

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



10

Decision Making

Data - facts or pieces of information

Measured
collected
Analyzed

Eg: Weight of students, IQ of students

Eg: House Price Dataset

City	Area	No of Rooms	Price
Hamburg	1000	2	450€
Berlin	1250	2.5	500€

Analyze

Data Scientist → Model - Price

Data Analyst → Report → Visualization - Meaningful Decision

Importance of statistics

- 1) Data Exploration
- 2) Model Building and Validation
- 3) Statistical Analysis → [Sample data \Rightarrow Conclusions]
- 4) Hypothesis Testing
- 5) Optimization and efficiency.
- 6) Report Making

Types of Statistics

1) Descriptive Statistics :- Descriptive statistics involve methods for summarizing and organizing data to make it understandable. This type of statistics helps to describe the basic features of data in a study.

1. Measure of Central Tendency

mean, median, mode

• ..

2. Measure of Dispersion

Variance, Standard Deviation

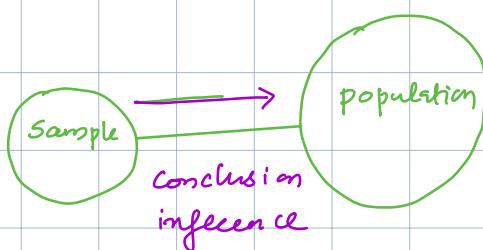
3. Data Distribution

Histogram, boxplot, Pie chart, PDF, PMF

4 Summary Statistics

Five number summary $\rightarrow Q_1, Q_2, Q_3, \text{ Max Value}$

2) Inferential Statistics : Inferential statistics involves methods for making predictions or inferences about a population based on a sample of data. It allows for hypothesis testing, estimation and drawing conclusions.



1) Hypothesis Testing

2) P value

3) Confidence Interval

4) Statistical analysis Test

① z test-

② t test

③ ANOVA \rightarrow F test

④ Chi Square

⑤ Regression Analysis

Example

Let say there are twenty statistics class in your university. You have collected the height of students in the class.

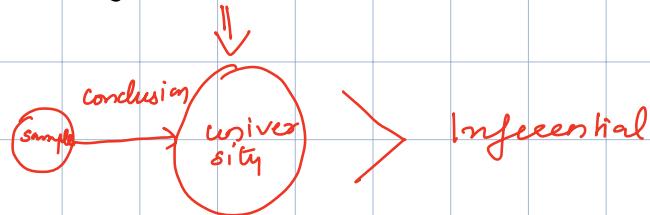
Heights are recorded as [175cm, 135cm, 180cm, 140cm - 120cm]

Descriptive Question - What is the average height of entire class?



Mean = measure of Central Tendency

Inferential Questions . Are the height of the sample students in class room similar to what you expect from the entire university?



Population and Sample Data

Population

| Sample Data

entire set of individuals or objects of interest in a particular study. It includes all members of defined group that we are studying or collecting information on

characteristics

1) complete set : contains all the observations of interest

2) parameter : A numerical value summarizing the entire population.

Eg: [population mean, population variance]
 μ , σ^2

Example 1) population in school study

- * All students in school
- * Determine the avg height of students = population mean

2) Population in market research

- *) All consumers in a city
- *) To understand the purchasing behavior of all consumers

3) Population in a medical study

- * All patients with a specific disease
- * To study the effectiveness of a drug

subset of the population that is used to represent entire group.

sampling - selecting a group of individuals / obs from the population to draw conclusion

characteristics

1) subset : Represent a portion of the population

2) statistic : A numerical value summarizing the sample data

[sample mean, sample variance]

3) Random Sampling :

samples should be randomly selected to avoid bias

Examples

1) Sample in school study

- * Group of 50 students from school

need care = estimate the avg height of students in a school

2) Market Research sample

- * 500 consumers
- * Behavior \rightarrow population generalisation

3) Medical Study Sample

- * 200 patients

* test effectiveness of the drug.

(4)

Types of Sampling Techniques

Eg:- After election an exit poll is conducted



- 1) Probability Sampling
- 2) Non Probability Sampling

1. Probability Sampling

a) simple random sampling

every member of the population has an equal chance of being selected

Eg: selecting sample randomly

drawing names randomly from a class of students

b) systematic sampling

select every n^{th} member of the population after a random starting point.

Airport — Credit card shops → Every 5th person they pitch for credit card (5th, 10th, 15th, ...)

Feedback Survey → Selected every 11th member
↓
Dropped mail for feedback.

c) Stratified sampling

Divide the population into strata (groups) based on specific characteristics and then randomly sample from each strata

Eg: ① Divide the employees by department and then randomly select a proportional number from each department to form a survey sample

② Population divided on age group

< 12 12 - 18 > 18 { survey on food

d) Cluster Sampling

Divide the population into clusters, randomly selecting clusters then sampling all the members

from the selected clusters.

Eg: Randomly selected several schools from a district and surveying all teachers within those schools.

e) Multi Stage Sampling

combining several sampling methods,

Selecting clusters → randomly sampling within those clusters

Eg:- Randomly selecting cities, each selected city randomly selecting households to survey.

2. Non probability Sampling

Select individuals who are easiest to reach.

Eg: surveying people at mall

a) Convenience Sampling

Select individuals who are easiest to reach

b) Judgemental (Purposive) Sampling

Select individual based on the researcher's judgement → useful or Representative

Eg: choose experts in a field to participate
{ Data Science }

(c) Snowball sampling

Existing study subjects recruit future subjects from among their acquaintances

Eg: Survey members of a rare disease

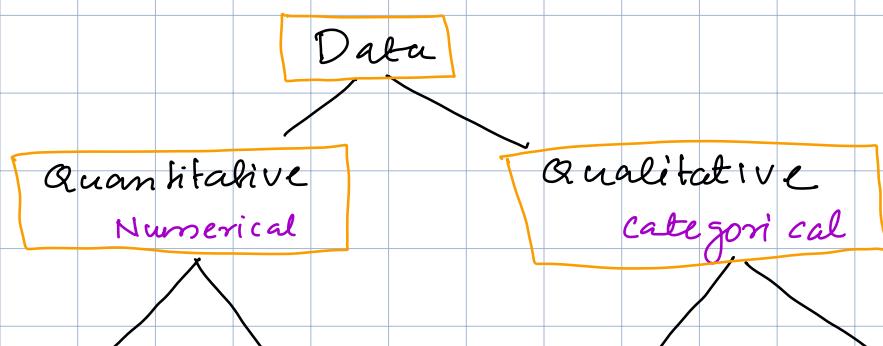
d) Quota Sampling

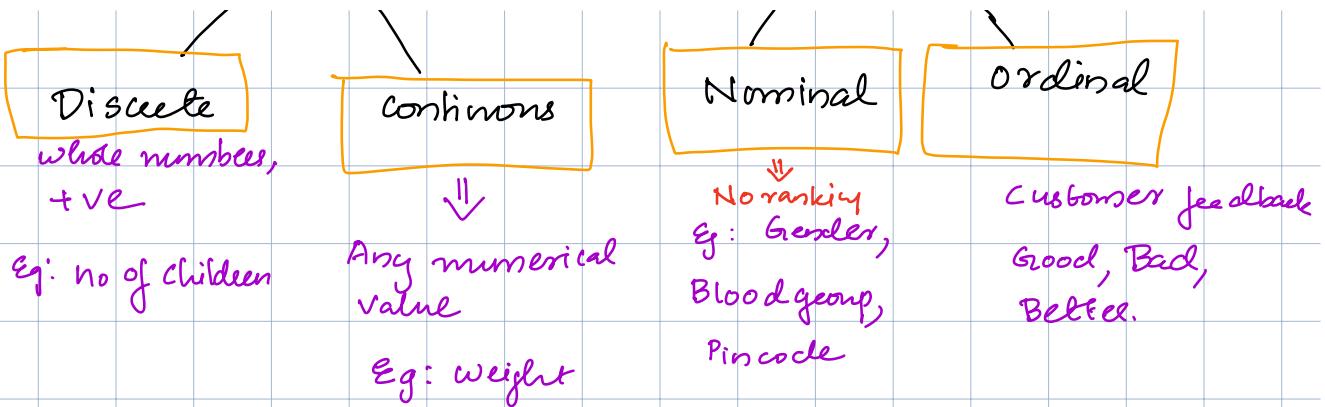
Age, group, gender, caste.

Selecting the sampling technique depends on use cases.

(5) Types of Data

Data is of different type → important while analysis and model selection





⑥ Scales of measurement of Data

The nature of information within the value assigned to variables

4 Primary scales of Measurements

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

D) **Nominal Scale** - classify data into distinct categories that do not have an intrinsic order.

Qualitative / Categorical data
characteristics

- 1) data is categorized based on labels, names or qualities
- 2) These categories are mutually exclusive
- 3) No logical order among categories [No rank]

Eg: Gender M, F

Color R, B, P, G

Types of Cuisine Italian, Indian, Mexican

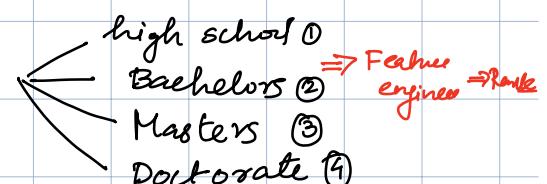
(2) Ordinal Scale

classifies the data into categories that can be ranked or ordered

Characteristics

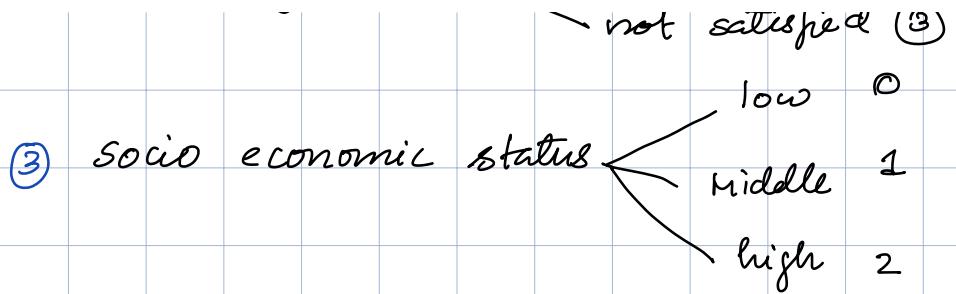
- i) data is categorized and ranked in a specific order
- ii) Interval between Ranks are not necessarily equal.

Example :- ① Education level



satisfied ②

② Customer feedback very satisfied ①



(3) Interval scale

The interval scale not only categorizes and orders but also specify the exact difference between intervals. It lacks a true zero point.

characteristics

- 1) Data is ordered with consistent interval between values
- 2) Allows for meaningful comparison of differences [ratio cannot be measured]
- 3) No true zero point.

Example : Temperature in Fahrenheit

10°F , 20°F , 30°F

$$20 - 10 = 10$$

⇒ differences same

$$30 - 20 = 10$$

Temp 0°F not possible

IQ scores

90 100 110
 \ \ /
 10 10

10 > 0 is not possible

Calender years

2024, 2020, 2016,

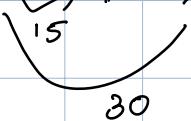
0 year is not defined

4) Ratio Scale

- ① The order matters
- ② Differences are measurable
- ③ contains a zero starting point

Eg: student marks in a class

0, 90, 60, 30, 75, 45

ASC = 0, 30, 45, 60, 75, 90


$$\text{Ratio} = 90/30 = 3:1$$

Assignment

- ① lengths of different rivers in world \Rightarrow Interval
- ② favorite food based on gender \Rightarrow Ordinal

③ Marital status = Nominal

④ IQ measurement = Ratio

Measure of Central Tendency

Measures of central tendency are statistical metrics that describe the center point or typical value of dataset. They provide a single value that summarizes a set of data by identifying the central position within that dataset.

① Mean or Average

② Median

③ Mode

Eg : Ages = [24, 32, 12, 48, 16, 20]

① Mean

Mean is sum of all values divided by the number of values.

Population mean (μ)

population (N)

$$\mu = \frac{\sum_{i=1}^N x_i}{N} = \text{population size}$$

Here X is a random variable

$$X = \{5, 8, 12, 15, 20\}$$

$$N = 5$$

$$\mu = \frac{60}{5} = 12$$

Characteristics

= Mean is usually affected by extreme outliers

= Usually used for interval and ratio data.

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

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

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

$$\mu = \frac{115}{6} = 19.1$$

② Median

Middle value in a dataset when the values are arranged in either ascending or descending order.

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

No of elements = 5

5 is odd

Middle element = 3 = Median

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

Median = 3.5

No much change

Characteristics

- = Not affected by extreme outliers
- Used specifically for ordinal, interval and ratio data

③ Mode

Appears most frequently in a dataset

Dataset - 2, 4, 4, 6, 7, 7, 7, 9

frequency

2	1
4	2
6	1
7	3

some datasets have multiple modes

2 modes - bimodal

more modes - multimodal

Characteristics

- ① Not affected by outliers
- ② Used for all the four scales nominal, ordinal, interval, ratio

Choosing the Appropriate Measure

Mean = data when symmetrical and without outliers, provide some mathematical average which is useful for further statistical calculations.

Median :- Best used when data is skewed or contains outliers. Provides the middle value, which better represents the center of a skewed dataset

Mode :- Best used for categorical data to identify the most common category. Also

useful for identifying the most frequent value in ordinal, interval or ratio data.

Real World Applications

Feature Engineering

	Age	Weight	Salary	Gender	Degree
24	70	40k	M	BE	
25	80	70k	F	-	
27	95	45k	F	-	
29	-	50k	M	PhD	
32	-	60k	-	BG	
-	60	-	-	Master	
-	65	55k	-	BSC	
40	F2	-	M	BE	

① Handling Missing Values



if Mean \Rightarrow outliers create problem

Median \Rightarrow more efficient

Interval scale, Ratio Data

Categorical \Rightarrow Mode

Nominal

Mode

Nominal & Ordinal

Measure of Dispersion

Measure of dispersion describe the spread or variability of dataset. They indicate how much the values in a dataset differ from the central tendency.

Common Measure of Dispersion

① Range

② Variance

③ Standard Deviation

④ Interquartile Range (IQR)

① Range

Defⁿ :- Difference between the maximum value and minimum values in a dataset.

$$\text{Range} = \text{Max Value} - \text{Min Value}$$

e.g.: Ages { 14, 13, 10, 20, 25, 45, 15 }

$$\text{Range} = 45 - 10 = 35$$

Characteristics

- 1) simple to calculate
- 2) sensitive to outliers
- 3) Provides only rough measure of dispersion.

$$\text{Weight} = \{ 35, 40, 45, 39, 30, \dots, 40 \}$$

$$\begin{aligned}\text{Range without FO} &= 40 - 30 = 10 \\ \text{Range with FO} &= 40 - 30 = 10\end{aligned}$$

(2) Variance

Defn :- Variance measures the average squared deviation of each value from the mean. It provides sense of how much the values in a dataset vary

Population Variance

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

Sample Variance

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

x_i \Rightarrow DATA POINTS

μ = Population mean

N = population size

x_i = Data points

\bar{x} = Sample mean

n = Sample size

Example \rightarrow size of a flower petals

$\{5, 8, 12, 15, 20\} \Rightarrow$ Variance of this distribution

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

$$\mu = \frac{5+8+12+15+20}{5} = 12$$

$$\text{Variance} = \frac{7^2 + 6^2 + 0^2 + 3^2 + 8^2}{5} = 49 + 36 + 9 + 64 \\ = \underline{\underline{27.6}} = \sigma^2$$

Characteristics

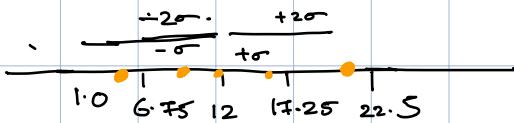
- ⇒ Provides a precise measure of variability
- ⇒ Units are squared of the original data unit
- ⇒ More sensitive to outliers

(3) Standard Deviation

Defn:- The standard deviation is the square root of the variance

$$\sigma = \sqrt{27.6} \approx 5.25$$

$\{5, 8, 12, 15, 20\}$



z -score

\rightarrow standard normal distribution

Characteristics

- ① provides a clear measure of spread in same units as the data
- ② Also sensitive to outliers

Key Differences and Similarities

Relationship

Standard deviation is the square root of variance. If you have the variance, you can find the standard deviation by taking the square root of the variance.

