

Statistics!

→ Statistics is the science of collecting, organizing and analysing data.

Data: facts or pieces of information

It can be measured, collected, analyzed.

Eg:- weights of students in the class

{60kg, 50kg, 45kg, 30kg, ...} Types of Statistics:

IQ of students in class. ① Descriptive Statistics:

{100, 90, 95, 99, ...}

House Price Dataset:

<u>Area</u>	No. of houses	Price
Bengaluru	1000	45L
New York	1250	50L
Mumbai	800	60L

After analysing this data, A Data scientist need to build a model which predicts a price for new houses based on past data.

Data Analyst → Report → visualizations

→ Data Analyst will build the reports and present the visualizations to the stakeholders. They will take the meaningful decisions.

① Descriptive Statistics:

→ This involves methods for summarizing & organizing data to make it understandable. This type of statistics helps to describe the basic features of the data in a study.

(i) Measure of Central Tendency
[Mean, Median, Mode]

(ii) Measure of Dispersion
[Variance, Standard Deviation]

(iii) Data Distribution

- Histograms
- Box Plot
- Pie chart
- PDF, PMF..

(iv) Summary Statistics

- Five Number Summary
- Q_1, Q_2, Q_3

② Inferential statistics

→ This 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.

- i) Hypothesis Testing
- ii) P-value
- iii) Confidence Interval
- iv) Statistical Analysis Test
 - Z-test
 - t-test
 - ANOVA → F-test
 - Chi Square

Eg:- Let say there are 20 statistics class in your college, and you have collected the height of students in the class.

Heights of [175cm, 180cm, 140cm, 135, 160cm, 120cm]

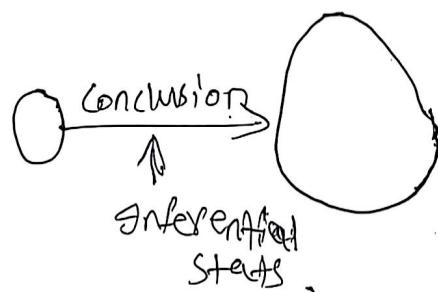
Descriptive question

"What is avg height of entire classroom?"

↑
Measure of Central Tendency.

Inferential Question

"Are the height of the sample students in classroom similar to what you expect in the entire college?"



Population

- It is the entire set of individuals.
- It includes all members of a defined group.

Characteristics

- (1) Complete Set : Containing all the observation of interest.
- (2) Parameter : A numerical value summarizing the entire population,

Population Mean (μ)

Population Variance (σ^2)

Ex:-

- (1) Population in a school study
 - * All students enrolled in a school.
 - * Determine avg height of student, population mean.

Sample Data

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

Characteristics!

① Subset: Represents 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.

Types of Sampling Techniques

1. Probability Sampling

2. Non Probability " "

① Probability Sampling

ⓐ Simple Random Sampling:

Every member of the population has an equal chance of being selected.

Eg: Selecting people randomly.

ⓑ Systematic Sampling:

Select every n th member of the population after a random starting point.

Eg! Airport → credit card

↳ 5th person, 10th person, 15th person ...

→ Feedback Survey → select every 11th member

ⓒ Stratified Sampling

Divide the population into strata (groups) based on specific characteristics & then randomly sampling from each strata.

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

Eg: Age → <12 12-18 >18

{Politics}

(d) Cluster Sampling

→ Divide the population into clusters, randomly selecting clusters, then sampling all the members from the selected clusters.

Eg! - Randomly selecting several schools from a district, & surveying all teachers within those schools.

(e) Snowball Sampling

Existing study subjects recruit future subjects from among their acquaintances.

Eg! - Survey members of a rare disease.

(f) Quota Sampling

Age, group, gender, caste.

(2) Non Probability Sampling

Select individuals who are easiest to reach.

Eg! - Surveying people at mall.

(g) Convenience Sampling

→ Selecting individual who are easiest to reach.

