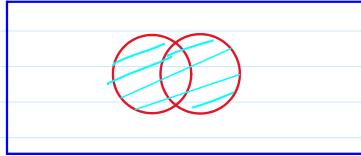


2. Union

$$A = \{1, 2, 3, 4, 5, 6, 7, 8\}$$

$$B = \{3, 4, 5, 6, 7, 9, 10\}$$

$$A \cup B = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$$

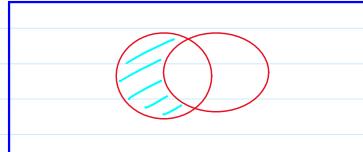


3. Difference

$$A = \{1, 2, 3, 4, 5, 6, 7, 8\}$$

$$B = \{3, 4, 5, 6, 7\}$$

$$A - B = \{1, 2, 8\}$$



4. Subset

$A \subset B \times$
 $B \subset A \checkmark$

5. Superset

$A \supset B$

$B \supset A$

Histogram And Skewness

Histogram is related to frequency.

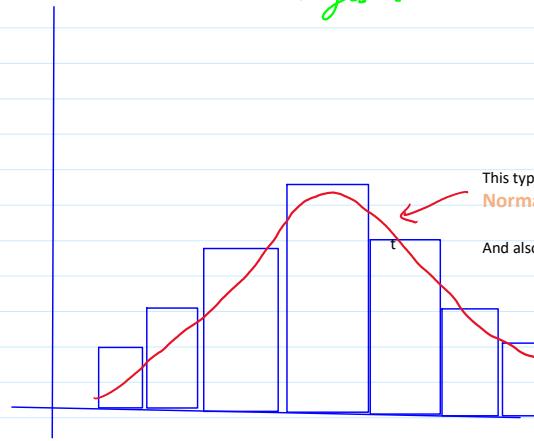
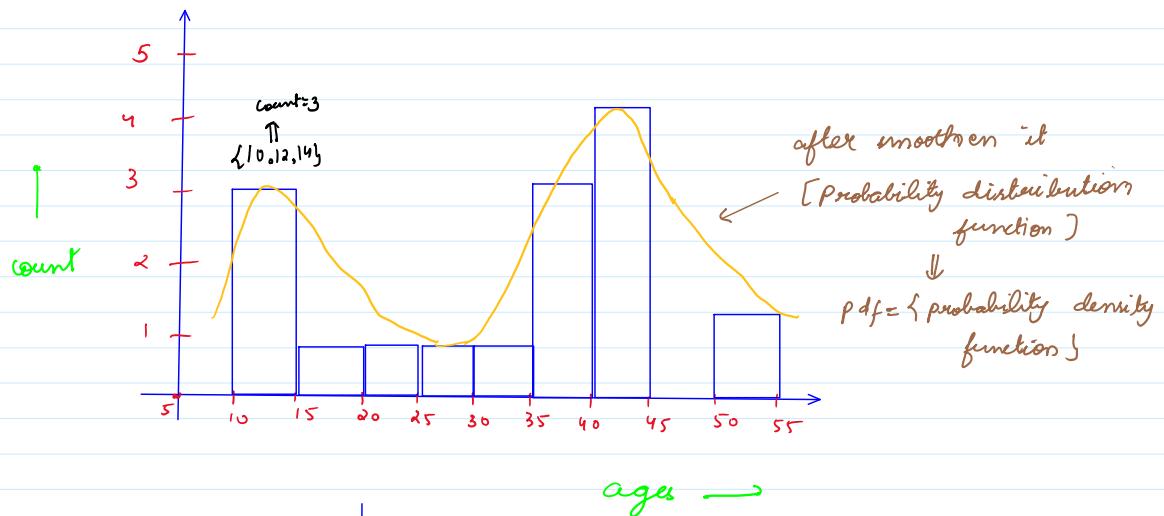
Ages = {10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51}