Conversely if you have the standard deviation, you can find the variance by squaring the standard deviation.

units:

Variance :- The unit of variance = square of unit of

original data

- provides measure of dispersion of data points in

squared units, which is difficult to interpret directly

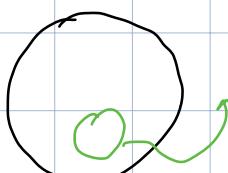
Standard deviation = same unit as original data.

= provides a measure of dispersion in the same units as the original data, making it easier to interpret and understand.

Sample Variance

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

$n-1$ = Bassel correction



Inferences

Population Variance

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

$$\begin{cases} \bar{x} \approx \mu \\ s^2 \approx \sigma^2 \end{cases} \quad \begin{array}{l} \text{Might come up} \\ \text{case 1} \end{array}$$

$$\begin{cases} \bar{x} \ll \mu \\ s^2 \ll \sigma^2 \end{cases} \quad \begin{array}{l} \text{case 2} \end{array}$$

if

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

then in case of case 2, we

are understanding the true population variance

$$\left(s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \right) \quad \begin{array}{l} \text{Bessel correction} \\ \text{degree of freedom } n-1 \end{array}$$

Random Variables X

$$y = 5x + 2$$

$$y = 7$$

$$x = 1$$

$$y = 12$$

$$x = 2$$

$$y = 17$$

$$x = 3$$

X \rightarrow function whose values are derived from different process or experiment.

e.g.: Tossing a coin.

$$X = \begin{cases} 0 & H \\ 1 & T \end{cases}$$

Rolling a fair dice

$$X = \begin{cases} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{cases}$$

Different types of Random Variables

Random Variables

Discrete R.V

Eg: Tossing coin
rolling dice

Continuous R.V

Eg: Tomorrow how many inches it
is going to rain
[0, 1.1, 5.5, ... 10.75]

X - function $\xrightarrow[\text{process}]{}$ value

Percentiles And Quartiles

Percentage :- $\{1, 2, 3, 4, 5, 6\}$

No of odd numbers = 3

percentage of odd numbers in this group, $\frac{3 \times 100}{6} = 50\%$.

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

$\{2, 2, 3, 4, 5, 5, 6, 7, 8, 8, 8, 9, 9, 9, 10\}$

Percentile of Value $x = \frac{\# \text{ of values below } x \times 100}{n}$

if $x = 9$

$$= \frac{11}{14} \times 100$$

= 78.57% of value 9



78.57% of entire distribution is less than 9

$$\Rightarrow \text{Value} = \frac{\text{Percentile}}{100} \times (n+1)$$

if percentile = 25

$$\text{Value} = \frac{25}{100} \times (15)$$

$$= \frac{1}{4} \times 15 = 3.75$$

3.75 not in the data so 3rd quartile number

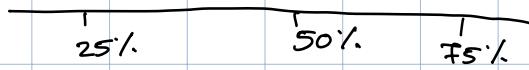
$$\text{average} = \frac{3+4}{2} = \underline{\underline{3.5}} = 25 \text{ percentile}$$

② Quartiles

25% = 1st quartile

50% = 2nd quartile

75% = 3rd quartile



③ 5 Number summary and Boxplot [To find outliers]
- in python using seaborn.

(i) Minimum

2) First Quartile (25 percentile) Q_1

3) Median

4) Third Quartile (75 percentile) Q_3

5) Maximum

Removing the outliers

$$X = \{ \underline{1}, \underline{2}, \underline{2}, \underline{2}, \underline{3}, \underline{3}, \underline{4}, \underline{5}, \underline{5}, \underline{6}, \underline{6}, \underline{6}, \underline{7}, \underline{8}, \underline{9}, \underline{\underline{29}} \}$$

\downarrow
 5^{th}

\downarrow
outlier

[Lower Fence \longleftrightarrow Higher Fence]

$$\text{Lower Fence} = Q_1 - 1.5 (\text{IQR})$$

, Inter Quartile Range

$$\text{Higher Fence} = Q_3 + 1.5 (\text{IQR})$$

$$Q_1 = 25 \text{ percentile} = \frac{25}{100} \times (19+1) = \frac{25}{100} \times (n+1)$$

= 5th value

$$Q_1 = 3$$

$$Q_3 = \frac{75}{100} \times (20) = \frac{3}{4} \times 20 = 15^{\text{th}} \text{ value}$$

= 7

$$IQR = Q_3 - Q_1 = 7 - 3 = 4$$

$$\text{Lower Fence} = Q_1 - 1.5(IQR)$$

$$= 3 - 1.5(4) = 3 - 6 = -3$$

$$\text{Higher Fence} = Q_3 + 1.5(IQR)$$

$$= 7 + 6 = 13$$

$$[-3 \quad 13]$$

||

Values outside this range is an outlier

so here outlier is 29

How do we create a box plot

Minimum = 1

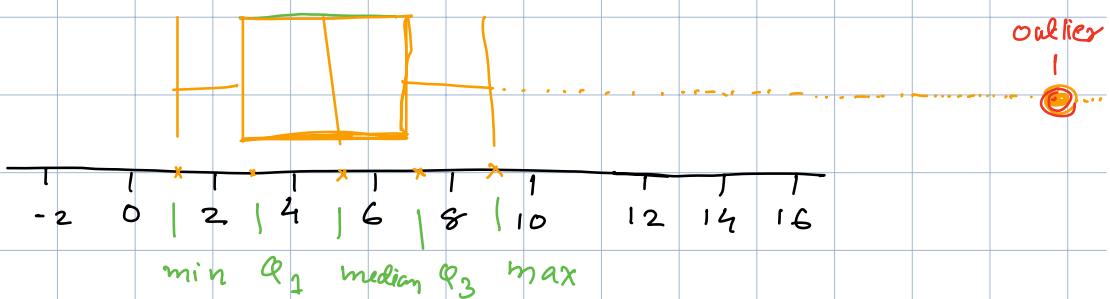
$Q_1 = 3$

$Q_3 = 7$

Median = (19 elements so 10th element)
= 5

Maximum = 9 (excluding the outlier)

Box plot



⇒ Outlier play a major role in datascience project-

Assignment

$$Y = \{ -13, -12, -5, -6, 3, 4, 5, 6, 7, 8, 10, 10, \\ 11, 55 \}$$

$$n = 15$$

$$Q_1 = \frac{25}{100} \times 16 = \underline{\underline{4}}$$

$$Q_3 = \frac{75}{100} \times 16 = 12$$

$$IQR = 12 - 4 = 8$$

$$\text{Lower Fence} = Q_1 - 1.5(8) = 4 - 1.5(8) \\ = 4 - 12 = -8$$

$$\text{Higher Fence} = 12 + 12 = 24$$

$$[-8, 24]$$

$$\text{Outliers} = [-13, -12, 55]$$

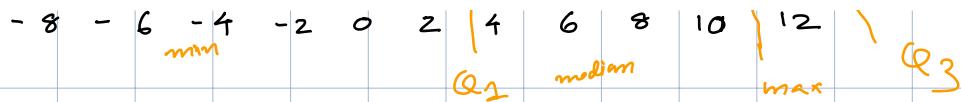
Box plot

$$\min = -5$$

$$\max = 11$$

$$\text{median} = 6$$





Histograms And Skewness

A histogram is a graphical representation of the distribution of numerical data. It is an estimate of the probability distribution of a continuous variable and is used to visualize the shape, central tendency and variability of a dataset.

$$\text{Age} = \{ 11, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50 \}$$

Range here = (11, 50)
consider all the histograms between 0-50

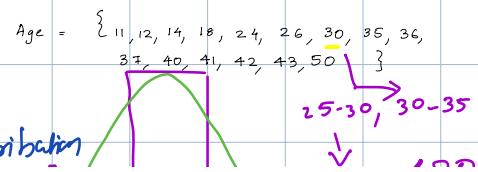
① No of bins = 10 $\frac{50}{10} \rightarrow 5$ = 5 bin size

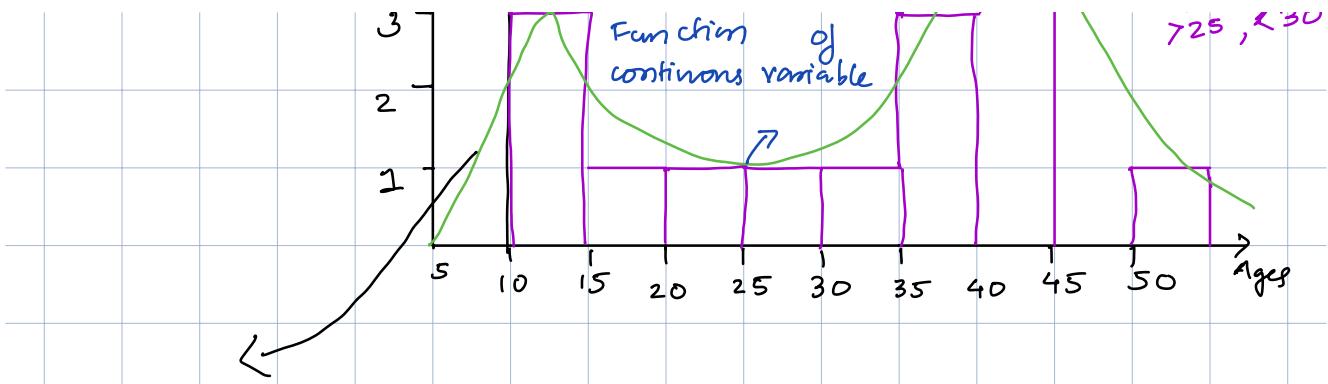
i.e. 10 bins each have size will be 5.

Bins $\rightarrow [0-5, 5-10, 10-15, 15-20, \dots, 45-50]$

frequency = count
5
4

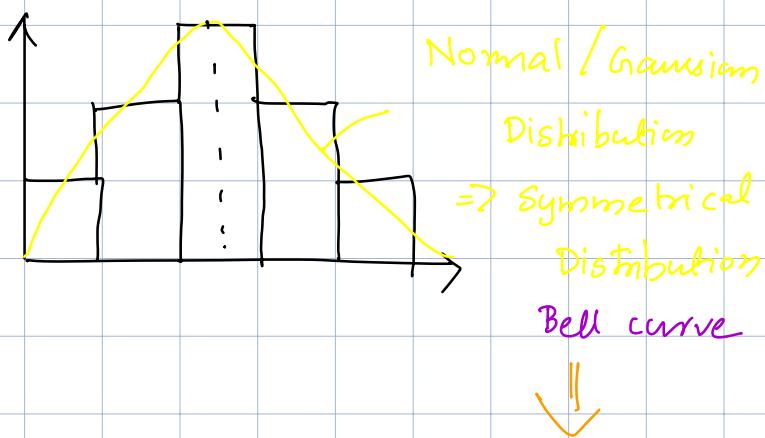
Probability distribution





Skewness

If we create a histogram



No skewness in symmetrical distribution.

①



The mean, median and mode are all perfectly at the center in the symmetrical distribution

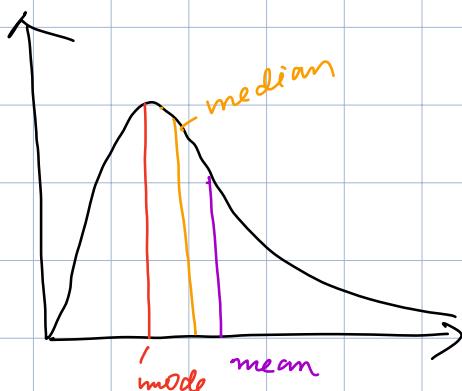
$$\text{Mean} = \text{Median} = \text{Mode}$$

(2)

Right Skewed



Right skewed / log Normal distribution.

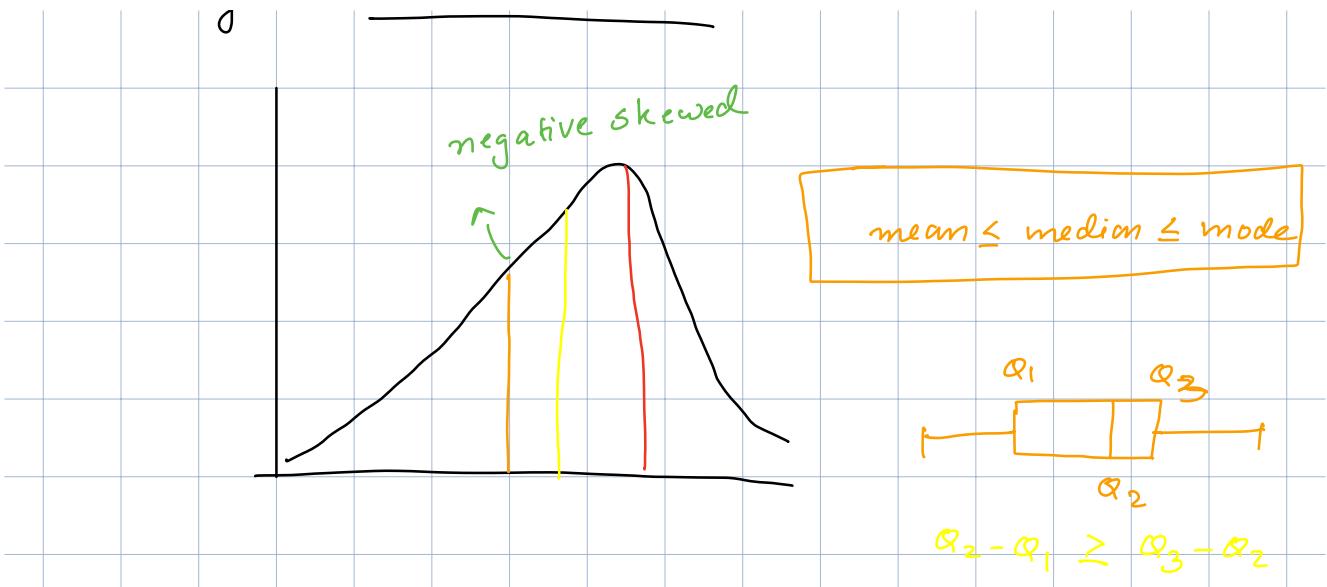


$$\boxed{\text{mean} \geq \text{median} \geq \text{mode}}$$



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

(3) Left Skewed Distribution



Covariance and Correlation

Covariance and correlation are two statistical measures used to determine the relationship between two variables. Both are used to understand how changes in one variable are associated with changes in another variable.