(h) Judgmental (purposive Sampling)

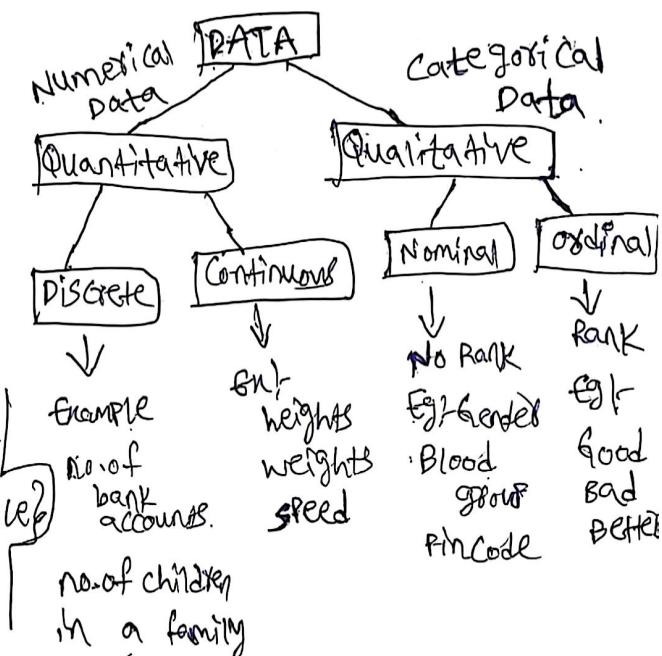
→ Select individual based on the researcher's judgement

Eg! choose experts in a field

to participate {Data Science}

* Selecting the sampling technique depends on use cases.

Types of DATA



Scales of Measurement of Data:

→ This describes the nature of information within the values assigned to variables.

(2) Ordinal Scale

→ This scale classifies the data into categories that can be ranked or ordered.

4 Primary Scales of Measurement

① Nominal Scale

② Ordinal Scale

③ Interval "

④ Ratio "

Characteristics:-

- (i) Data is categorized & ranked in a specific order.
- (ii) The interval b/w ranks are not necessarily equal.

① Nominal Scale:

→ This scale classifies data into distinct categories that do not have an intrinsic order.

→ Qualitative/Categorical data.

Characteristics:-

- (i) Data is categorized based on labels, names or qualities.
- (ii) Categories are mutually exclusive.

Ex:- Gender

- male
- female.
- (iii) No logical order among categories. [No Rank]

Examples:-

<u>Education level</u>	\Rightarrow Feature	\Rightarrow Ranks
High School	Engineering	1

Bachelor	2
----------	---

Master's	3
----------	---

Doctorate	4
-----------	---

Customer Feedback

Satisfied	1
Very Satisfied	2
Not Satisfied	0

Socio Economic Status

Low	0
Medium	1
High	2

③ Interval Scale

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

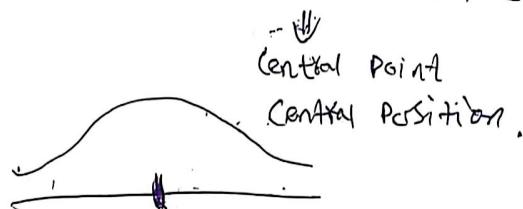
Characteristics:

- (i) Data is ordered with consistent interval b/w values.
- (ii) Allows for meaningful comparison of differences.
- (iii) No true zero point.

Measure of Central Tendency

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

$$\text{Ex:- Age} = [24, 32, 12, 48, 16, 20]$$



④ Ratio Scale

- The order matters.
- Differences are measurable.
- Contains a 0 starting point.
- Ratio can be measured.
Eg:- marks of students
0, -30, 40, 90.

5

0

0

1

-5

① Mean

→ Mean is the sum of all values divided by the no. of values.

Population Mean (μ)

Population (N)

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

Sample mean (\bar{x})

Sample (n)

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

Here x is a random variable

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

$$N = 5$$

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

$n \rightarrow$ sample size.

Characteristics:

- mean is affected by extreme outliers → most frequently appeared value in a dataset.
- used for interval and ratio data

③ Mode:

most frequently appeared value in a dataset.

Ex: 2, 4, 4, 6, 7, 7, 7, 9

mode = 7 (most frequent value)

3, 5, 5, 6, 6, 8

Mode = 5, 6 (bimodal)

② Median:

→ The median is the middle value in a dataset after arranging the values in a ascending, or descending order.

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

no. of elements = 5,

5 is odd

median = 3

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

$$\Leftrightarrow \{1, 2, 3, 4, 5, 100\}$$

no. of elements = 6

6 is even

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

Characteristics:

* Not affected by outliers.

* used for ordinal, interval and ratio data.

Characteristics:

* This is not affected by extreme values.

* used for nominal, ordinal, interval and ratio data.

Choosing the Appropriate Measure

1. Mean: Best used when data is symmetrically distributed without outliers.

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

3. 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.