$$\frac{50}{10} = 5 \rightarrow \text{bin size}$$

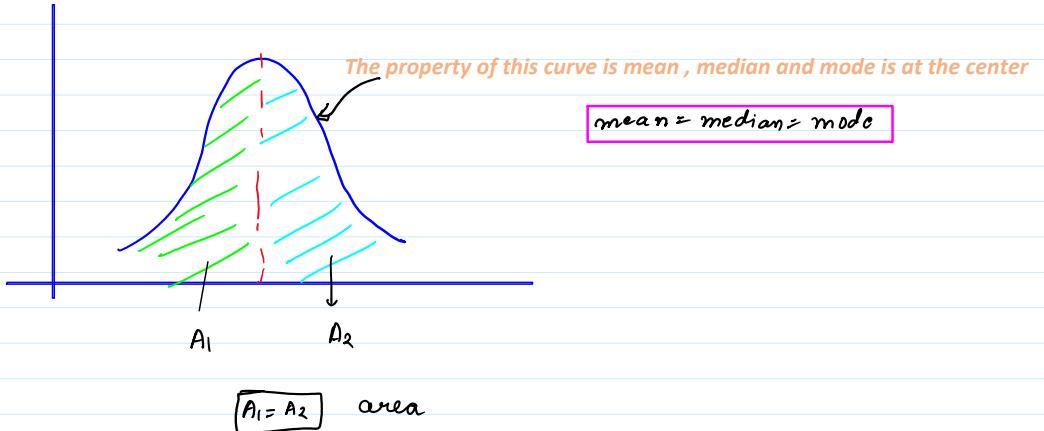
{ No. of bins = 10 }

$$\frac{50}{20} = 2.5 \text{ bin size}$$

{ No. of bins = 20 }



1. Superset

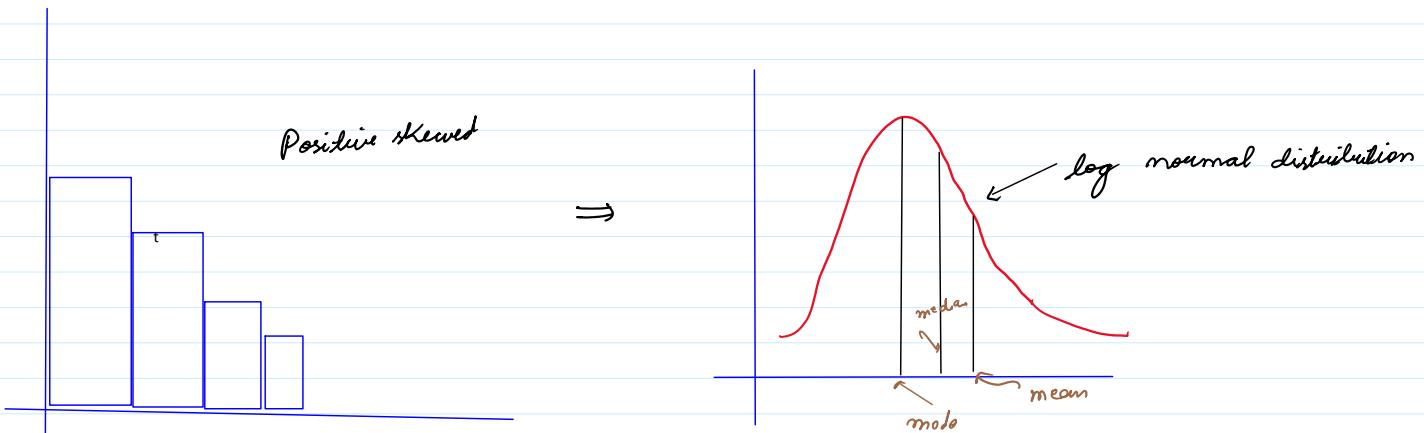


Bon plot: used for finding outliers



$$Q_3 - Q_2 \approx Q_2 - Q_1$$

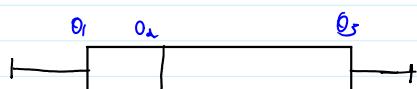
2. Right Skewed



In this what is the relation b/w mean, median, mode

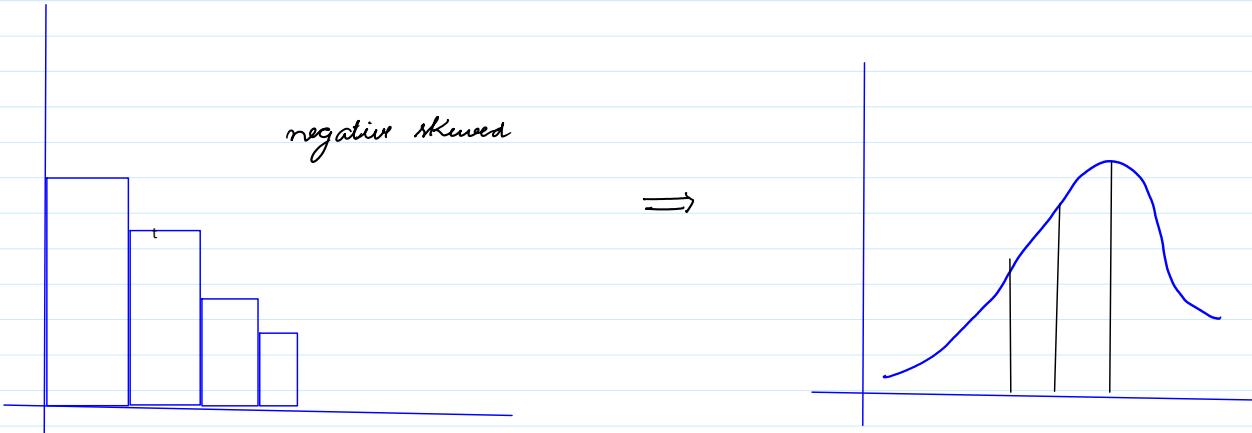
$$\text{mean} > \text{median} > \text{mode}$$

Bon plot:



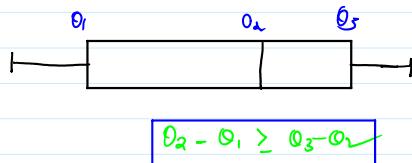
$$Q_3 - Q_2 \geq Q_2 - Q_1$$

3. Left Skewed/Negative Skewed



Mean < median ≤ mode

Bonferroni:



Covariance And Correlation

X	Y
2	3
4	5
6	7
8	9

What is the Relationship between X & Y

x↑	y↑
x↑	y↓
x↓	y↑
x↓	y↓

← we can find this relationship / dependency using correlation & covariance

Covariance

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

$x_i \rightarrow$ Data points of x

$\bar{x} \rightarrow$ Sample mean

$y_i \rightarrow$ Data points of y

$\bar{y} \rightarrow$ Sample mean of y

$$\text{var}(n) = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1}$$

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

= cov(n, n) ⇒ it will tell about spread

↓
Observe the covariance and try to find the relationship

cov(x, y)

$$\begin{array}{|c|c|} \hline x↑ & y↑ \\ \hline \end{array} \rightarrow +ve$$

$$\begin{array}{|c|c|} \hline x↑ & y↓ \\ \hline \end{array} \rightarrow -ve \text{ covariance}$$

$\text{Cov}(x, y)$

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

+ve covariance

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

-ve covariance

Ques :

$$\begin{array}{r} x \\ 2 \\ 4 \\ 6 \\ \hline \bar{x} = 4 \end{array} \quad \begin{array}{r} y \\ 3 \\ 5 \\ 7 \\ \hline \bar{y} = 5 \end{array}$$

$$\begin{aligned} \text{cov}(x, y) &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1} \\ &= \frac{[(2-4)(3-5) + (4-4)(5-5) + (6-4)(7-5)]}{3-1} \end{aligned}$$