Covariance

Covariance is a measure of how much two random variables change together. If the variables tend to increase and decrease together, the covariance is positive. If one tends to increase when the other decreases, the covariance is negative.

to quantify the relationship between
x and y

Eg :-	x	y	
	2	3	$x \uparrow y \uparrow$
	4	5	$x \downarrow y \uparrow$
	6	7	$x \uparrow y \downarrow$
	8	9	$x \downarrow y \downarrow$

Data Set

size of house

1200

1300

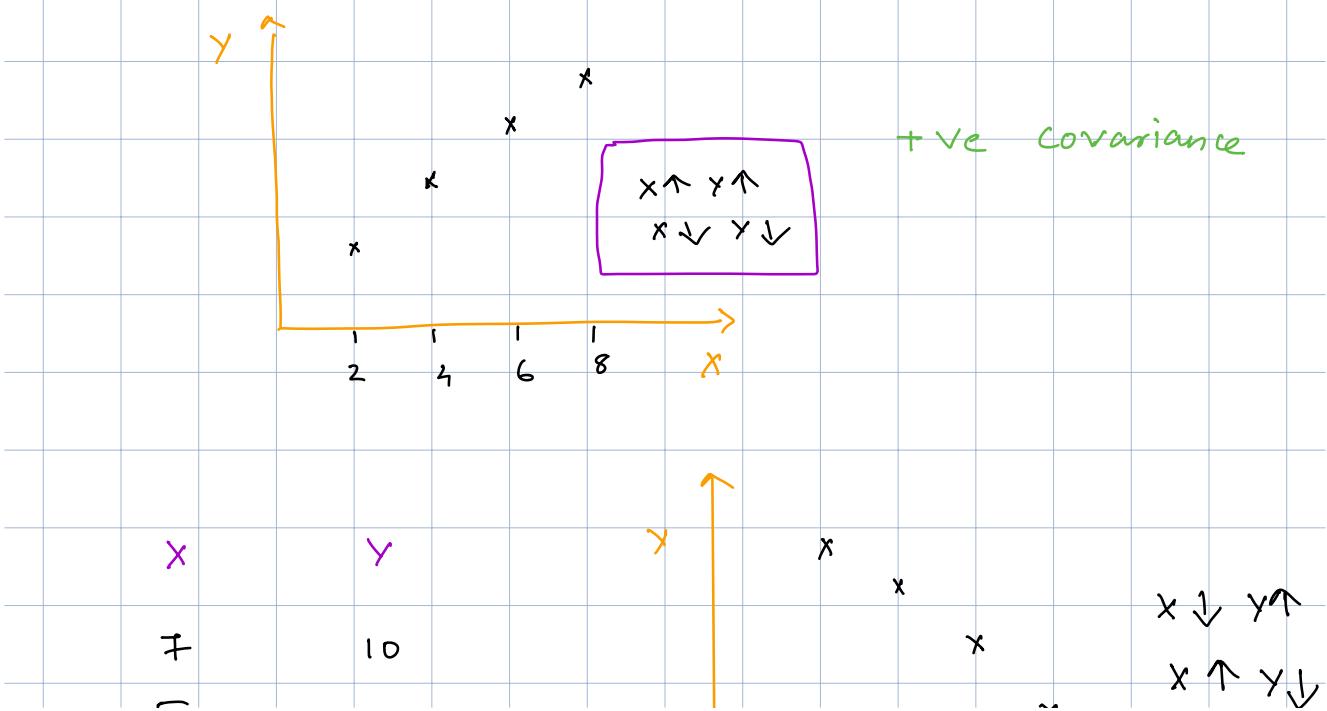
1500

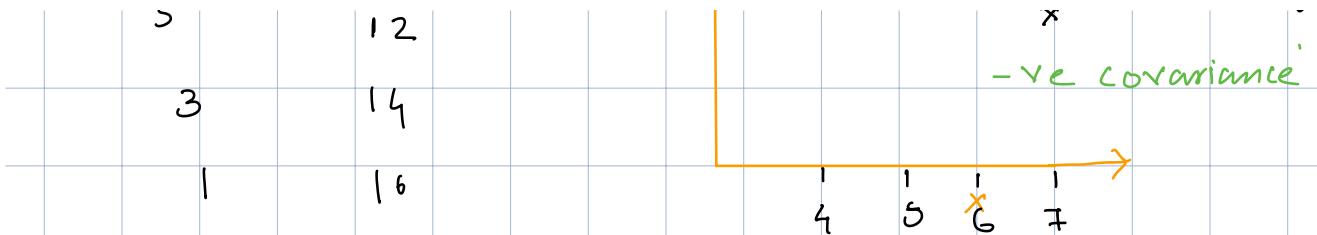
Price

45000 Euro

50000 Euro

75000 Euro





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

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

\Downarrow
Spread of data

x_i = Datapoint of random variable X

\bar{x} = Sample mean of n

y_i = Datapoint of random variable Y

\bar{y} = Sample mean of Y

Eq:

Students

Hours Studied

2

3

4

Exam Score (Y)

50

60

70

$x \uparrow x \uparrow$

			$X \downarrow$	$Y \downarrow$	
	5				80
	6				90

Calculate the Covariance

$$\bar{x} = \frac{2 + 3 + 4 + 5 + 6}{5} = \frac{20}{5} = 4$$

$$\bar{y} = \frac{50 + 60 + 70 + 80 + 90}{5} = 70$$

$$\text{Cov}(x, y) = \frac{-2(-20) + -1(-10) + 0(0) + 1(10) + 2(20)}{4}$$

$$= \frac{40 + 10 + 10 + 40}{4} = \frac{100}{4} = 25$$

Positive Covariance

Indicates the number of hours studied increases the exam score.

Advantages

⇒ Quantify the relationship between X and Y

Disadvantage

① Covariance does not have a specific limit value

$(-\infty, \infty)$
limit value is not restricted

Correlation

→ Pearson Correlation Coefficient

→ Spearman Rank Correlation.

Pearson Correlation Coefficient $\Rightarrow [-1 \text{ to } 1]$

$$\rho_{x,y} = \frac{\text{Corr}(x, y)}{\sigma_x \cdot \sigma_y}$$

① The more the value towards +1 the more positive correlated x and y is.

② The more the value towards -1 the more -ve correlated it is.

For the earlier problem

$$\rho_{x,y} = \frac{\text{Corr}(x, y)}{\sigma_x \cdot \sigma_y} = \frac{25}{\sigma_x \cdot \sigma_y}$$

$$\sigma_x = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}} = \sqrt{\frac{1^2 + 2^2 + 1^2 + 2^2}{4}} = 5$$

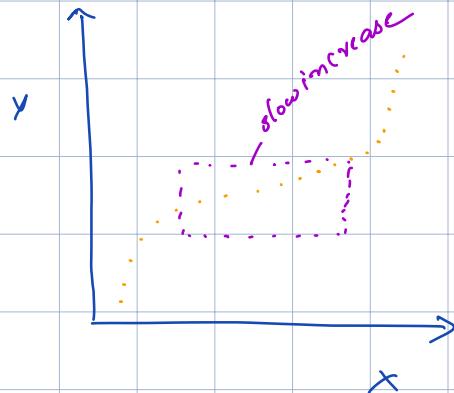
$$\sigma_y = \sqrt{\frac{20^2 + 10^2 + 20^2 + 10^2}{4}} = \sqrt{\frac{1000}{4}} = 250$$

$$\rho_{x,y} = \frac{25}{5(250)} = \frac{1}{10(5)} = \frac{1}{50} = 0.02$$

U to 1
↓
+ve

Spearman Rank Correlation

In pearson correlation = $-1/1$ in a linear
in between = scattered)



Spearman Correlation 1
Pearson Correlation 0.88

Pearson correlation cannot capture relationship
in nonlinear data.

Spearman Correlation $\gamma_s = \frac{\text{Cor}(R(x), R(y))}{\sigma(R(x)) * \sigma(R(y))}$

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

Eq

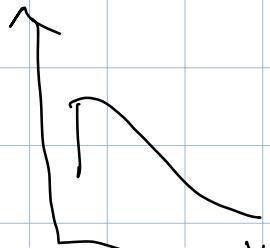
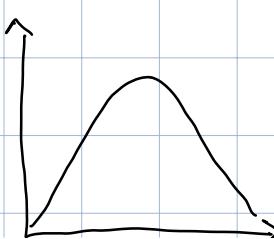
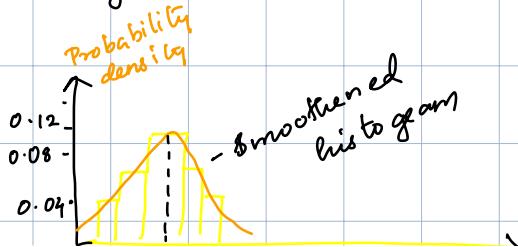


Probability Distribution Function and Types of Distribution.

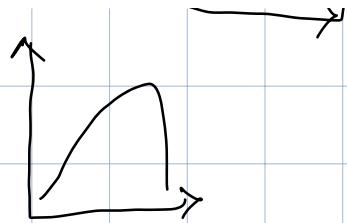
Probability distribution function

Probability distribution functions describe how the probabilities are distributed over the values of a random variable.

$$\text{Age} = \{ \dots \}$$



PDF helps us to understand how the probabilities are distributed over random variables



2 Main type of probability distribution functions

① Probability Mass Functions (PMF): used for discrete random variables

② Probability density function (PDF): used for continuous random variables.

③ Cumulative density Function (CDF)

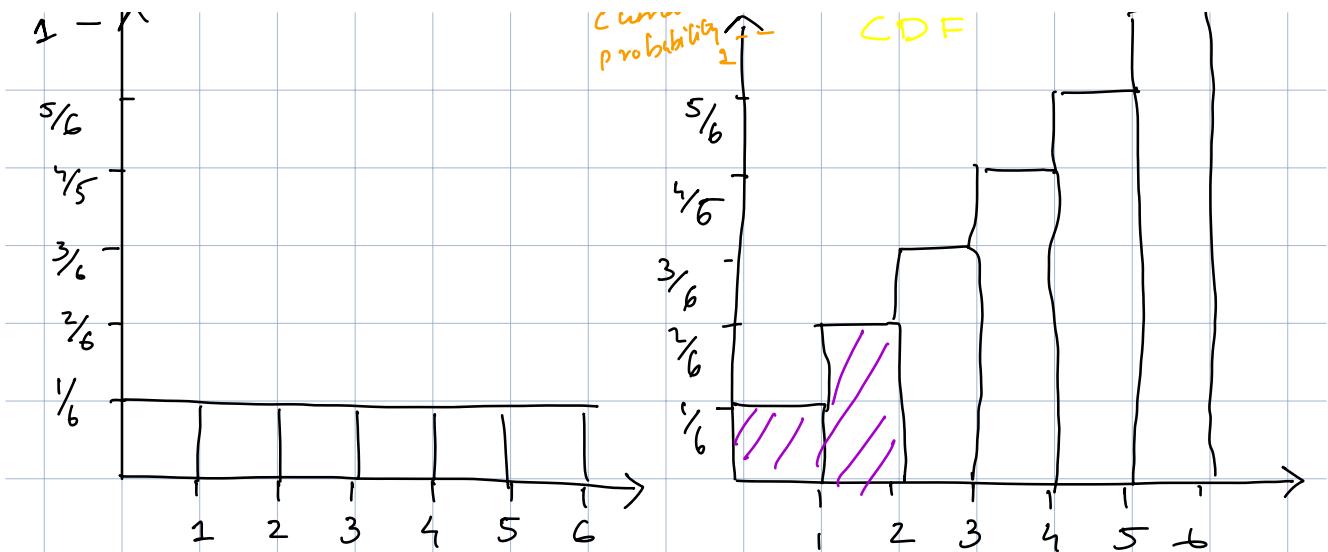
① Probability Mass Function [specifically for discrete random variable]

Eg: Rolling a dice $\{1, 2, 3, 4, 5, 6\}$ =
Fair Dice

$$P(1) = P(2) = P(3) = P(4) = P(5) = P(6) = \frac{1}{6}$$

Probability
PMF

Fair dice
by
Cumulative Density Function



combine probability
as we go from 1 to 6

Why CDF?

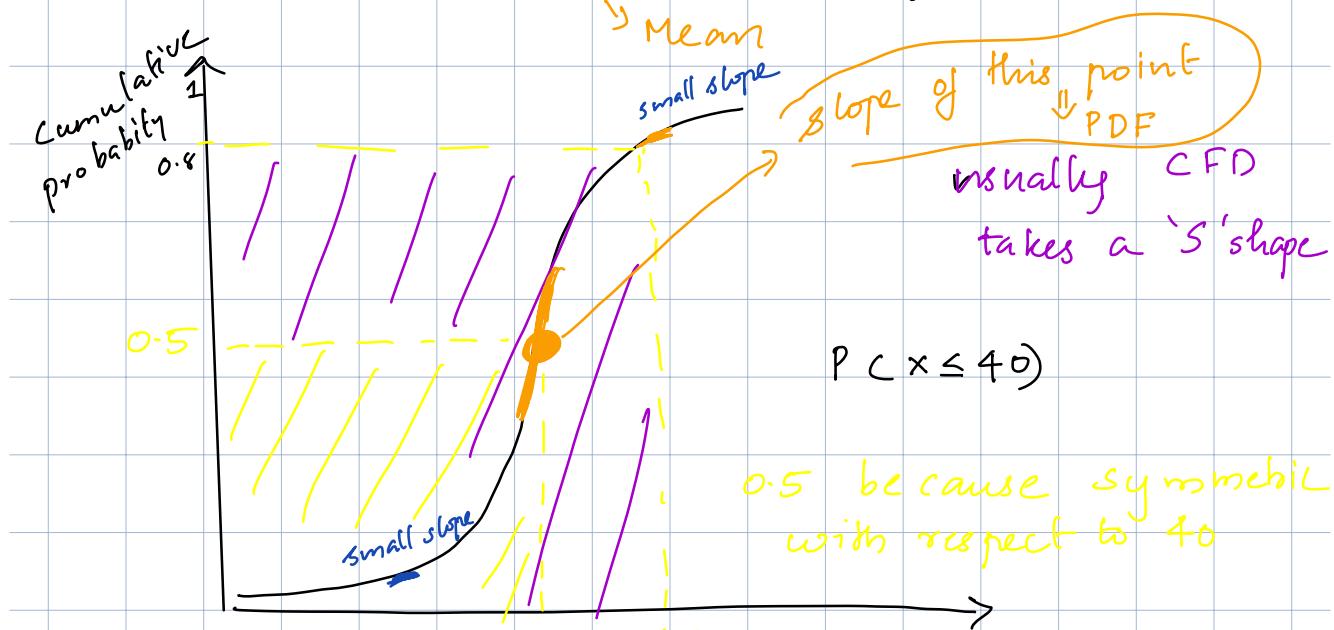
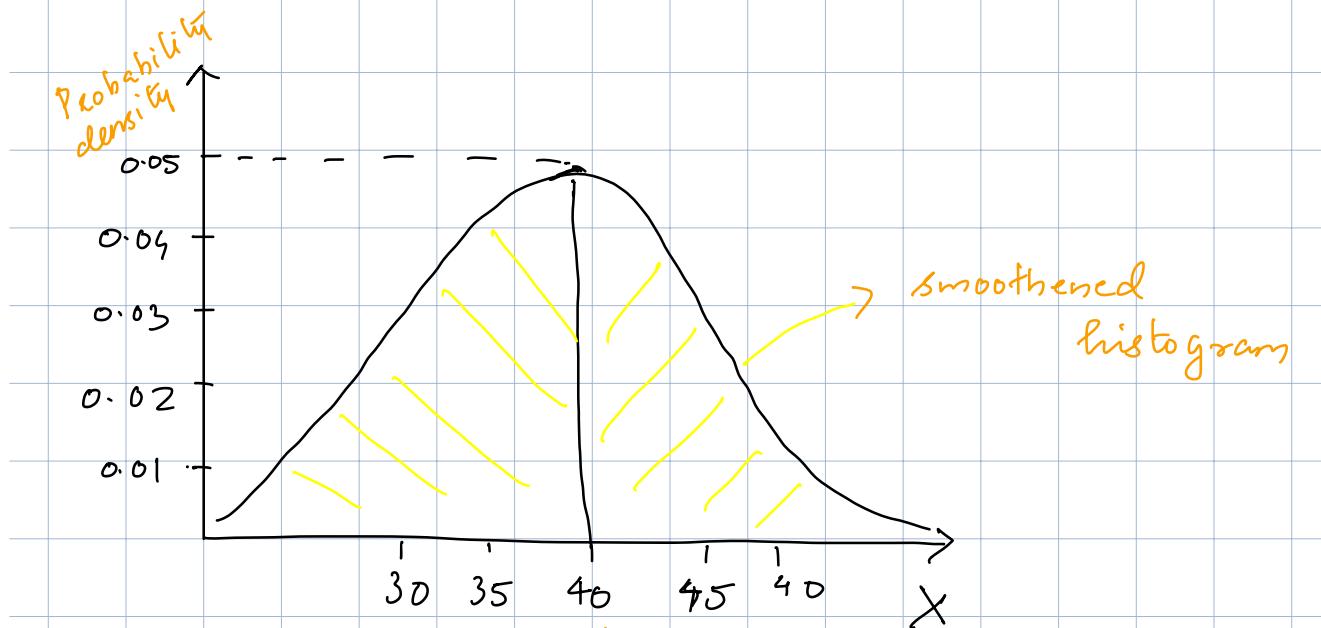
$$\begin{aligned}
 P(X \leq 2) &= P(x=1) + P(x=2) \\
 &= P(X \leq 2) = \frac{2}{6} = \frac{1}{3}
 \end{aligned}$$

$$\begin{aligned}
 P(X \leq 6) &= P(x=1) + P(x=2) + P(x=3) \\
 &\quad + P(x=4) + P(x=5) + P(x=6) \\
 &= 1
 \end{aligned}$$

② Probability Density Function (PDF)

① Distribution of continuous random variable

$$X = \text{Ages} = \{ \dots \}$$



40 45

x

① Area under the curve

② Probability density

So slope of CDF give PDF

$$\Pr(X \leq 45) = 0.8 = 80\%.$$

So Probability density is gradient of
cumulative density function.