Measures of Dispersion

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

Common measures of dispersion

- (1) Range
- (2) Variance
- (3) Standard Deviation
- (4) Interquartile Range (IQR)

② Variance:

→ It measures the average squared deviation of each value from the mean. It provides a 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:

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

Ex:- size of flower petals

{5, 8, 12, 15, 20} → Variance of this distribution.

$$N = 5$$

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

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

$$\text{Variance} = \frac{(5-12)^2 + (8-12)^2 + \dots + (20-12)^2}{5}$$

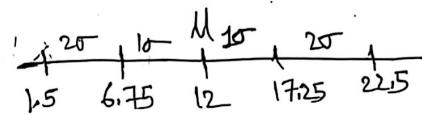
$$\text{Variance} = 27.6$$

③ Standard Deviation:

→ It is the square root of the variance.

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

→ By using standard deviation we can say how far the data point is from mean



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

Where 20 will locate?

→ It is 2σ (two standard deviations) to the right of the mean

① Range:

→ It is the difference b/w the maximum & minimum value in a dataset.

$$\text{Range} = \text{Maximum value} - \text{minimum value}$$

Ex:- Ages {14, 13, 10, 12, 25, 75, 15}

$$\text{Range} = 75 - 10 = 65$$

Characteristics:

- Simple to calculate
- Sensitive to outliers
- Rough measure of dispersion.

Characteristics:

→ Provides a clear measure of spread in the same unit as the data.

→ Sensitive to outliers.

Why Sample Variance is $\frac{1}{n-1}$?

Characteristics:

- Provide a precise measure of variability.
- Units are squared of the original data units.
- More sensitive to outliers than the Range.

Ex:- Ages = { }

Sample



Population me

Mean

$$\bar{x}$$

M

gap will

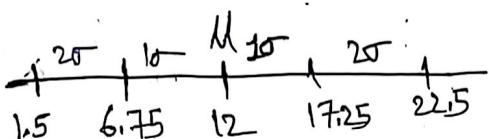
→ When you take σ of the mean of sample

③ Standard Deviation!

→ It is the square root of the variance.

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

→ By using standard deviation we can say how far the data point is from mean.



Where 20 will locate?

→ It is 20 (two standard deviations) to the right of mean.

Characteristics?

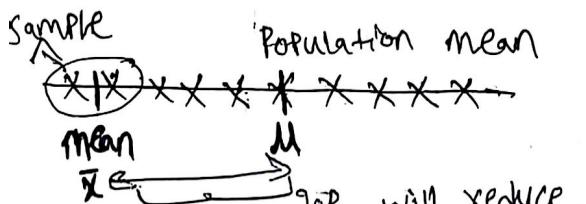
→ Provides a clear measure of spread in the same units as the data.

→ Sensitive to outliers.

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

Ex:-

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



→ When you take a sample, the mean of sample is

at more distance compared to population mean.

$$\bar{x} \ll \mu$$

$$s^2 \ll \sigma^2$$

→ If we divide sample mean with 'n', there is so much difference b/w sample mean & population mean.

→ To do correction & to decrease difference somewhat, Bessel's correction is introduced.

'n-1'

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

understanding the true population mean.

Degrees of freedom = $n-1$

Random Variables!

X → function → Values



Process of experiments

$$y = 5x + 2$$

$x = 1, 2, 3, \dots$ are random numbers variables

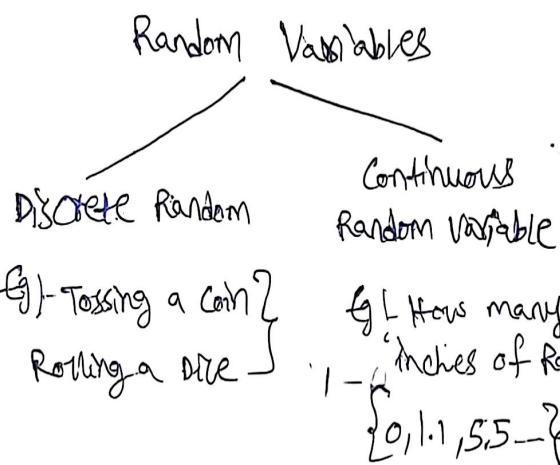
$$X = \begin{cases} 0 & H \\ 1 & T \end{cases} \quad \text{Tossing a coin}$$

$$X = \begin{cases} 1 & \text{1, 2, 3} \\ 2 & \text{4, 5, 6} \end{cases} \quad \text{Rolling a fair dice}$$

Percentiles)

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

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



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

$$\begin{aligned} \text{Percentile of value 9} &= \frac{11}{14} \times 100 \\ &= 78.57\% \text{ of value 9} \end{aligned}$$

↙ Percentile Ranking.