= 9
+ve

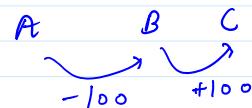
Conclusion $\rightarrow x \& y$ having a positive covariance

Advantages

- Relationship between X and Y +ve or -ve value

Disadvantages

- Covariance does not have a specific limit [Interview Question]



To find the specific limit of covariance we used :-

2. Pearson Correlation Coefficient [-1, 1] output

$$P_{x,y} = \frac{\text{cov}(x, y)}{\sigma_x \cdot \sigma_y}$$

- The more the value towards +1 the more +ve correlation it is (x, y).
- The more the value towards -1 the more -ve correlated it is (x, y)

3. Spearman Rank Correlation [-1, +1]

$$r_s = \frac{\text{cov}(R(x), R(y))}{\sigma(R(x)) \cdot \sigma(R(y))}$$

Rank

Rank means more higher the value more higher the rank

x	y	$R(x)$	$R(y)$
1	2	5	5
3	4	4	4
5	6	3	3
7	8	2	1
0	7	6	2
8	1	1	6

Where we really used these above in Machine Learning i.e In **Feature Selection**

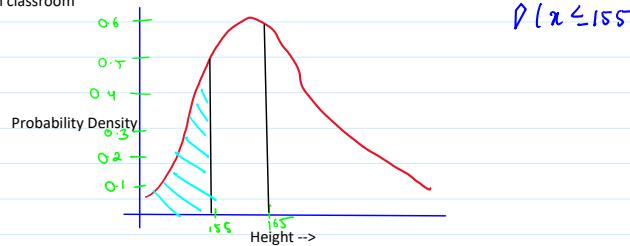
Code

```
import seaborn as sns  
df=sns.load_dataset('healthexp') #Load_dataset is a method in java and healthexp is a inbuilt database  
df.head()  
  
import numpy as np  
np.cov() # cov mean covariance  
##Correlations using Pearson correlation coefficient  
df.corr(method='pearson')
```

Probability Distribution Function / Density Function

1. Probability Density Function [PDF]

Ex- Heights of students in classroom



2. Continuous Random Variable

Variable → Discrete Random variable

Ex- Rolling a Dice $\{1, 2, 3, 4, 5, 6\}$

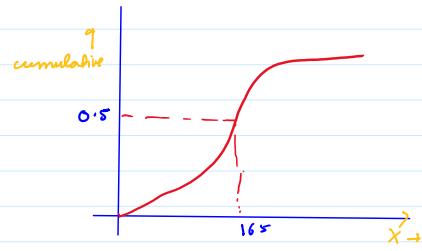
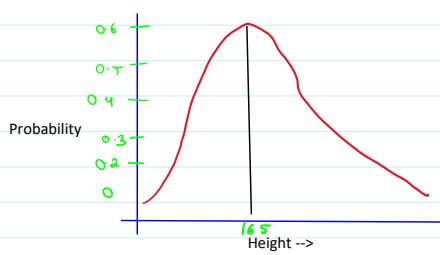


$$P_n(1) = \frac{1}{6} \quad P_n(2) = \frac{1}{6}$$

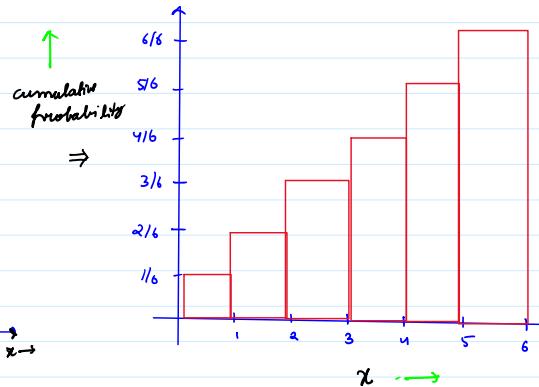
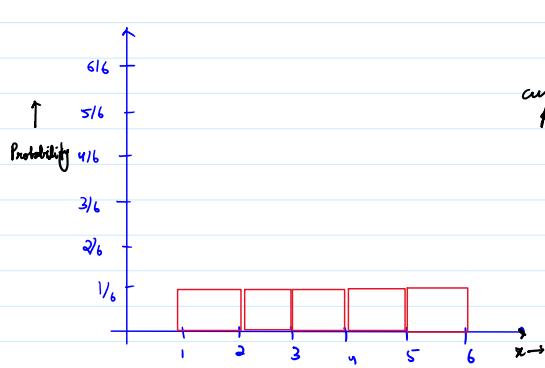
$$P_n(x \leq 4) = P_n(x=1) + P_n(x=2) + P_n(x=3) + P_n(x=4)$$

$$= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{4}{6} = \frac{2}{3}$$

3. Cumulative Distribution Function [CDF]

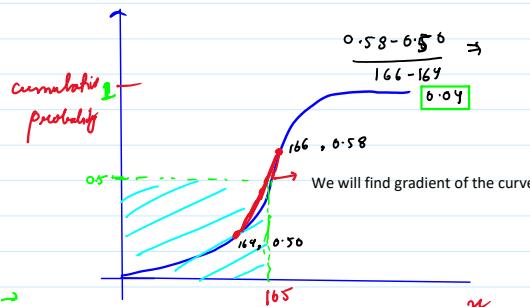
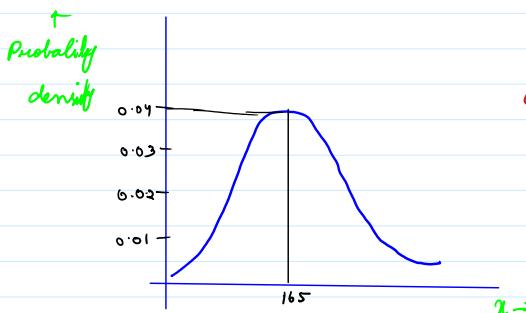


Rolling a Dice $\Rightarrow \{1, 2, 3, 4, 5, 6\}$



Probability distribution Function

① Distribution of continuous random variable



Probability density \Rightarrow Gradient of Cumulative Curve

Types of Probability Distribution

1. Normal / Gaussian Distribution [pdf]
2. Bernoulli Distribution [pmf]
3. Uniform Distribution [pmf]
4. Poisson Distribution [pmf]
5. Log Normal Distribution [pdf]
6. Binomial Distribution [pmf]

➤ Bernoulli Distribution

In probability theory and statistics, the Bernoulli distribution, named after Swiss mathematician Jacob Bernoulli, is the discrete probability distribution of a random variable which takes the value 1 with probability p and the value 0 with probability $q=1-p$. Less formally, it can be thought of as a model for the set of possible outcomes of any single experiment that asks a yes-no question.