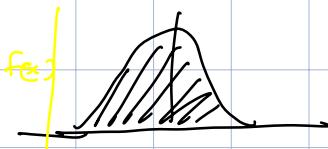
PDF properties \Rightarrow

① Always non negative $f(x) \geq 0$ for all x

② The total area under the PDF curve

is equal to 1

$$\int_{-\infty}^{\infty} f(x) dx = 1$$



With respect to different distribution function is going to change.

Different distribution types [pdf, pmf, cdf]

Dataset often follows different types of distribution:

(1) Bernoulli Distribution \rightarrow Outcomes are binary

(pmf) = Discrete Random Variable

(2) Binomial Distribution \rightarrow pmf

(3) Normal / Gaussian Distribution \rightarrow pdf

= some sort of bell curve

= used in assumptions

(4) Poisson Distribution = pmf

(5) Log Normal Distribution = (pdf)

(6) Uniform Distribution = (pmf)

Eg: House Price Predictions Dataset

size of house	No of Rooms	location	Floor	Sea side	Price
continuous	Discrete		Discrete	Bernoulli 0 or 1	Continuous pdf

This step is important in explanatory data analysis and feature engineering.

① Bernoulli Distribution

Definition:- The Bernoulli Distribution is the simplest discrete probability distribution. It represents the probability distribution of a random variable that has exactly two

possible outcomes: success with probability p) and failure (with probability $1-p$). It is used to model binary outcomes, such as a coin flip or a yes/no question.

① Discrete Random Variable (contd)

② Outcomes are Binary

Eg: ① Tossing a coin $\{H, T\}$

$$\Pr(X = H) = 0.5 = p$$

$$\Pr(X = T) = 1 - 0.5 = 0.5 = q$$

$$p, q \Rightarrow p = 1 - q \quad q = 1 - p$$

② Whether the person will pass/Fail

$$\Pr(X = \text{pass}) = 0.4$$

$$\Pr(X = \text{Fail}) = 1 - 0.4 = 0.6$$

Parameters

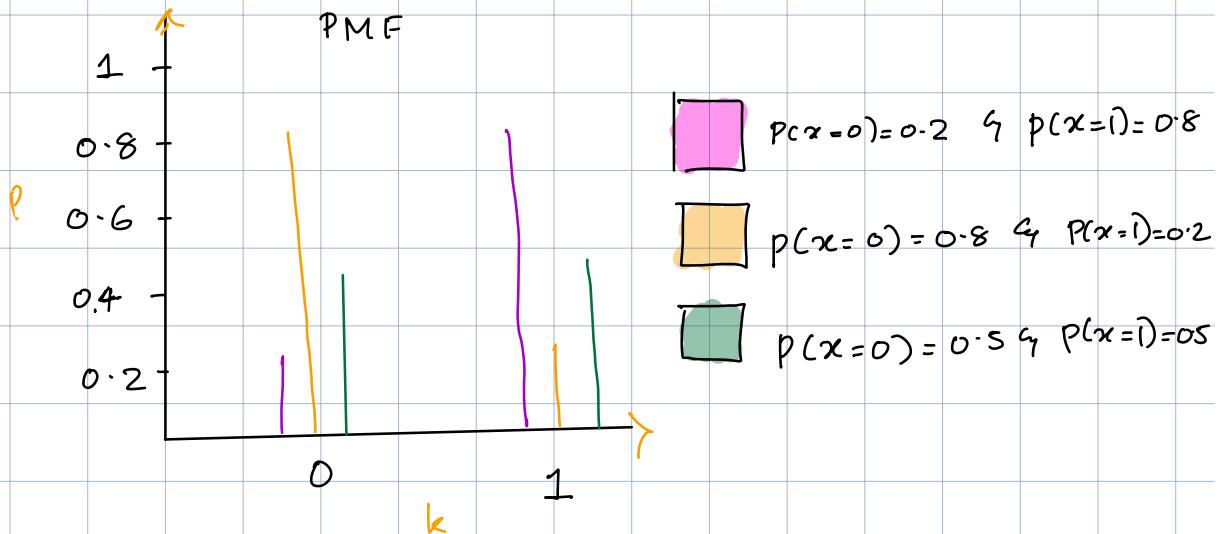
$$0 \leq p \leq 1$$

$$q = 1 - p$$

$$K = \{0, 1\} \Rightarrow 2 \text{ outcomes}$$

$\Pr(\text{Success}) \Rightarrow k=1$

$\Pr(\text{Fail}) \Rightarrow k=0$



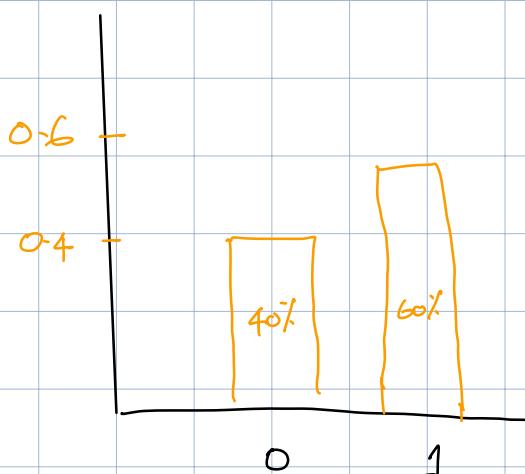
e.g/

PMF :- Company has launched a new smartphone 'A'

k

(1) use = 60%

(2) not use = 40%



$$\text{PMF} = P^k * (1-p)^{1-k}$$

$$\text{if } k=1 \quad P = 0.6$$

$$\text{ie } p(k=1) = p^1(1-p) = p$$

$$p(k=0) = p^0(1-p)^{1-0} = 1-p = q$$

Simplified

$$P_m f \left\{ \begin{array}{ll} q = 1-p & \text{if } k=0 \\ p & \text{if } k=1 \end{array} \right.$$

Mean of Bernoulli Distribution

$$E(x) = \sum_{i=0}^{k=1} k \cdot p(k) \quad k \in \{0, 1\}$$

$$= 0 \cdot (0.4) + 1 \cdot (0.6)$$

$$= 0 + 0.6 = 0.6 \Rightarrow P$$

Median of Bernoulli Distribution

Median

{	0	if $p < \frac{1}{2}$
	$[0, 1]$	\downarrow
	1	if $p > \frac{1}{2}$

Median

{	0	if $q > p$
	0.5	if $q = p$
	1	if $q < p$

* Mode

$p > q \Rightarrow p$ will be the mode
else q will be the mode.

* Variance

$k=0$ and 1

$$P(k=0) = 0.4$$

$$P(k=1) = 0.6$$

$$\sigma^2 = 0.4 * (0 - 0.6)^2 + 0.6 * (1 - 0.6)^2$$

$$= \sum_{i=1}^n P_i (x_i - \mu)^2 \quad \text{here } \mu = 0.6$$

$$= 0.4(0.36) + 0.6(0.16)$$

$$\sigma^2 = 0.24 \Rightarrow P(k=0) * P(k=1)$$

$$\sigma^2 = P * V$$

$$\sigma = \sqrt{P * V}$$

② Binomial Distribution

In probability theory and statistics the binomial distribution with parameters n and p is the discrete probability distribution of the numbers 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 (with probability $q = 1-p$). A single success/failure experiment is also called a Bernoulli trial or 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.

① Discrete Random Variable

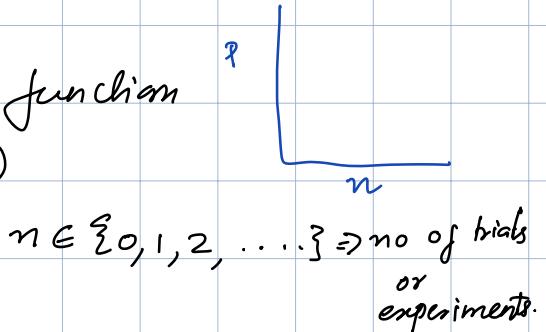
- (*) Every outcome of the experiment is binary
- (*) These experiments are performed for n trials

Eg: Tossing a coin 10 times, $n=10$

\Downarrow
 $\{H, T\}$

PMF function of Binomial function

Notation: $B(n, p)$



$p \in [0, 1] \Rightarrow$ success probability for each trial.

$$q = 1-p$$

Support :- $k \in \{0, 1, 2, 3, \dots, n\} =$ Number of success.

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

$${}^n C_k = \frac{n!}{k!(n-k)!} \Rightarrow \text{Binomial Coefficient}$$

$$\text{mean} = n \cdot p$$

$$\text{variance} = n \cdot p \cdot q$$

$$\sigma = \sqrt{n p q}$$

Eg: Coin Flip

Number of trials (n) = 5

Probability of success (p) = 0.5

No of success (k) = varies from 0 to 5

i) What is the probability of getting exactly 3 heads in 5 flips?

$$n = 5$$

$$k = 3$$

$$\Pr(X=3) = {}^5C_3 (0.5)^3 (1-0.5)^{5-3}$$
$$= 0.3125$$

Example: Quality control

Scenario: Inspecting 10 items in a factory where each item has a 10% chance of being defective.

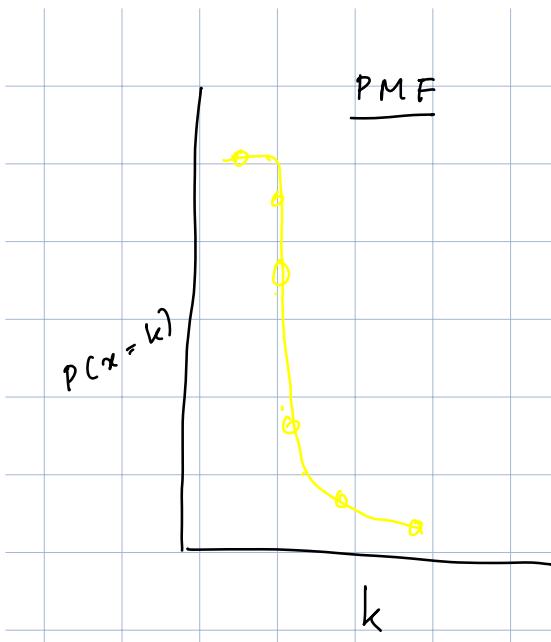
- * No of Trials (n) = 10
- * Probability of success (p) = 0.1 (defective item)
- * No of successes (k) = varies from 0 to 10

question : What is the probability of finding exactly 2 defective items in a sample of 10?

$$\begin{aligned}
 P(X=2) &= {}^{10}C_2 (0.1)^2 (0.9)^{10-2} \\
 &= \frac{10!}{8! 2!} (0.1)^2 (0.9)^8 = 0.1937 //
 \end{aligned}$$

Poisson Distribution

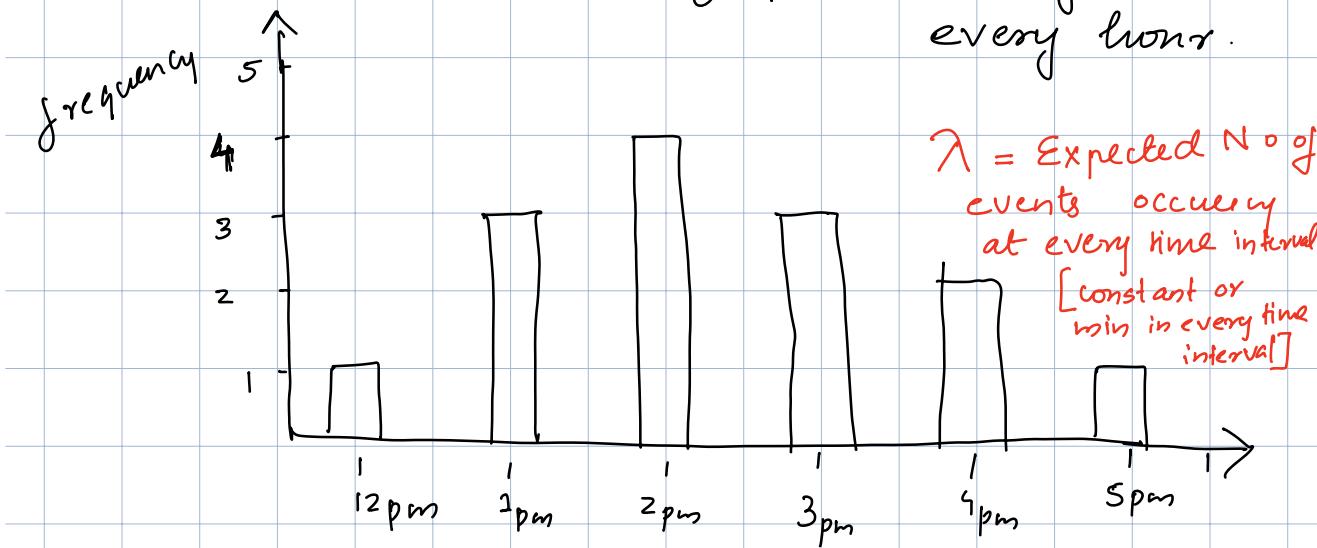
In probability theory and statistics, the Poisson distribution is a discrete probability distribution that expresses the probability of a given number of events occurring in a fixed interval of time if these events occur with a known constant mean rate and independently of the time since the last event.



(1) Discrete random variable (Pmf)

(2) Describe the numbers of events occurring in a fixed time intervals

Eg: No of people visiting hospital every hour
No of people visiting banks every hour.



λ = Expected No of events occurring at every time interval
[constant or min in every time interval]

PMF

$$P(X=5) = \frac{e^{-\lambda} \lambda^x}{x!}$$

$$\lambda = 3$$

$$= \frac{e^{-3} 3^5}{5!} = 0.101 = 10.1\%$$

Mean of Poisson Distribution

$$\text{Mean} = E(X) = \mu = \lambda * t$$

λ = expected No of events occur in every time interval

t = time interval

Normal / Gaussian Distribution

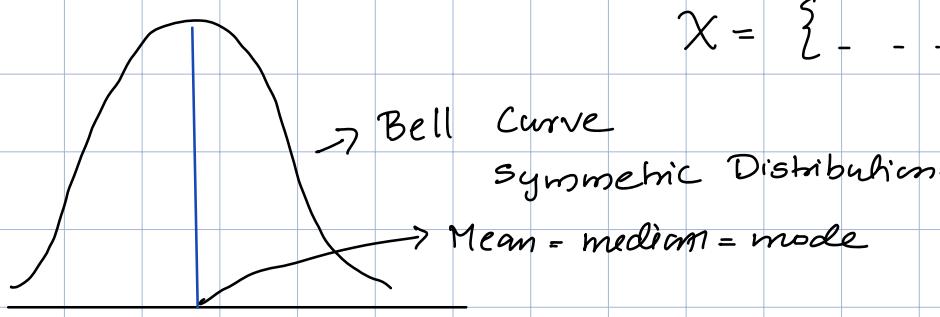
In probability theory and statistics, a normal distribution or Gaussian distribution is a type of continuous probability distribution for a real-valued random variable

⇒ most data follows Normal distribution

⇒ For continuous R.V (P.d.f)

⇒

$$X = \{ \dots \}$$



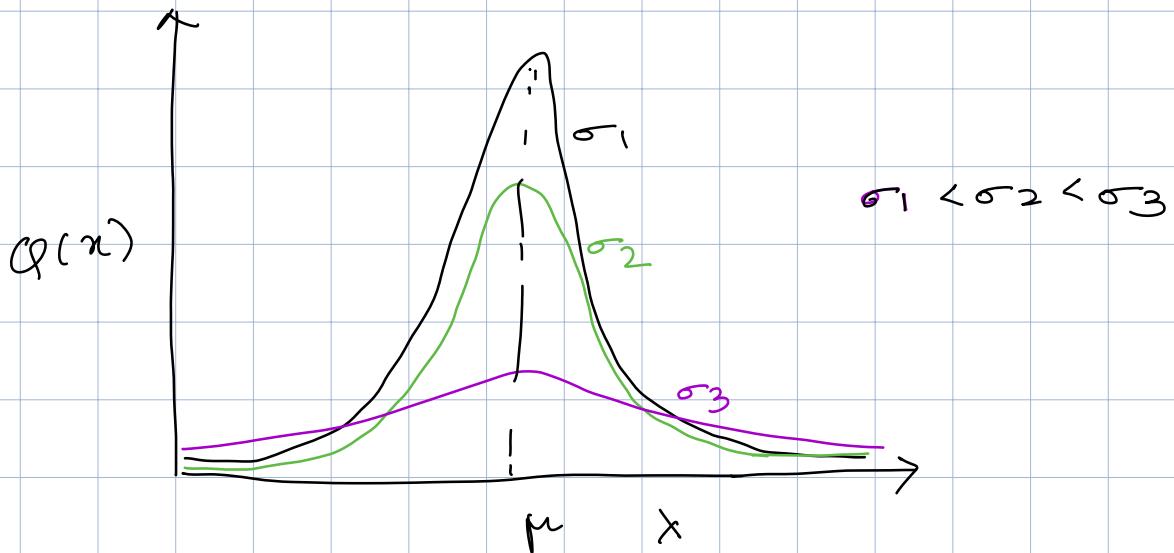
Notation

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

Parameters :- $\mu \in \mathbb{R}$ = mean

$\sigma^2 \in \mathbb{R} > 0$ = variance

$x \in \mathbb{R}$



Example : IRIS Dataset : Petal length, sepal length, Petal width, Sepal width.