→ This says that 78.57% of entire distribution is less than 9.

→ If value want to find?

Percentiles and Quartiles

Percentile:

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

no. of odd numbers = 3

% of odd numbers
in this group

$$= \frac{3}{6} \times 100 = 50\%$$

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

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

$$\boxed{\text{Value} = 3.75}$$

→ 25% percentile values are below $3.5 = \frac{3+4}{2}$

→ The values that are less than 3.5 are fall in 25 percentile.

② Quartiles

25% = 1st quartile

50% = 2nd quartile

75% = 3rd quartile

$$Q_3 = 75 \text{ percentile} = \frac{75}{100} * (19+1)$$

= 15th value

$$Q_3 = 7,$$

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

$$\text{lower fence} = Q_1 - 1.5(IQR)$$

$$= 3 - 6 = -3,$$

$$\text{upper fence} = Q_3 + 1.5(IQR)$$

$$= 7 + 1.5(4)$$

$$= 7 + 6$$

$$= 13$$

$$[-3, 13]$$

5 Number Summary and Boxplot

[To find outliers]

(i) ~~minimum~~ minimum

(ii) First quartile (25 percentile) Q₁

(iii) median

(iv) Third quartile (75 percentile) Q₃

(v) Maximum.

Removing the outlier?

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

[Lower fence \rightarrow Higher fence]

→ Any value that is lower than lower fence & any value that is higher than higher fence is determined as outlier.

outlier = 29 in the above data

5 Number Summary

Minimum = 1

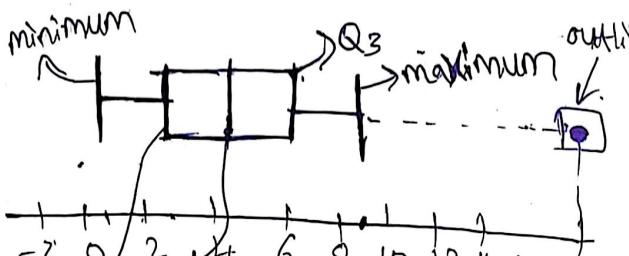
Q₁ = 3

median = 5

Q₃ = 7

maximum = 9

Box Plot



$$IQR = Q_3 - Q_1 = \text{Interquartile Range}$$

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

$$= 5^{\text{th}} \text{ value}$$

Q₁ Median

Q₁ = 3

Histogram 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.

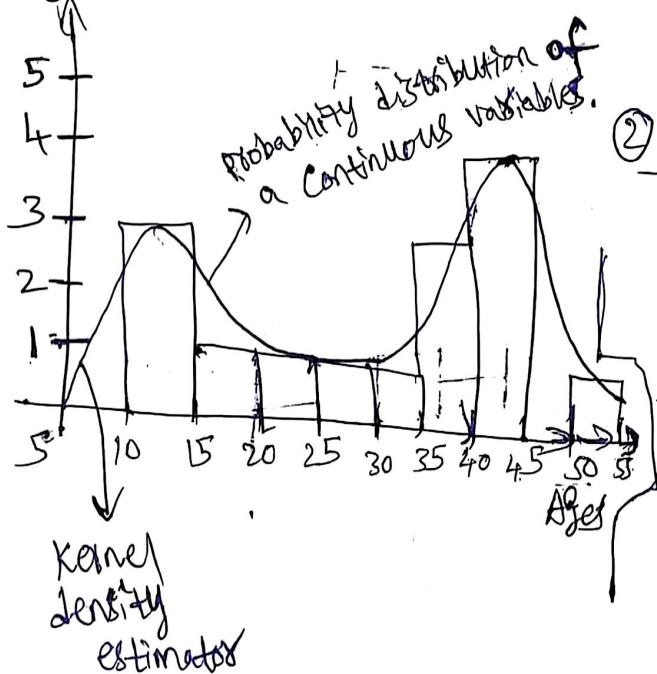
$$\text{Ages} = \{11, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50\}$$