→ Discrete Random variable (pmf)

a yes-no question.

→ Discrete Random variable (pmf)

→ outcomes are binary

Ex - Tossing a coin {H, T}

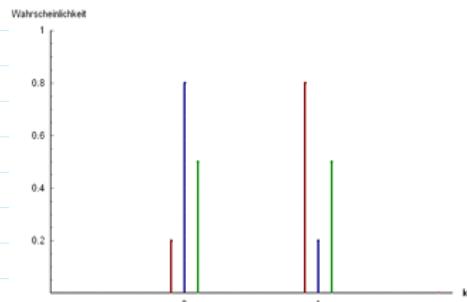
$$P(X(\text{head})=0.5 = p)$$

$$P(X(\text{tail})=1-p=q)$$

② whether the person will Pass / fail

$$P(X(\text{Pass})=0.7 = p)$$

$$P(X(\text{Fail})=1-p=0.3 = q)$$



$0 \leq p \leq 1$

$$q = 1 - p$$

$$K = \{0, 1\}$$

Probability Mass function

$$PMF = p^k * (1-p)^{1-k} \quad k \in \{0, 1\}$$

if ($k=1$)

$$P(X=k=1) = p^1 * (1-p)^{1-1}$$

$$= p$$

if $K=0$

$$P(X(K=0)) = p^0 * (1-p)^{1-0}$$

$$= (1-p) = q$$

We can write this in simplified Manner

$$\text{Probability} = \begin{cases} q = 1-p & \text{if } K=0 \\ p & \text{if } K=1 \end{cases}$$

➤ Mean of Bernoulli Distribution

$$E(K) = \sum_{k=0}^K k \cdot p(k)$$

➤ Median of Bernoulli Distribution

$$\text{Median} = \begin{cases} 0 & \text{if } p \leq \frac{1}{2} \\ [0, 1] & \text{if } p = \frac{1}{2} \\ 1 & \text{if } p > \frac{1}{2} \end{cases}$$

➤ Variance of Bernoulli Distribution

$$\text{Var} = p * (1-p) \\ = pq$$

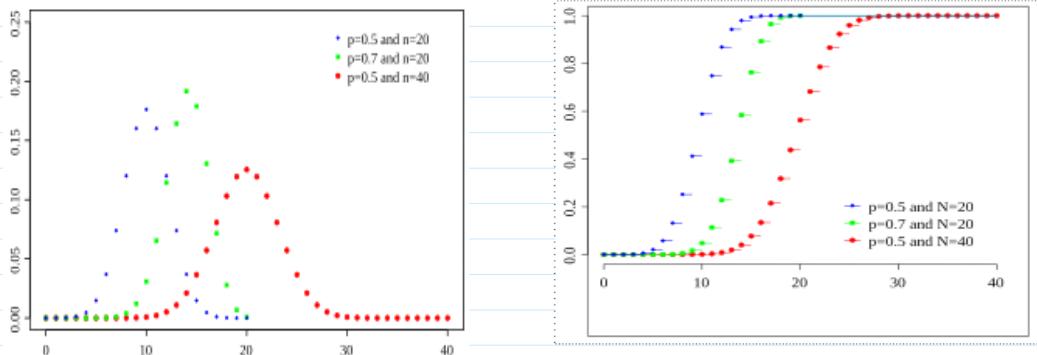
Standard Deviation of Bernoulli Distribution

$$\text{std} = \sqrt{\text{Var}} = \sqrt{pq}$$

➤ Binomial Distribution

In probability theory and statistics, 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 (with probability p) or failure $q=1-p$. A single success/failure experiment is also called a Bernoulli experiment, and a sequence of outcomes is called a Bernoulli process; for a single trial, i.e. $n=1$, the binomial distribution is a Bernoulli distribution. The binomial distribution is the basis for the popular binomial test of statistical significance.



- Discrete random variable
- Every experiment outcome is binary
- These experiment is performed for n trials.

Notation : $B(n, p)$

Parameters: $n \in \{0, 1, 2, 3, \dots\} \rightarrow$ no. of trials

$P \in [0, 1] \rightarrow$ success probability for each trial

$$q = 1 - p$$

Output: $K \in \{0, 1, 2, \dots, n\} \rightarrow$ No. of success

$$P_n(K, n, p) = {}^n C_k p^k (1-p)^{n-k}$$

Mean of Binomial Distribution

$$\text{Mean} = np$$

Variance of Binomial Distribution

$$\text{Var} = npq$$

Standard deviation of Binomial Distribution

$$sd = \sqrt{pq}$$

➤ Poisson Distribution

→ Discrete Random variables (Pmf)

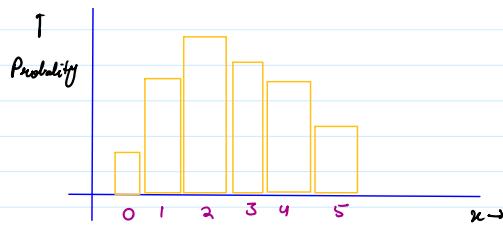
→ Describes the number of events occurring in a fixed time interval.