For Gaussian distribution : PDF = $\frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x_i - \mu}{\sigma} \right)^2}$

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

Mean

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

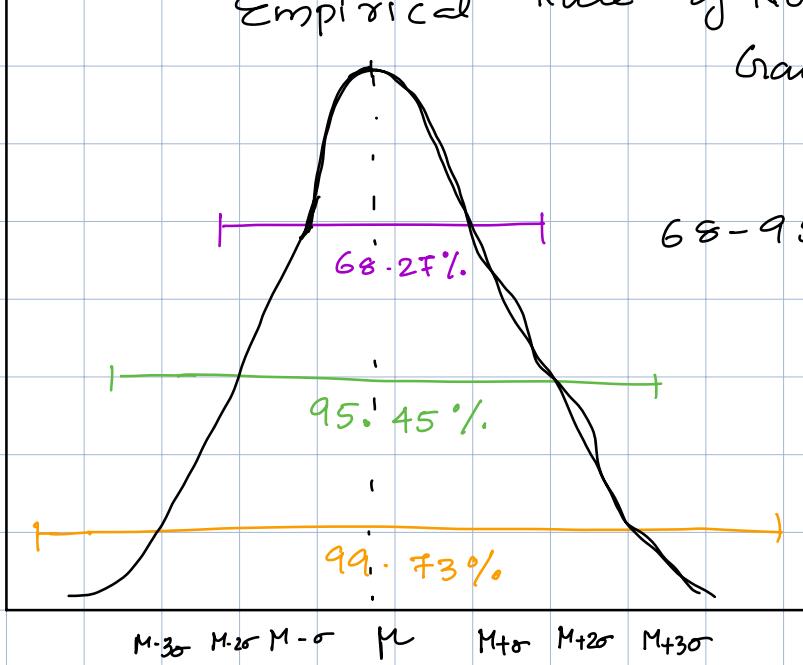
Variance

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

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

Empirical Rule of Normal/Gaussian Distribution

68-95-99.7 Rule



$X = \{ \dots \} \Rightarrow$ Normal / Gaussian Distribution

Probability

$$P(\mu - \sigma \leq X \leq \mu + \sigma) = 68\%$$

$$P(\mu - 2\sigma \leq X \leq \mu + 2\sigma) = 95\%$$

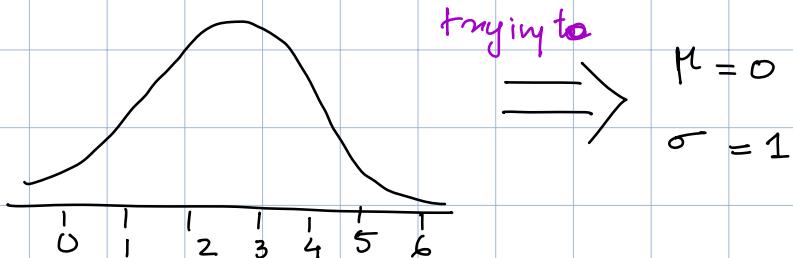
$$P(\mu - 3\sigma \leq X \leq \mu + 3\sigma) = 99.7\%$$

Standard Normal Distribution

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

$$\mu = 3$$

$$\sigma = 1.414 \approx 1$$



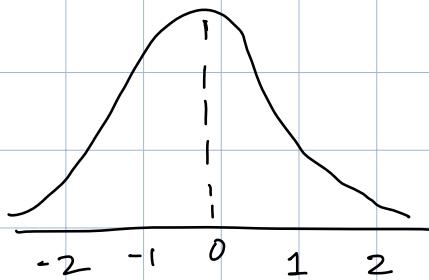
If $\mu = 0$ and $\sigma = 1$ for Gaussian distribution
then we call it as Standard Normal distribution

Every Gaussian distribution can be converted
to standard Normal,
 \Downarrow

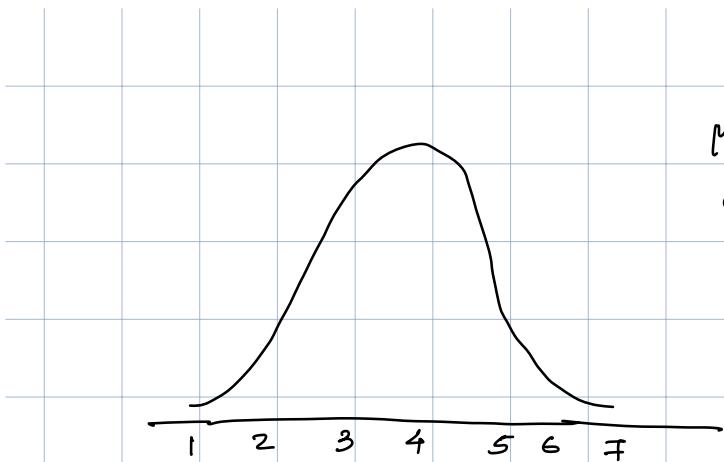
Transformation is done by using Zscore

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

$$x = \{1, 2, 3, 4, 5\} \quad Z = \frac{x_i - \mu}{\sigma} \quad \begin{matrix} \mu = 3 \\ \sigma = 1 \end{matrix}$$
$$\Downarrow \quad x_i - \mu / \sigma$$
$$y = \{-2, -1, 0, 1, 2\}$$



$$X \approx \text{SND } (\mu = 0, \sigma = 1)$$



$$\mu = 4$$

$$\sigma = 1$$

How many standard deviation 4.25 is away from the mean.

$$x_i = 4.25$$

$$Z = \text{Score} = \frac{4.25 - 4}{1} = 0.25 //$$

i.e. 4.25 is 0.25σ way from the mean.

Why Z-score / Standard Normal Distribution?

Eg	Dataset	years	kg	cms	INR
	Age	24	70	175	40k
		25	60	160	50k
		26	55	180	60k
		27	40	130	30k
		30	30	175	20k
		31	35	180	70k

To bring all data points into same unit to make better utilization by ML algorithms \Rightarrow Standardization

Standardization \Rightarrow uses Z score in every feature

$$Z \text{ score} = \frac{x_i - \text{Mean}}{\sigma_{\text{age}}} \quad , \quad \frac{x_i - \text{weight}}{\sigma_{\text{weight}}}$$

\Rightarrow Better model efficiency & performance

Uniform Distribution.

(1) Continuous Uniform Distribution (pdf)

(2) Discrete uniform Distribution (pmf)

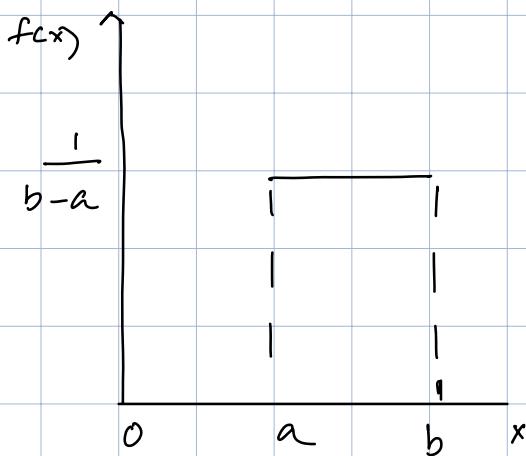
(1) Continuous Uniform Distribution

In probability theory and statistics, the continuous

uniform distributions or rectangular distribution are a family of symmetric probability distribution.

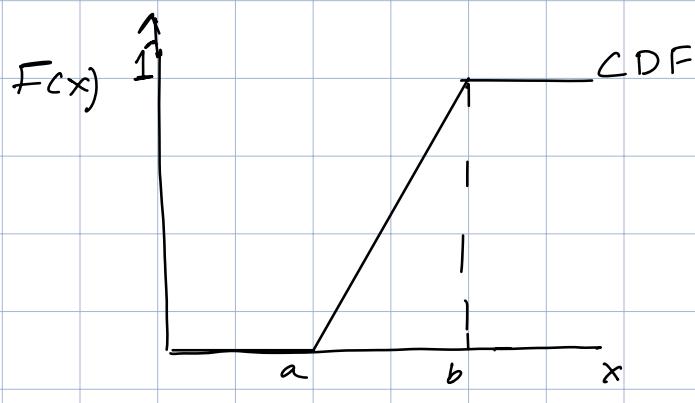
Such a distribution describes an experiment where there is an arbitrary outcome that lies between certain bounds. The bounds are defined by the parameters, a and b , which are the minimum and maximum values

Notation : $U(a, b)$



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

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



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

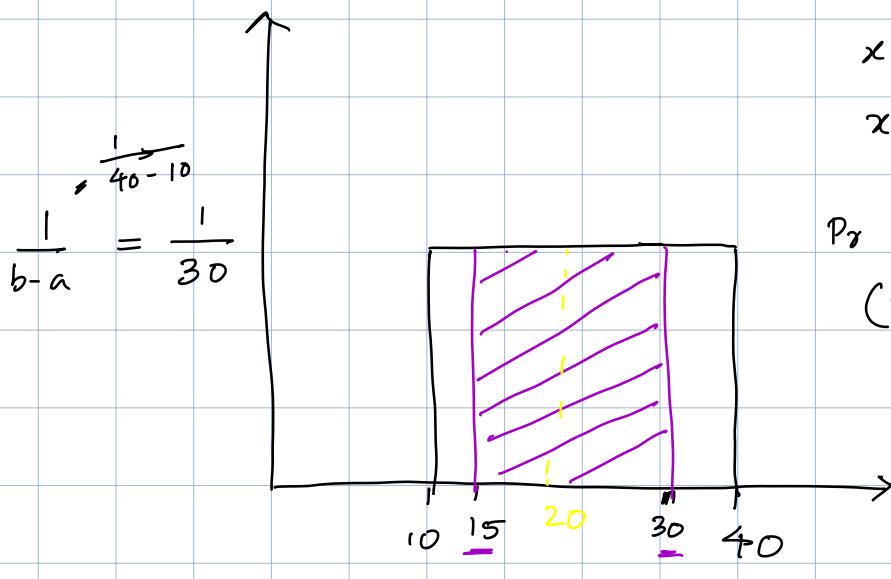
$$\text{Mean} = \frac{1}{2}(a+b)$$

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

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

Eg :- The number of candies sold daily at a shop is uniformly distributed with a maximum of 40 candies and minimum of 10 candies.

(i) Probability of daily sales to fall between 15 and 30?



$$P_r(15 \leq x \leq 30) =$$

$$(x_2 - x_1) * \frac{1}{b-a}$$

$$= (30 - 15) * \frac{1}{30} = \underline{\underline{0.5}}$$

$$\Pr(x \geq 20) = (40-20) \times \frac{1}{30} = 66\% = 0.66$$

Discrete Uniform Distribution

④ Discrete Random Variable

④ Pmf

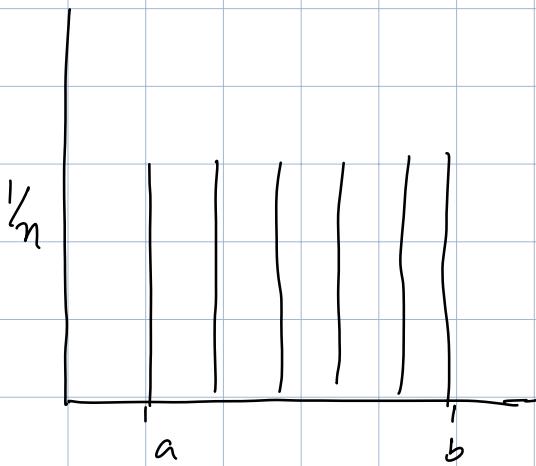
Eg: Rolling a dice $\{1, 2, 3, 4, 5, 6\}$

$$\Pr(1) = \frac{1}{6}$$

$$\Pr(2) = \frac{1}{6}$$

$$\Pr(3) = \frac{1}{6}$$

$$\vdots$$



$$\frac{1}{n} \Rightarrow n = b - a + 1 \\ = 6 - 1 + 1 = 6$$

Notation $U(a, b)$

Parameters a, b where $b > a$

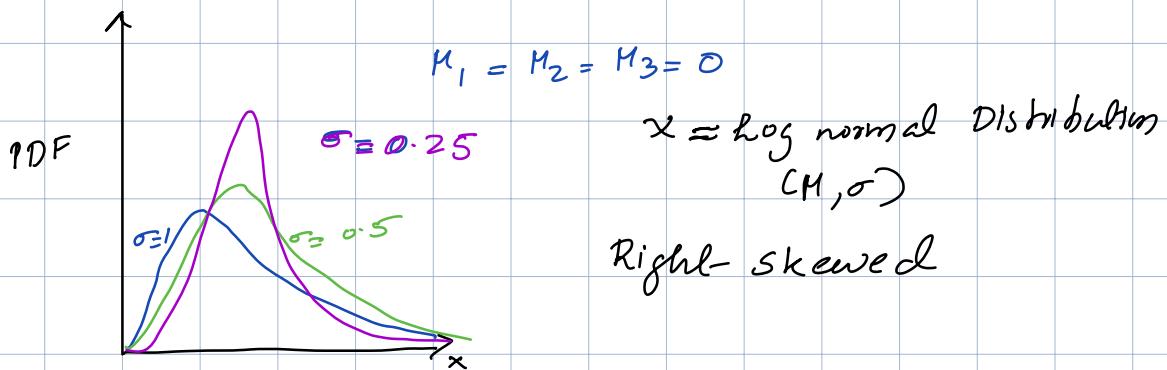
$$PMF = \frac{1}{n}$$

In probability theory and statistics, the discrete uniform distribution is a symmetric probability distribution wherein a finite number of values

are equally likely to be observed; every one of n 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."

Log Normal Distribution

In probability theory, a log-normal (or lognormal) distribution is a continuous probability distribution of R.V whose logarithm is normally distributed. Thus, if the random variable X is a 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.