Ex

Ex- No. Of People visiting hospital/bank/airport every hour

$\lambda = 3$ Expected no. of event occur at every time interval





formula:-

$$p(x) = \frac{e^{-\lambda} \lambda^x}{x!}$$

Ques : Based on the graph what is the probability that the patient will come at 5th hour ?

$$\begin{aligned} p(x=5) &= \frac{e^{-3} 3^5}{5!} \\ &= 10.1\% \end{aligned}$$

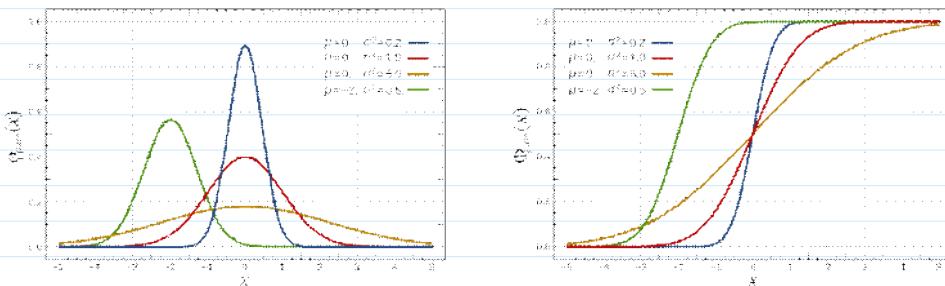
Mean of Poisson Distribution

$$\text{Mean} = E(x) = \lambda t$$

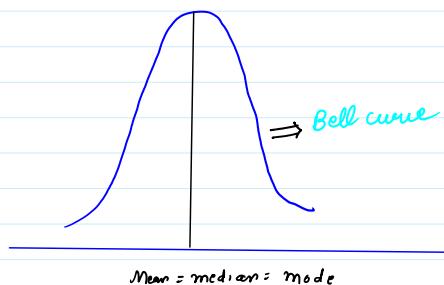
λ = Expected no. of events to occur at every time interval
t = Time interval
notation

Normal/Gaussian Distribution [pdf]

In statistics, a normal distribution or Gaussian distribution is a type of continuous probability distribution for a real-valued random variable



In this curve is most symmetrical



$$\text{PDF} = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma} \right)^2}$$

Mean of Normal Distribution

Mean = μ : Average

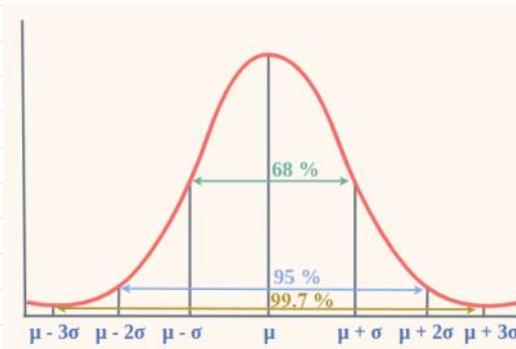
Variance of Normal Distribution

$$\text{Var} = \sigma^2$$

Standard deviation of Normal Distribution

$$\sigma = \sqrt{\text{Var}}$$

Empirical Rule for normal distribution \rightarrow V. Imhardt \rightarrow Intervall



The **empirical rule**, also known as the **68-95-99.7 rule**, describes where most values lie in a normal distribution:

- Around **68%** of values are within **1 standard deviation** from the mean.
- Around **95%** of values are within **2 standard deviations** from the mean.
- Around **99.7%** of values are within **3 standard deviations** from the mean.

Probability

$$P_n(\mu - \sigma \leq n \leq \mu + \sigma) \approx 68\%$$

$$P_n(\mu - 2\sigma \leq n \leq \mu + 2\sigma) \approx 95\%$$

$$P_n(\mu - 3\sigma \leq n \leq \mu + 3\sigma) \approx 99.7\%$$

Ex- Weight of students
Heights of students

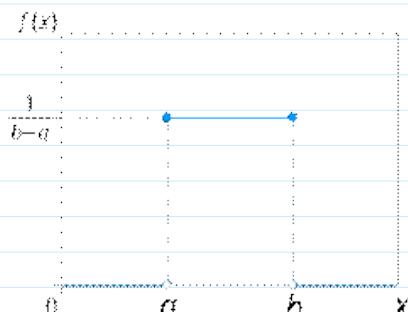
Uniform distribution

1. Continuos Uniform Distribution

2. Discrete Uniform Distribution

1. Continuos Uniform distribution

In probability theory and statistics, the continuos uniform distribution or rectangular distribution is a family of symmetric probability distribution. The distribution describes an experiment where there is an arbitrary outcome that lies between certain bounds. The bounds are defined by the parameters, a and b , which are the minimum and maximum values.



Notation :- $U(a,b)$

Parameters :- $-\infty < a < b < \infty$

$$\text{pdf} = \begin{cases} \frac{1}{b-a}, & \text{for } x \in [a, b] \\ 0, & \text{otherwise} \end{cases}$$



$$\text{Mean} = \frac{1}{2}(a+b)$$

$$\text{Median} = \frac{1}{2}(a+b)$$

$$\text{Variance} = \frac{1}{12}(b-a)^2$$

Ex- The no. of candies sold daily at a shop is uniformly distributed with a maximum of 40 and a minimum of 10.

o) What is the probability of daily sales to fall between 15 and 30?

ans :-

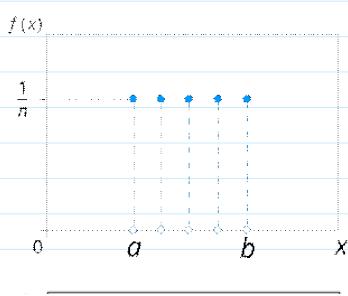


$$\begin{aligned} P(15 \leq n \leq 30) &= (x_2 - x_1) \times \frac{1}{b-a} \\ &= 30 - 15 \times \frac{1}{30} \\ &= 0.5 \approx 50\% \end{aligned}$$

$$P(n \geq 20) = 40 - 20 \times \frac{1}{30} = \frac{20}{30} = 0.6$$

2. Discrete Uniform distribution (pmf) ✓

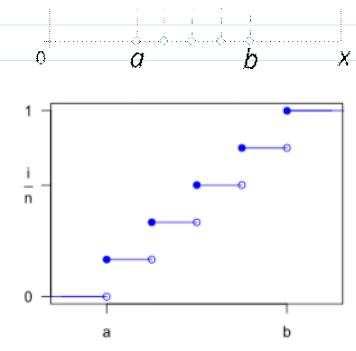
In probability theory and statistics, the discrete uniform distribution is a symmetric probability distribution wherein each of some finite whole number n of outcome values are equally likely to be observed. Thus every one of the n outcome values has equal probability $1/n$. Another way of saying "discrete uniform distribution" would be "a known, finite number of outcomes equally likely to happen".



Ex- Rolling a dice = {1, 2, 3, 4, 5, 6}

$$\frac{1}{n} \Rightarrow n = b - a + 1$$

Notations :-



Notation :-

Parameters a, b with $b \geq a$

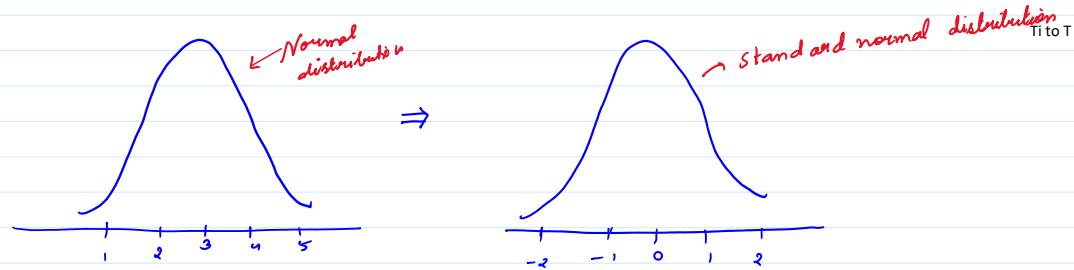
$$PMF = \frac{1}{n}$$

$$\text{Mean} = \frac{a+b}{2}$$

$$\text{Median} = \frac{a+b}{2}$$

Standard Normal distribution and Z-Score [z-stats]

$X = \{1, 2, 3, 4, 5\}$ Normally distributed assume
 Random variable $\mu = \text{Mean} = 3$
 $\sigma = 1.411 \approx 1$ [assume]



How to convert normal distribution to standard normal distribution uske liye use hota hai z-score

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

$$= \frac{1-3}{1} = -2$$

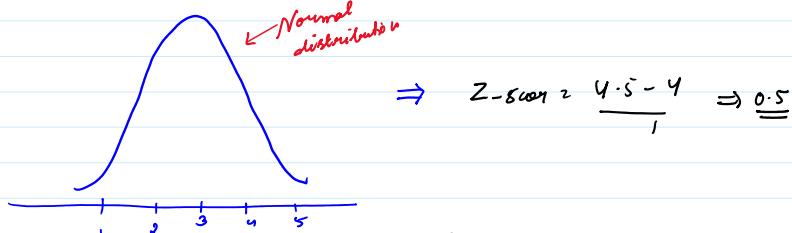
$$= \frac{2-3}{1} = -1$$

$$3 \rightarrow -2$$

graph 1 graph 2

$$2 \rightarrow -1$$

Ques : What percentage of data is falling above 4.5 with $\mu = 4$ and $\sigma = 1$



Area under the curve (> 4.5)

$$= 1 - 0.69146$$

from table

z	0	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
-0	.5000	.49601	.49202	.48803	.48405	.48006	.47608	.47210	.46812	.46414
-0.1	.46017	.45620	.45224	.44828	.44433	.44034	.43640	.43251	.42858	.42465
-0.2	.42074	.41683	.41294	.40905	.40517	.40129	.39743	.39358	.38974	.38591
-0.3	.38209	.37828	.37448	.37070	.36693	.36317	.35942	.35569	.35197	.34827
-0.4	.34458	.34090	.33724	.33360	.32997	.32636	.32276	.31918	.31561	.31207
-0.5	.30854	.30503	.30153	.29800	.29460	.29116	.28774	.28434	.28096	.27760
-0.6	.27425	.27093	.26763	.26435	.26109	.25785	.25463	.25143	.24825	.24510
-0.7	.24196	.23885	.23576	.23270	.22965	.22663	.22363	.22065	.21770	.21476
-0.8	.21186	.20897	.20611	.20327	.20045	.19766	.19489	.19215	.18943	.18673
-0.9	.18406	.18141	.17879	.17619	.17361	.17106	.16853	.16602	.16354	.16109
-1	.15866	.15625	.15386	.15151	.14917	.14686	.14457	.14231	.14007	.13786
-1.1	.13567	.13350	.13136	.12924	.12714	.12507	.12302	.12100	.11900	.11702