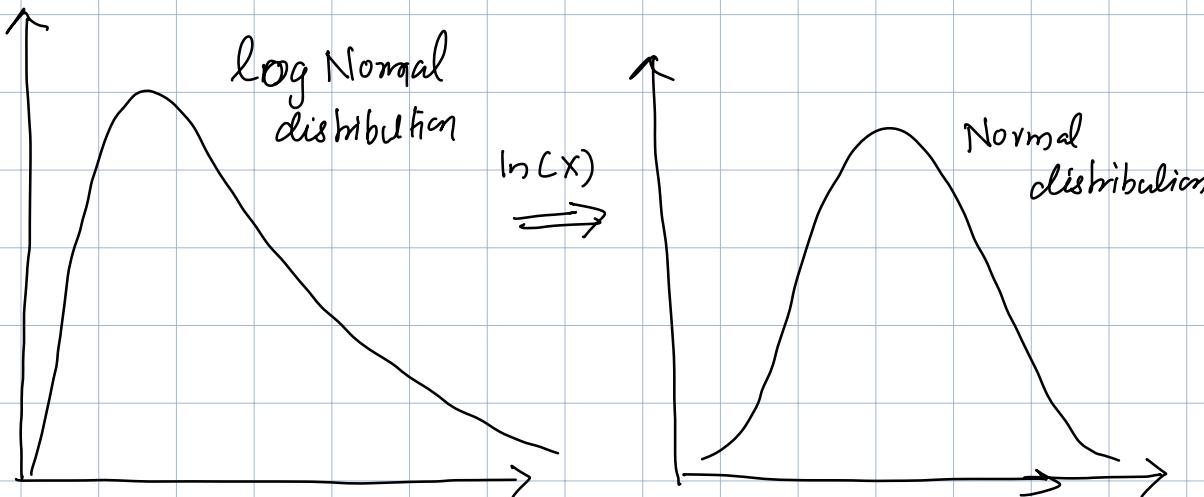
$Y \approx \ln(x)$ = Normal Distribution.

↓
natural log
↓
 \log_e

If $Y \approx \ln(x)$

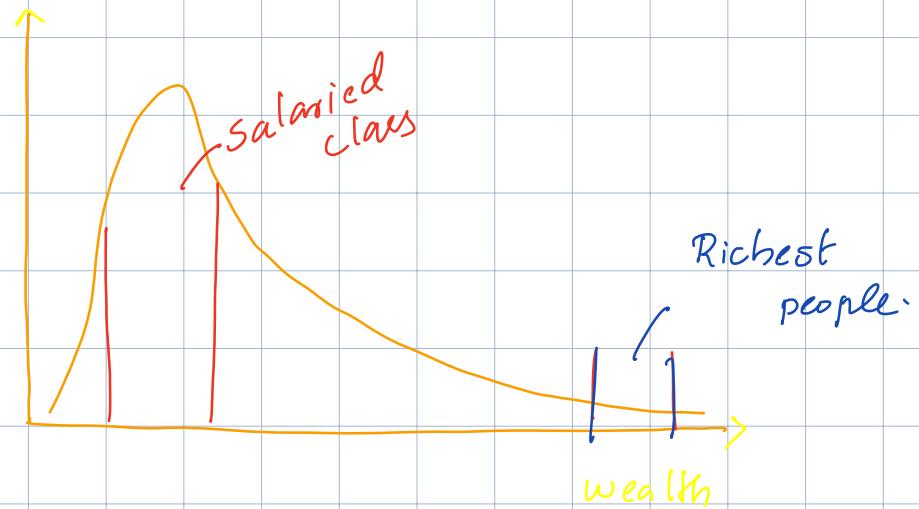
↓

$X \approx \exp(Y) \Rightarrow$ log Normally Distributed



Q-Q plot used to check log normal distribution

Eg: ① health distribution of the world



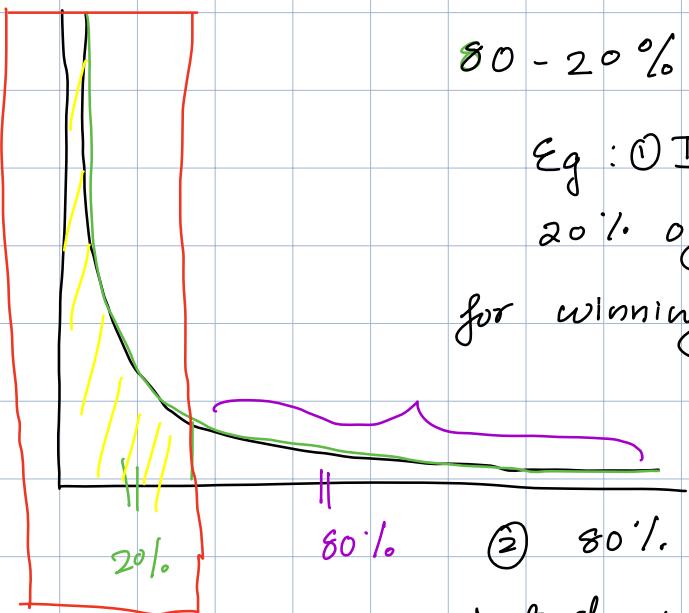
(2) Discussion Forum \Rightarrow lengths of the comments

<u>long</u>	<u>comments</u>	<u>Medium sized comment</u>	<u>small</u>
few		more people	large number

③ Time Spend in Websites / Dwell time.

Power law Distribution

In statistics, a power law is a functional relationship between two quantities, where a relative change in one quantity results in a proportional relative change in the other quantity, independent of the initial size of those quantities. One quantity varies as a power of another.



80 - 20 % Rule

Eg : ① I Ph

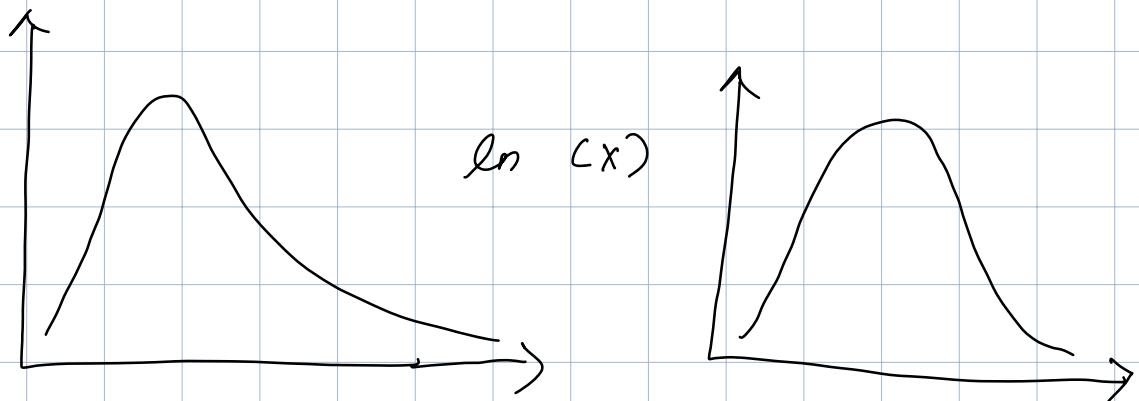
20% of team is responsible
for winning 80% of matches.

② 80% of wealth are distributed with 20% of the total population

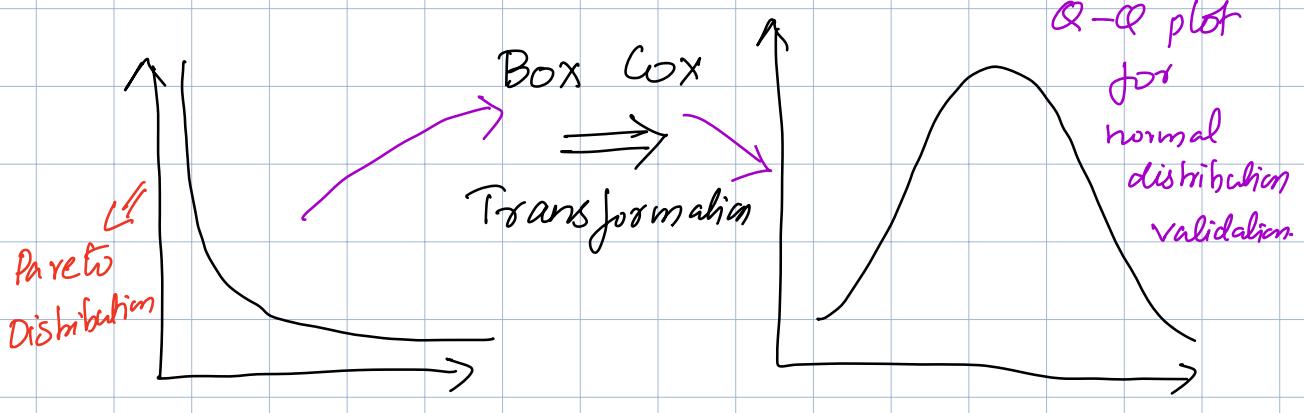
③ 80% of the total oil is with 20% of the nation.

(4) Frequency of words in most languages

(5) 20% of major defects fixes the 80%.
of incoming defects in software product



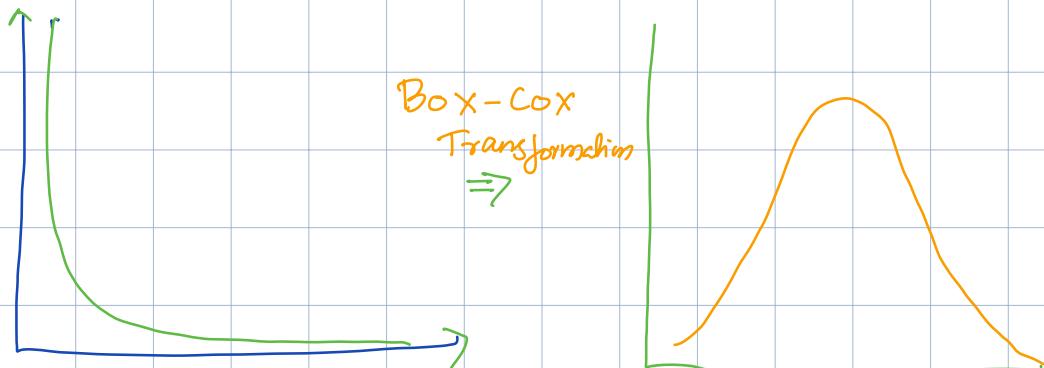
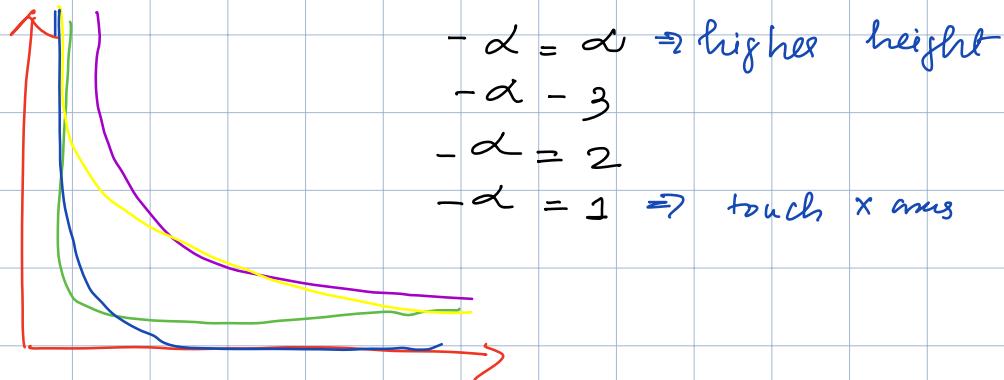
Normal distribution because many ML
algorithms performs well with normal
distribution.



Pareto Distribution

Any distribution which follow power law (80-20 rule) is called Pareto Distribution.

⇒ Non gaussian distribution



Eg :- IT Industry \Rightarrow 80% of entire project is done by 20% of teams.

② 80% of defects can be solved if we solve 20% of defects

Central limit theorem

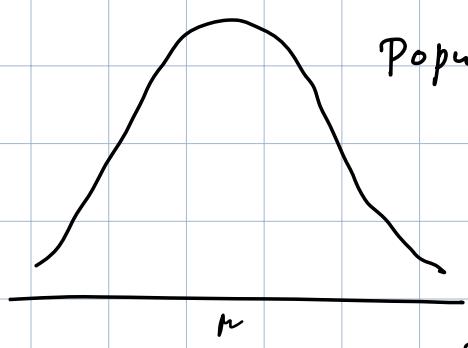
The central limit theorem relies on the concept of a sampling distribution, which is the probability distribution of a statistic for a large number of samples taken from a population.

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

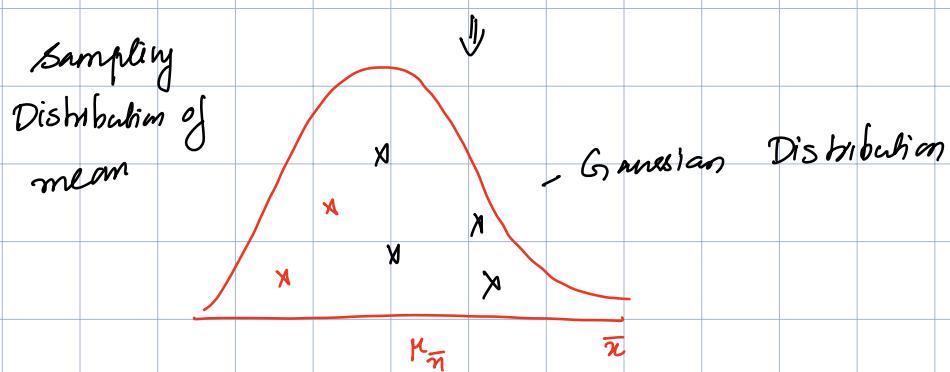
$$\textcircled{1} \quad X \approx N(\mu, \sigma)$$

$n = \text{sample size} = \text{any value}$

Population DATA



$$\begin{aligned} S_1 &= \{x_1, x_2, \dots, x_n\} = \bar{x}_1 \\ S_2 &= \{x_2, x_3, \dots, x_n\} = \bar{x}_2 \\ S_3 &= \vdots \\ x_i &: \\ s_m &= \bar{x}_m \end{aligned}$$



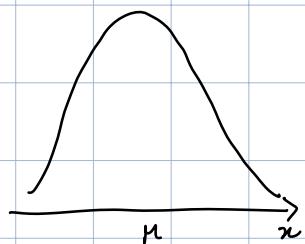
if R.V is not following a gaussian distribution then

n -value should ≥ 30 ($n \geq 30$) sample size.

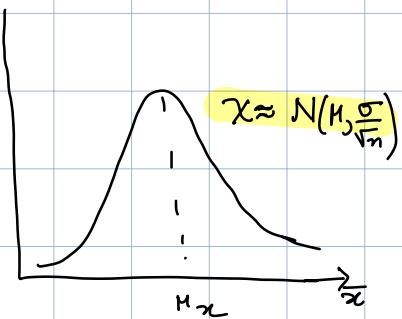
For Normal distribution n can be any value.

① Normal

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



\Rightarrow Sampling distribution of mean



After sampling new mean = population mean

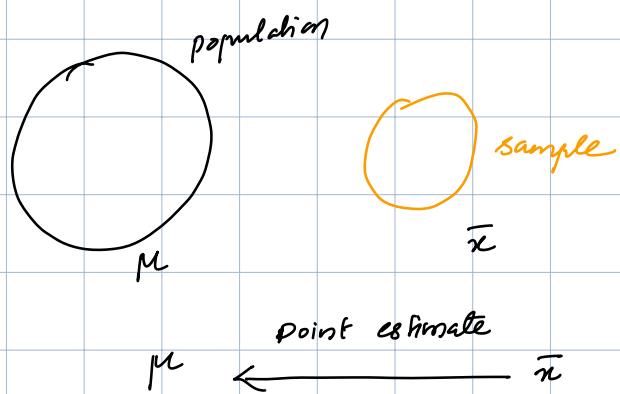
$$\text{new std} = \frac{\sigma}{\sqrt{n}}$$

Estimates

Estimate :- Is specified observed numerical value used to estimate an unknown population parameter.

(1) Point estimate : Single numerical value used to estimate an unknown population parameter

Eg: sample mean is a point estimate of population mean

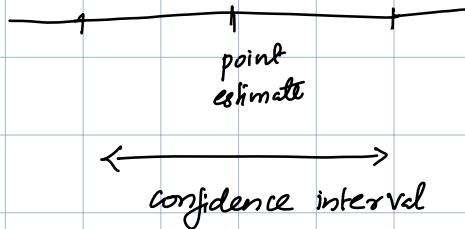


(2) Interval estimate

A range of values is used to estimate unknown

population parameters.