<i>z</i>	0	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
-0	.50000	.49601	.49202	.48803	.48405	.48006	.47608	.47210	.46812	.46414
-0.1	.46017	.45620	.45224	.44828	.44433	.44034	.43640	.43251	.42858	.42465
-0.2	.42074	.41683	.41294	.40903	.40517	.40129	.39743	.39358	.38974	.38591
-0.3	.38209	.37828	.37448	.37070	.36693	.36317	.35942	.35569	.35197	.34827
-0.4	.34458	.34090	.33724	.33360	.32997	.32636	.32276	.31918	.31561	.31207
-0.5	.30854	.30503	.30153	.29806	.29460	.29116	.28774	.28434	.28096	.27760
-0.6	.27425	.27093	.26763	.26435	.26109	.25785	.25463	.25143	.24825	.24510
-0.7	.24196	.23885	.23576	.23270	.22965	.22663	.22363	.22065	.21770	.21476
-0.8	.21186	.20897	.20611	.20327	.20045	.19766	.19486	.19215	.18943	.18673
-0.9	.18406	.18141	.17879	.17619	.17361	.17106	.16853	.16602	.16354	.16109
-1	.15866	.15625	.15386	.15151	.14917	.14686	.14457	.14231	.14007	.13786
-1.1	.13567	.13350	.13136	.12924	.12714	.12507	.12302	.12100	.11900	.11702
-1.2	.11507	.11314	.11123	.10935	.10749	.10565	.10383	.10204	.10027	.09853
-1.3	.09680	.09510	.09342	.09176	.09012	.08852	.08692	.08534	.08379	.08226
-1.4	.08076	.07927	.07780	.07636	.07493	.07353	.07215	.07078	.06944	.06811
-1.5	.06681	.06552	.06426	.06301	.06178	.06057	.05938	.05821	.05705	.05592
-1.6	.05480	.05370	.05262	.05155	.05050	.04947	.04846	.04746	.04648	.04551
-1.7	.04457	.04363	.04272	.04182	.04093	.04006	.03920	.03836	.03754	.03673
-1.8	.03593	.03515	.03438	.03362	.03288	.03216	.03144	.03074	.03005	.02938
-1.9	.02872	.02807	.02743	.02680	.02619	.02559	.02500	.02442	.02385	.02330
-2	.02275	.02222	.02169	.02118	.02068	.02018	.01970	.01923	.01876	.01831
-2.1	.01786	.01743	.01700	.01659	.01618	.01578	.01539	.01500	.01463	.01426
-2.2	.01390	.01355	.01321	.01287	.01255	.01222	.01191	.01160	.01130	.01101
-2.3	.01072	.01044	.01017	.00990	.00964	.00939	.00914	.00889	.00866	.00842
-2.4	.00820	.00794	.00776	.00755	.00734	.00714	.00695	.00676	.00657	.00639
-2.5	.00621	.00604	.00587	.00570	.00554	.00539	.00523	.00508	.00494	.00480
-2.6	.00466	.00453	.00440	.00427	.00415	.00402	.00391	.00379	.00368	.00357
-2.7	.00347	.00336	.00326	.00317	.00307	.00298	.00289	.00280	.00272	.00264
-2.8	.00256	.00248	.00240	.00233	.00226	.00219	.00212	.00205	.00199	.00193
-2.9	.00187	.00181	.00175	.00169	.00164	.00159	.00154	.00149	.00144	.00139
-3	.00135	.00131	.00126	.00122	.00118	.00114	.00111	.00107	.00104	.00100
-3.1	.00097	.00094	.00090	.00087	.00084	.00082	.00079	.00076	.00074	.00071
-3.2	.00069	.00066	.00064	.00062	.00060	.00058	.00056	.00054	.00052	.00050
-3.3	.00048	.00047	.00045	.00043	.00042	.00040	.00039	.00038	.00036	.00035
-3.4	.00034	.00032	.00031	.00030	.00029	.00028	.00027	.00026	.00025	.00024
-3.5	.00023	.00022	.00022	.00021	.00020	.00019	.00019	.00018	.00017	.00017
-3.6	.00016	.00015	.00015	.00014	.00014	.00013	.00013	.00012	.00012	.00011
-3.7	.00011	.00010	.00010	.00010	.00009	.00009	.00008	.00008	.00008	.00008
-3.8	.00007	.00007	.00007	.00006	.00006	.00006	.00006	.00005	.00005	.00005
-3.9	.00005	.00005	.00004	.00004	.00004	.00004	.00004	.00004	.00003	.00003
-4	.00003	.00003	.00003	.00003	.00003	.00003	.00002	.00002	.00002	.00002

from table

Real World Problem : In India the average IQ is 100, with a standard deviation of 15 . What is the percentage of the population would you expect to have an IQ lower than 85%

$$\mu = 100$$

$$\sigma = 15$$

$$x_i =$$



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

$$= \frac{85 - 100}{15}$$

$$= -1$$

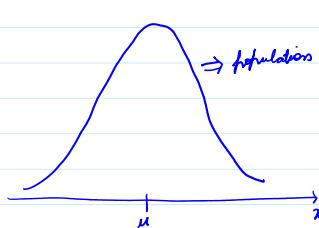
$$\text{Area} = 0.15866$$

$$= 15.866$$

Central Limit Theorem [Interview Ques]

The central limit theorem says that the sampling distribution of the mean will always be normally distributed, as long as the sample size is larger enough. Regardless of whether the population has a normal, Poisson, binomial , or any other distribution, the sampling distribution of the mean will be normal .

[Interview Question]



$$X \sim N(\mu, \sigma)$$

σ = population std

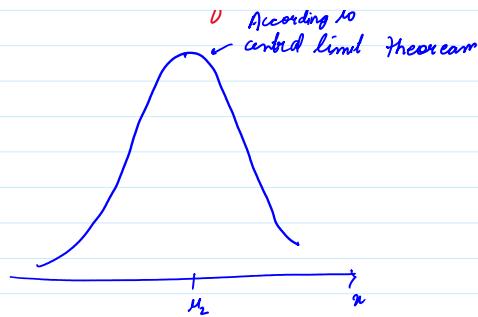
μ = population mean

n = sample size



$$\bar{x} \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

sample distribution of mean
According to
central limit theorem



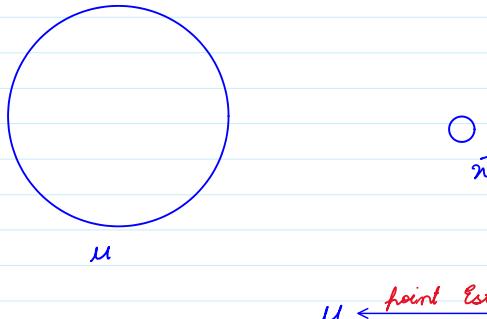
$$\bar{x} \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

n can be any value

Estimate It is an observed numerical value used to estimate an unknown population parameter

1. Point Estimate Single numerical value used to estimate the unknown population parameter.

Ex- Sample mean is a point estimate of population mean.