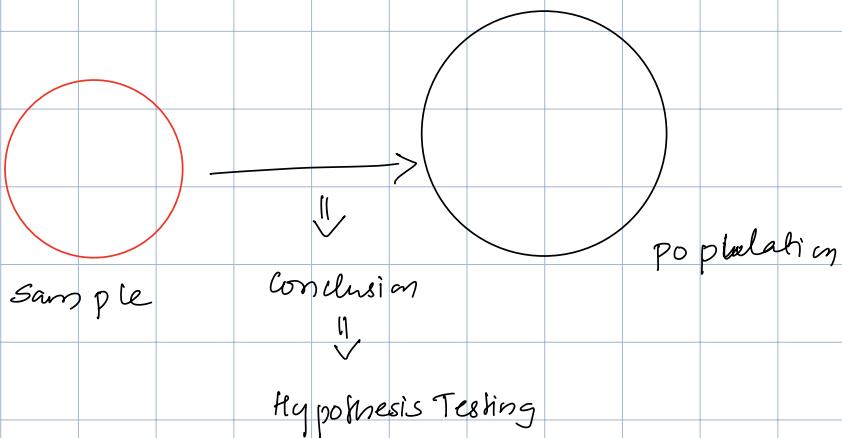
55 - 65



Inferential Statistics

↳ Conclusion or Inferences from sample data.
about population parameter.

A hypothesis and Hypothesis Testing Mechanism



Hypothesis Testing Mechanism

① Null hypothesis (H_0) → Person is not guilty

- The assumption you are beginning with.

② Alternate hypothesis (H_1) - The person is guilty

Opposite of null hypothesis

③ Experiments \rightarrow Statistical Analysis

\rightarrow collect proofs (DNA, finger print)

④ Accept or reject the null hypothesis

Eg: Colleges at District A states its average passed 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%.

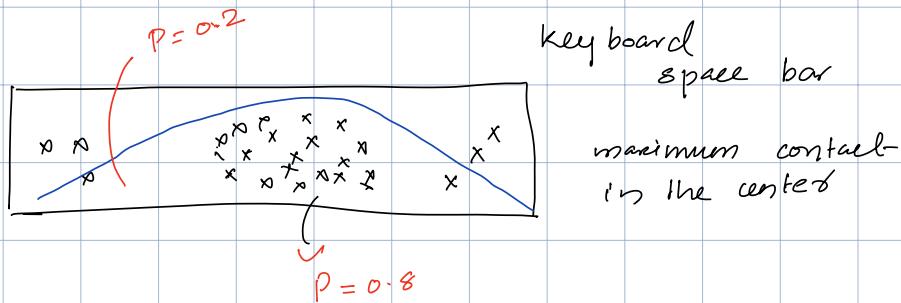
Does this college have a different passed percentage.

Ans). Null hypothesis (H_0) = $\mu = 85\%$

Alternate hypothesis (H_1) = $\mu \neq 85\%$

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.



Hypothesis Testing

Eg : Coin is Fair or Not { 100 Tosses } { 100 times }

$$P(H) = 0.5 \quad P(T) = 0.5$$

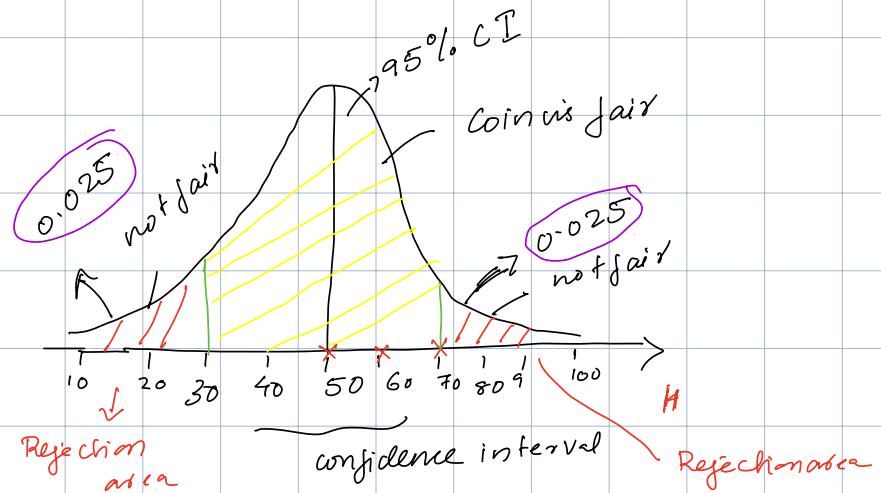
$$P(H) = 0.6 \quad P(T) = 0.4$$

$$P(H) = 0.7 \quad P(T) = 0.3$$

① Null hypothesis (H_0) = Coin is fair

② Alternate hypothesis (H_1) = Coin is not fair

③ Experiment \Rightarrow toss for 100 times



④ significance value : $\alpha = 0.05$

$$\text{Confidence interval} = 1 - 0.05 = 0.95$$

⑤ Conclusion

If $P <$ significance value
↓
Reject the null hypothesis

Rejection area

Else \Rightarrow Fail to reject the null hypothesis

Hypothesis Testing and Statistical Analysis

Z Test } Data dealing with average data Z table (zscore, pval)
t test }

chi square \Rightarrow Categorical data

Anova \Rightarrow For variance of data

Z test

Conditions to use Z-test
= (i) population std known (ii) $n \geq 30$

The average height of all residents in a city is 168cm with $\sigma = 3.9$. A doctor believes the mean to be different. He measured the height of 36 individuals and found the average height to be 169.5cm.

(a) State null and Alternate Hypothesis

(b) At a 95% confidence level, is there enough evidence to reject the null hypothesis.

$$\text{Ans}) \quad \mu = 168\text{cm} \quad \sigma = 3.9$$

$$n = 36$$

$$\bar{x} = 169.5\text{cm}$$

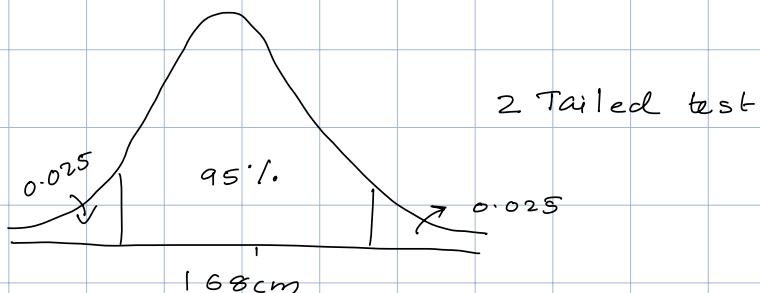
$$C.I = 0.95$$

$$\alpha = 1 - C.I = 0.05$$

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

② Alternate hypothesis (H_1) $\Rightarrow \mu \neq 168\text{cm}$

③ Based on C.I we will draw decision boundary



$$1 - 0.025 = 0.9750 \Rightarrow Z\text{-score}$$



Z-table pdf

↳ check the value 0.9750



Area $\Rightarrow +1.96$

similarly in left $= -1.96$

if Z is less than -1.96 or greater than $+1.96$, Reject the Null hypothesis

Z will be calculated using statistical analysis

④ Z-test

$$Z_d = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$

$$Z\text{-Score} = \frac{\bar{x} - \mu}{\sigma}$$

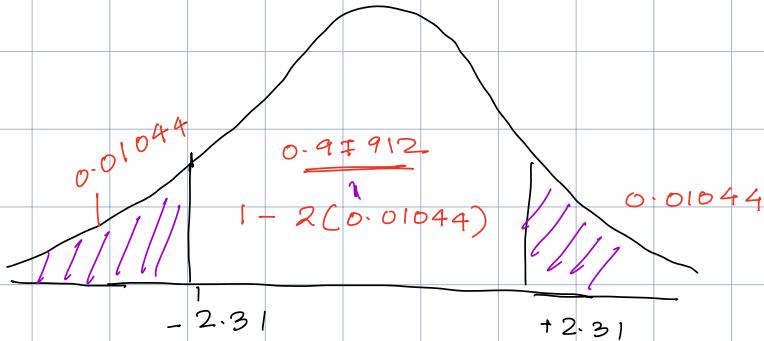
↓
↳ sample std according to CLT

$$= \frac{169.5 - 168}{3.9 / \sqrt{36}} = 2.31$$

2.31 is greater than 1.96

$2.31 > 1.96 \Rightarrow$ Reject the null hypothesis

$$P < 0.05$$



again go to Z-table

2.31 correspond to area

0.98956

so shaded area is $1 - 0.98956 = 0.01044$

$$P \text{ value} = 0.01044 + 0.01044$$

$$= 0.02088 \quad P < 0.05$$

$= 0.02088 < 0.05 = \text{reject null hypothesis}$

(*) Final conclusion the average height $\neq 168\text{cm}$

The average height seems to be increasing based on sample data.

Question

A factory manufactures bulbs with a average warranty of 5 years with standard deviation of 0.50.

A worker believes that the bulb will malfunction in less than 5 years. He test a sample of 40 bulbs and find the average time to be 4.8 years.

(a) State null and alternate hypothesis

(b) At a 2% significance level, is there enough evidence to support the idea that the warranty should be revised?

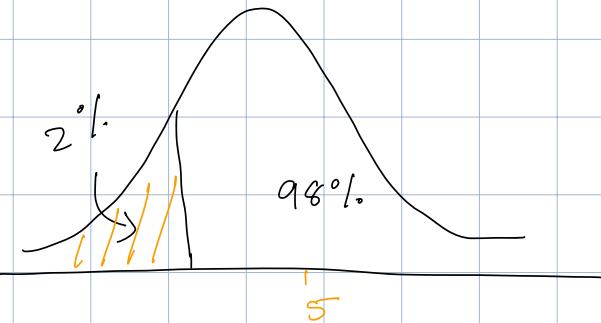
a) ① Null hypothesis $\mu = 5$

given $\sigma = 0.50$

② Alternate hypothesis is $\mu < 5$ { 1 Tail test}

given $n = 40$

b) decision boundary = one tail test



c) Z-test

$$Z_d = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{4.8 - 5}{0.5 / \sqrt{40}} \approx -2.53$$

Area under the curve with z score

$$-2.53 = 0.0570$$

$$P\text{-value} = 0.0570$$

Compare p-value with α

$$0.0570 < 0.02 \Rightarrow \text{False}$$

P value is not falling in the rejection area

We accept the null hypothesis

or failed to reject null hypothesis

Student test

In Z stats when we perform any analysis using Z-score we require σ population standard deviation → is already known

How do we perform any analysis when we do not know the population standard deviation?

Here we use student's t distribution.

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

population std

Z table

used here

$$t = \frac{\bar{x} - \mu}{s / \sqrt{n}}$$

sample std



t table

also need
degree of freedom

Degree of freedom

$$dof = n - 1$$

let's say 3 people



chairs

→ 1st selection from
3 option

[ie when 3 people
were there only two
people were having options.]



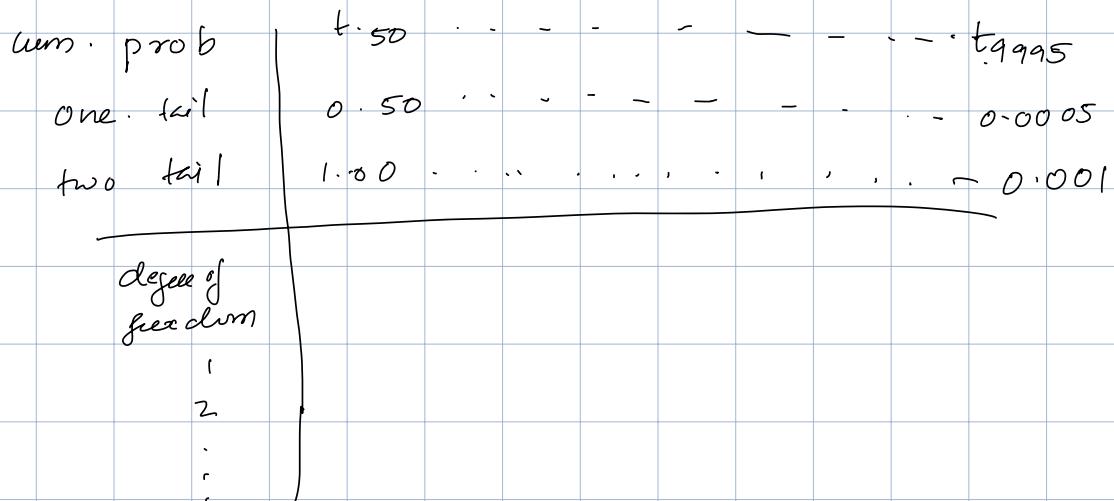
→ 2nd person selects
from 2 option



3rd person has no
option

$$\text{so } dof = 3 - 1 = 2$$

T table has



T-stats and T test \rightarrow One sample t-test

① In the population the average IQ is 100. A team of researchers want to test a new medication to see if it has either a positive or negative effect on intelligence, or no effect at all. A sample of 30 participants who have taken the medication has a mean of 140 with a standard deviation of 20. Did the medication affect intelligence.

\Rightarrow either positive or negative effect on intelligence \Rightarrow Two tail test

$$\text{Ans} \Rightarrow \mu = 100$$

$$n = 30$$

$$\bar{x} = 140$$

$$s_d = 20$$

$$\text{Initialize } C.I = 95\% \quad \alpha = 0.05$$

① Null hypothesis H_0 $\mu = 100$

Alternate hypothesis H_1 $\mu \neq 100$

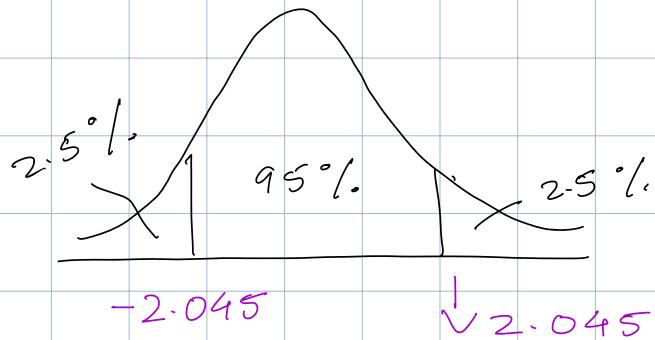
(two tailed test)

② $\alpha = 0.05$

③ Degree of freedom =

$$n-1 = 30-1 = 29$$

④ Decision rule



Go to t-table

1) 2 tail with df = 29

$\alpha = 0.05$ with $d = 29$

↓
2.045

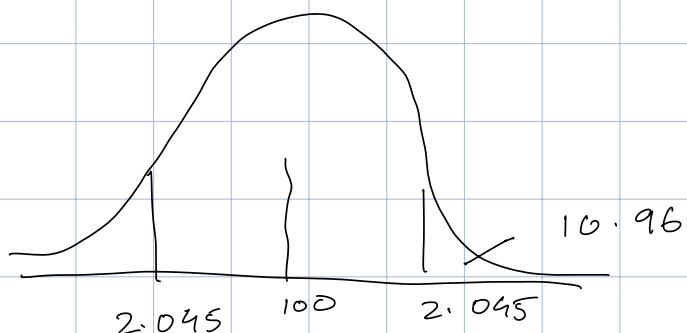
The assumption here is that if my t-test is less than -2.045 and greater than 2.045, reject the null hypothesis.

⑤ Calculating Test Statistics

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{140 - 100}{20/\sqrt{30}} = \frac{40}{3.65} = 10.96$$

$$10.96 > 2.045$$

so



10.96 falls in the rejection area

Conclusion

Medication used has affected the intelligence

Specifically medication has increased the IQ.

When to use T-test vs Z-test

Do you know the population

Std σ

✓ Yes

No \downarrow

Is the sample size
above 30

use t-test

Yes \downarrow

Z-test

No \downarrow

use t-test

Type I and Type 2 Errors

Reality : Null hypothesis is True or Null hypothesis is False

Decision : Null hypothesis is True or Null hypothesis is False

Outcome 1 :- We reject the null hypothesis when in reality it is false \rightarrow Good

Outcome 2 : We reject the null hypothesis when in reality it is true \Rightarrow Error
 \Downarrow
Type I Error

Outcome 3 : We retain the Null hypothesis, when in reality it is False
 \Downarrow
Type II error

Outcome 4 :- we retain the Null hypothesis when in reality it is True \Rightarrow Good / correct scenario

Bayes theorem (Bayes statistics)

Bayesian statistics is an approach to data analysis and parameter estimation based on Baye's theorem.

Baye's Theorem

Probability \rightarrow Independent events
 \rightarrow Dependent events

(1) Independent event

e.g.: Rolling a dice
 $\{1, 2, 3, 4, 5, 6\}$

(2) Dependent event



Eg Tossing a coin

{ H, T }

$$P(\text{red}) = \frac{2}{5} = 1^{\text{st}} \text{ event}$$

$$P(\text{yellow}) = \frac{3}{4} = 2^{\text{nd}} \text{ event}$$

One after the other
and there is dependency

$$\text{so } P(\text{Red} \text{ and Yellow}) =$$

$$P(R \& Y) = P(R) * P(Y|R)$$

||
conditional probability

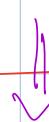
$$= \frac{2}{5} \times \frac{3}{4} = \frac{6}{20}$$

In general :

$$P(A \text{ and } B) = P(B \text{ and } A)$$

$$P(A) * P(B/A) = P(B) * P(A/B)$$

$$P(B/A) = \frac{P(B) * P(A/B)}{P(A)}$$



Baye's theorem

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

A, B = events

$P(A|B)$ = Probability of A given B is true

$P(A), P(B)$ = Independent probabilities of A and B

DATA SET

I/P
⇒ Independent feature

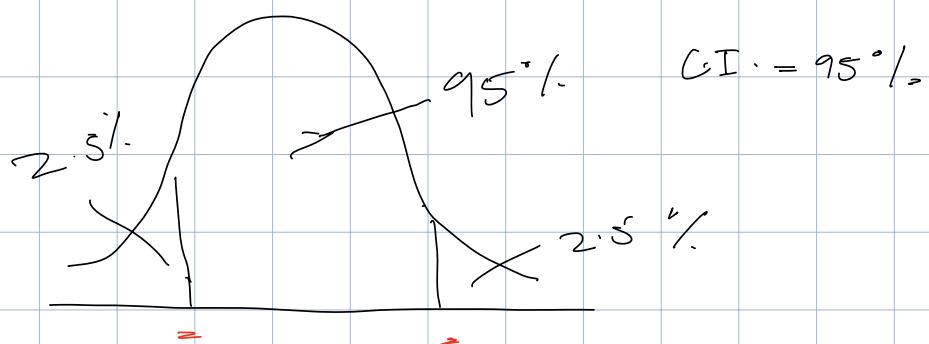
O/P
dependent
= feature

Size of house	Number of Rooms	Location	Price
x_1	x_2	x_3	y

Bayes theorem used here

$$P(Y/x_1, x_2, x_3) = \frac{P(Y) * P(x_1, x_2, x_3/Y)}{P(x_1, x_2, x_3)}$$

Confidence Intervals and Margin of Error



Point estimate

$$\bar{x}$$

$$\bar{x} = 2.5$$

$$\mu$$

$$\mu = 3$$

not always same

Confidence interval

Point estimate \pm margin of error

$$\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

↑
for z test

$$\bar{x} \pm t_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

for t test

Eg: On the verbal section of CAT exam,

The standard deviation is known to be 100.

A sample of 30 test takers has a mean of 520. Construct 95% CI about the mean.

$$\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

$$\alpha = 0.05$$

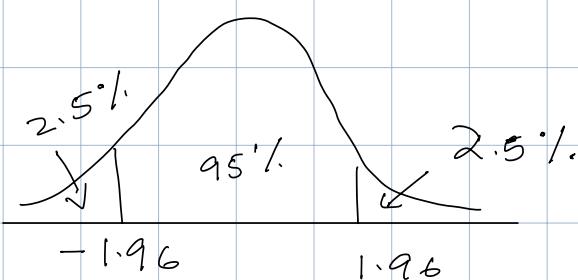
$$1 - 0.025 = 0.9750$$

↓

Z_{table}

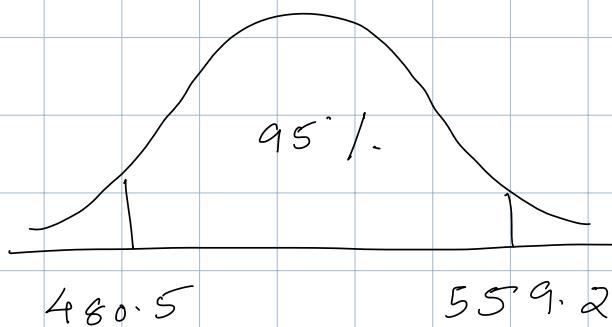
||

+1.96



$$\text{Lower C.I} = 520 - (1.96) \frac{100}{\sqrt{25}} = 480.5$$

$$\text{Higher C.I} = 520 + (1.96) \frac{100}{\sqrt{25}} = 559.2$$



So CAT score will be 95% confidently
between 480.5 and 559.2

CHI SQUARE TEST

⇒ for categorical variables

⇒ Chi-square test for goodness of fit test

Claims about population proportion

It is a non parametric test that is performed

on the categorical data. [ordinal and nominal]

There is a population of male who likes different colour bikes.

(Prob) known theory
about population Sample

Yellow Bike

$\frac{1}{3}$

22

Red Bike

$\frac{1}{3}$

17

Orange Bike

$\frac{1}{3}$

59

We have to see if the sample information is a good fit for population

* Goodness of fit test

Goodness of fit test

In a science class of 75 students, 11 are left handed. Does this class fit the theory that 12% of people are left handed.

	<u>Observed</u>	<u>Expected</u>
Ans) Left handed	11	$\frac{12}{100} \times 75 = 9$

Right handed	<u>64</u> 75	<u>66</u> 75
--------------	-----------------	-----------------

Chi-square test we are comparing 2 categorical variables

CHI SQUARE FOR GOODNESS OF FIT

In 2010 census of the city, the weight of the individuals in a small city were found to be the following

$\leq 50\text{kg}$	$50 - 75$	> 75
20%	30%	50%

In 2020, weight of $n = 500$ individuals were sampled. Below are the results

≤ 50	$50 - 75$	> 75
140	160	200

Using $\alpha = 0.05$, would you conclude the population difference of weights has changed in the last 10 years?

Ans) We use Chi-Square table here

1st take 2010 data

$\leq 50\text{kg}$	$50 - 75$	> 75	Expected
20%	30%	50%	

+ take 2020 observed data $n=500$

≤ 50	$50 - 75$	> 75
140	160	200

≤ 50	$50 - 75$	> 75	Expected info based 2010 scnus.
0.2×500	0.3×500	0.5×500	
110	150	250	

Apply Chi - Square

① Null hypothesis: H_0 : The data meets the expectation

Alternate hypothesis H_1 : The data does not meet the expectation

② $\alpha = 0.05$ C. I = 95%

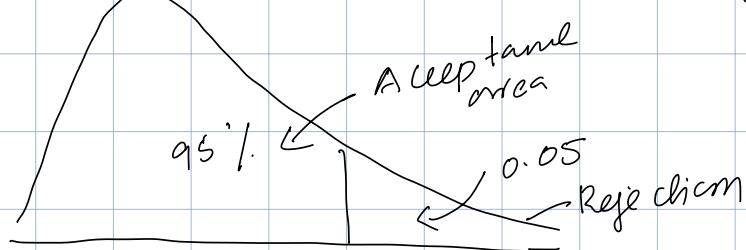
③ Degree of freedom

$$df = k - 1 = 3 - 1 = 2$$

number of
categories

④ Decision boundary [In Z, t we used symmetrical distribution)

But we use a one tailed statistic



chi-square Probability

A leap ^{tand}
area

0.05
Reject H₀

Critical
value

5.991 (by checking chi-square
table)

if χ^2 is greater than 5.99, Reject H_0

else

we fail to reject H_0

⑤ Calculate Chi-Square Test Statistics

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

$$= \underbrace{(140 - 100)^2}_{100} + \underbrace{(160 - 150)^2}_{150} + \underbrace{(200 - 250)^2}_{250}$$

$$= \frac{1600}{100} + \frac{100}{150} + \frac{2500}{250}$$

$$= 16 + 0.66 + 10 = 26.66$$

Decision rule was if χ^2 is greater than

5.99 Reject H_0

Here $26.66 > 5.99$ so Reject H_0

Answered

The weights of 2020 population are different than those expected in 2010 population.

What is ANOVA?

Analysis of Variance

Definition: ANOVA is a statistical method used to compare the means of 2 or more groups.

ANOVA

① Factors (variable)

② Levels

Eg: Medicine (Factor)

[Dosage] 5 mg 10 mg 15 mg [level]

Eg) Mode of Payment [Factor]
GPA Y PHONGPE PAYPAL [Levels]

Assumptions in ANOVA

① Normality of Sampling Distribution of Mean

The distribution of sample mean is normally

distributed

(2) Absence of Outliers

Outlier score need to be removed from dataset

(3) Homogeneity of Variance

Population variance in different levels of each independent variable are equal

$$[\sigma_1^2 = \sigma_2^2 = \sigma_3^2]$$

(4) Samples are independent and randomly selected

Types of ANOVA

3 types

1) One way ANOVA: One factor with at least 2 levels and these levels are independent

Eg :- Doctor wants to test a new medication to decrease headache.
They split the participants in 3 conditions
[10mg, 20mg, 30mg]

Doctor ask the participants to rate the headache between [1 - 10]

Medication \rightarrow Factor

10mg 20mg 30mg

5

7

2

3

4

6

-

-

-

-

-

-

Type II) Repeated measures ANOVA :-

one factor with atleast 2 levels, levels
are dependent.

Running \rightarrow Factor

Day 1

Day 2

Day 3

\leftarrow levels

8

5

4

7

4

9

-

-

-

③ Factorial ANOVA:

Two or more factors each of which with at least 2 levels, levels can be independent and dependent.



Hypothesis Testing In ANOVA

(Partitioning of Variance in the ANOVA)

1) Null hypothesis $H_0: \mu_1 = \mu_2 = \mu_3 = \dots = \mu_k$

Alternate hypothesis $H_1: \text{at least one of sample mean is not equal}$

Test statistics

$F = \frac{\text{Variance between samples}}{\text{Variance within sample}}$

Variance between samples

Variance within samples

Samples =	\bar{x}_1	\bar{x}_2	\bar{x}_3
	1	6	5
	2	7	6
variance within samples	4	3	3
	5	2	2
	3	1	1
	$\bar{x}_1 = 3$	$\bar{x}_2 = 19/5$	$\bar{x}_3 = 16/5 = 4$

$$H_0 : \bar{x}_1 = \bar{x}_2 = \bar{x}_3$$

H_1 atleast one sample mean is not equal

One way ANOVA

One factor with at least 2 levels, levels are independent

- ① Doctors want to test a new medication which reduces headache. They split the participants into 3 condition [15mg, 30mg, 45mg].
Later on the doctor ask the patient to rate the headache between [1-10]. Are there any differences between 3 conditions

$$\alpha = 0.005?$$

	15mg	30mg	45mg
Any	9 8 7 8 8 9 8	7 6 6 7 8 7 6	4 3 2 3 4 3 2

① Define null and alternate hypothesis

$$H_0: \mu_{15} = \mu_{30} = \mu_{45}$$

H_a: at least one mean is not equal

② Significant value

$$\alpha = 0.05$$

$$C.I = 95\%$$

③ Degree of freedom

$$N = 21 \quad [7 \times 3]$$

\downarrow \downarrow
rows columns

$a = 3 \Rightarrow$ categories (levels)

$n = 7$ samples in each category

$$df_{\text{between}} = a - 1 = 3 - 1 = 2$$

?

$$df_{\text{within}} = N - a = 21 - 3 = 18$$

$$df_{\text{total}} = N - 1 = 20$$

$$df_{\text{btwn}} + df_{\text{within}} = df_{\text{total}}$$

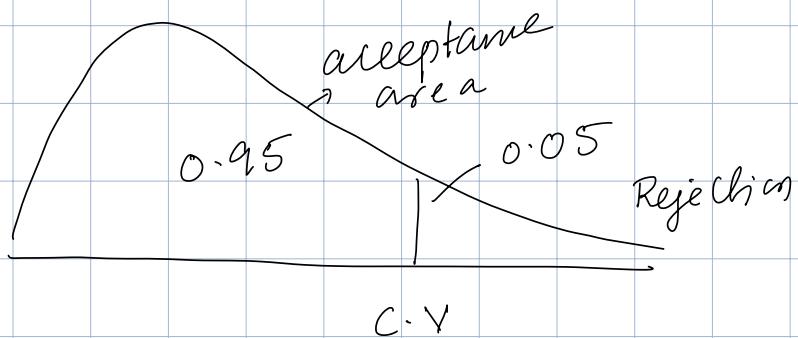
In Ftable

$$df_1 \quad df_2 \\ (21, 18)$$

$$\alpha = 0.05$$

critical value

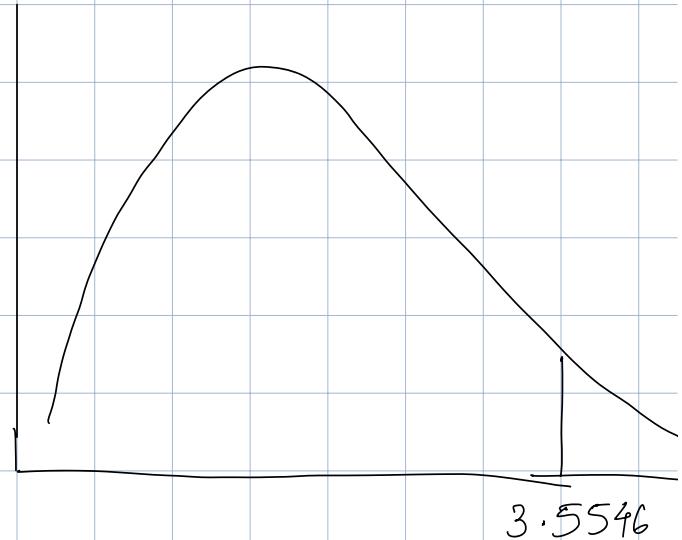
③ Decision Boundary



(5) F table

check $\alpha = 0.05$, $(2, 18)$
 df_1 , df_2
between, within

$$\alpha = 0.05$$



Decision Rule

If F is greater than 3.5546,

reject the Null hypothesis

(5) Calculate the F test statistics

$$F = \frac{\text{Variance between Sample}}{\text{Variance within Sample}}$$

	Sum of squares	df	M S	F
Between				
Within				
Total				

	15mg	30mg	45mg
9			4
8		6	3
7		6	2
8		7	3
8		8	4
9		7	3
8		6	2

$$S.S_{\text{between}} = \frac{\sum (\bar{x}_i)^2 - \bar{x}^2}{n}$$

$$\bar{x}_i$$

$$15mg \Rightarrow 9 + 8 + 7 + 8 + 8 + 9 + 8 = 57$$

$$30mg \Rightarrow 7+6+6+7+8+7+6 = 47$$

$$45mg \Rightarrow 4+3+2+3+9+3+2 = 21$$

$$S.S_{\text{between}} = \sum \frac{(\sum a_i)^2}{n} - \frac{\bar{Y}^2}{N}$$

$$= \frac{57^2 + 47^2 + 21^2}{7} - \frac{[57^2 + 47^2 + 21^2]}{21}$$

$$= 98.67$$

$$\begin{aligned} \textcircled{2} \quad S.S_{\text{within}} &= \sum y^2 - \sum \frac{(\sum a_i)^2}{n} \\ &= \sum y^2 = 9^2 + 8^2 + 7^2 + 8^2 + 8^2 + \dots \\ &= 853 - \frac{57^2 + 47^2 + 21^2}{7} \\ &= 10.29 \end{aligned}$$

	Sum of squares	df	mean squared	F
Between	98.67	2	49.34	
Within	10.29	18	0.54	
Total	108.96	20		

$$F \text{ test} = \frac{49.34}{0.54} = 86.56$$

If F is greater than 3.5546, we Reject H₀

Here 86.56 > 3.5546 so reject H₀