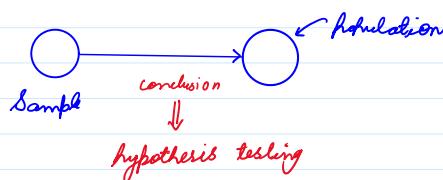


2. Interval Estimate Range of values used to estimate the unknown population parameters. Interval estimates of population parameters are called confidence intervals.



A hypothesis and Hypothesis testing Mechanism

Inferential stats :- conclusion or inferences



Hypothesis testing Mechanism

- ① Null hypothesis (H_0) → The assumption you are beginning with
- ② Alternate hypothesis :- Opposite of null hypothesis
- ③ Experiments → proof collect
- ④ Accept the null hypothesis or reject the null hypothesis

Ex- Colleges at District A states its passed average percentage of students are 85%. A new college opened in the district and it was found that a sample of student 100 have a pass percentage of 90% with a standard deviation of 4%. Does this school have a different passed percentage?

What is null and alternate hypothesis

null hypothesis $H_0 = 85$ (by default)
 alternate hypothesis $H_1 \neq 85$

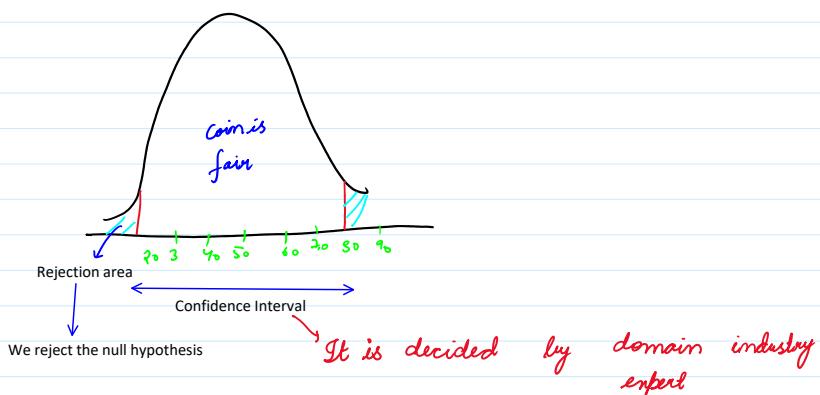
P value

The p value is a number, calculated from a statistical test, that describes how likely you are to have found a particular set of observations if the null hypothesis were true. P values are used in hypothesis testing to help decide whether to reject the null hypothesis.

Ex- Tossing a fair coin

- 1) Null Hypothesis \rightarrow coin is fair
- 2) Alternate hypothesis \rightarrow coin is not fair

Experiment:



If $p \leq$ significance value \rightarrow reject the null hypothesis

Hypothesis testing And Statistical Analysis

- 1) Z test \Rightarrow Average
- 2) t test \Rightarrow Variance
- 3) chi square \Rightarrow Categorical data
- 4) ANOVA \Rightarrow Variance

Z test: Used when \rightarrow (i) population std given
 (ii) $n \geq 30$

1. The average heights of all residents in a city is 168cm. A doctor believes the mean to be different. He measured the height of 36 individuals and found the average height to be 169.5cm.

- State null and alternate hypothesis
- At a 95% confidence level, is there enough evidence to reject the null hypothesis.

given $\mu = 168\text{ cm}$, $\sigma = 3\text{ cm}$, $n = 36$, $\bar{x} = 169\text{ cm}$

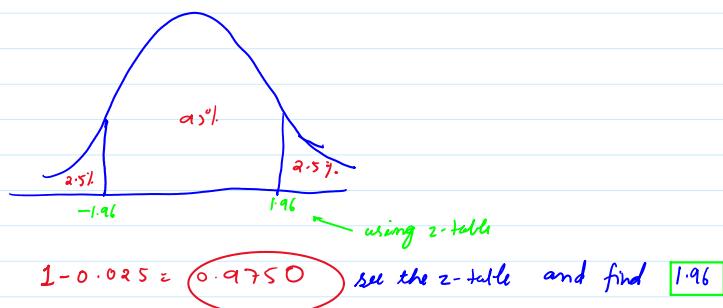
Null hypothesis (H_0) = $\mu = 168\text{ cm}$

alternate hypothesis (H_1) = $\mu \neq 168\text{ cm}$

$$CI = 0.95$$

$$\alpha = 1 - 0.95 = 0.05$$

Decision boundary



Statistical analysis

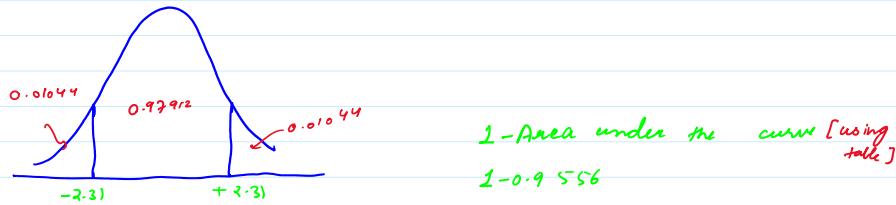
① By z-test

$$Z_{\text{test}} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{169.5 - 168}{\frac{3.9}{\sqrt{36}}} = 2.31$$

it should be in range of 1.96

↓
So we reject null hypothesis

② By p-value



$$p \text{ value} = 0.01044 + 0.01044 = 0.02088$$

if p-value < significance value

$0.02088 < 0.05 \rightarrow \text{so reject}$