0 → 50 \downarrow
distribution of numerical data

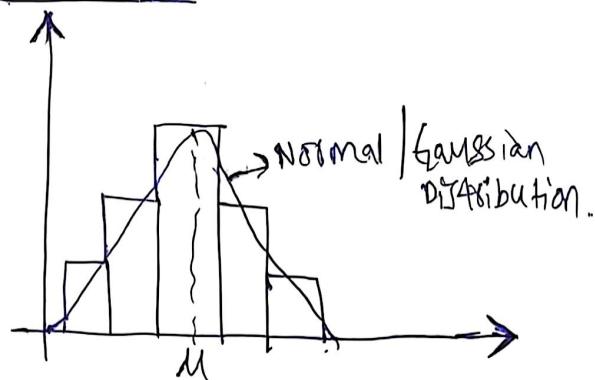
① No. of bins = 10

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

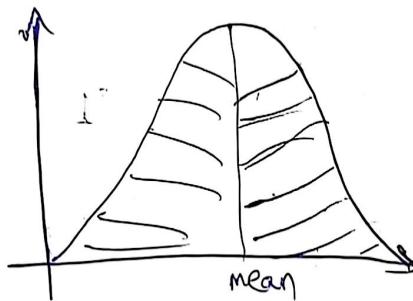
Bins → [0-5, 5-10, 10-15, 15-20,
20-25, 25-30, ..., 45-50]
(frequency count)



Skewness



→ Symmetric distribution \Rightarrow No skewness



Box Plot

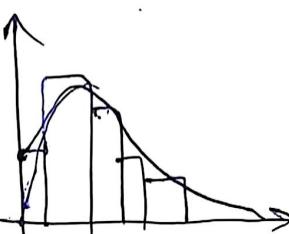


$$Q_3 - Q_2 \approx Q_2 - Q_1$$

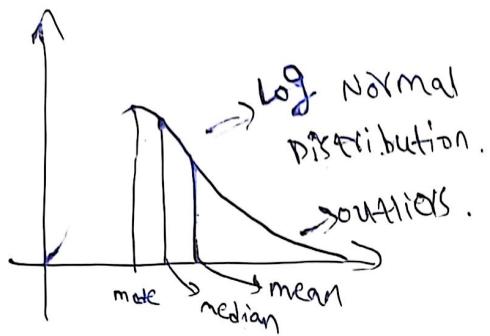
→ The mean, median and mode are all perfectly at the center.

$$\text{Mean} = \text{median} = \text{mode}$$

② Right Skewed



Positive skewed

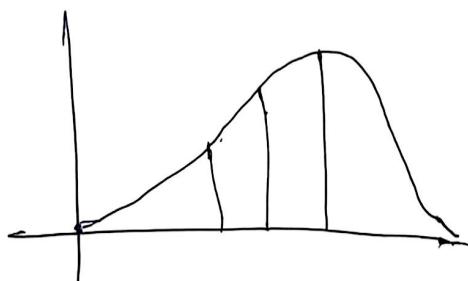


Box Plot :-



→ mean will shift to left.
median $Q_3 - Q_2 \geq Q_2 - Q_1$

③ Left skewed distribution



mean ≤ median ≤ mode

Negative skewed.



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

Covariance + 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

→ it is a measure of how much two random variables change together.

→ If the variables tend to increase or decrease together, the covariance is positive. If one tends to increase when the other decreases, the covariance is negative.

X	Y
2	3
4	5
6	7
8	9

Quantify the relationship b/w X & Y

X↑	Y↑
X↓	Y↑
X↑	Y↓
X↓	Y↓

Dataset

SIZE of House

1200

1300

1500

Price

45 lakhs

50 lakhs

75 lakhs

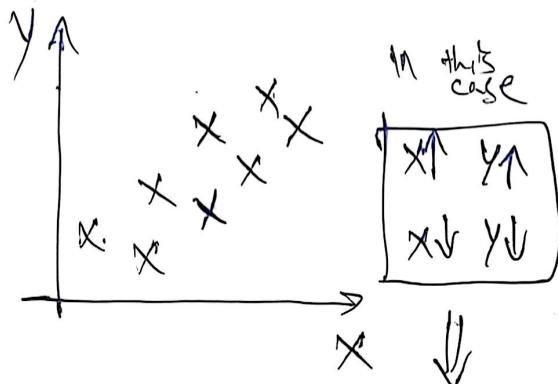
① Covariance

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

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

$$\text{Cov}(X, X) = \text{Var}(X)$$

$$\text{Var}(X) = \frac{(x_i - \bar{x})^2}{n-1}$$



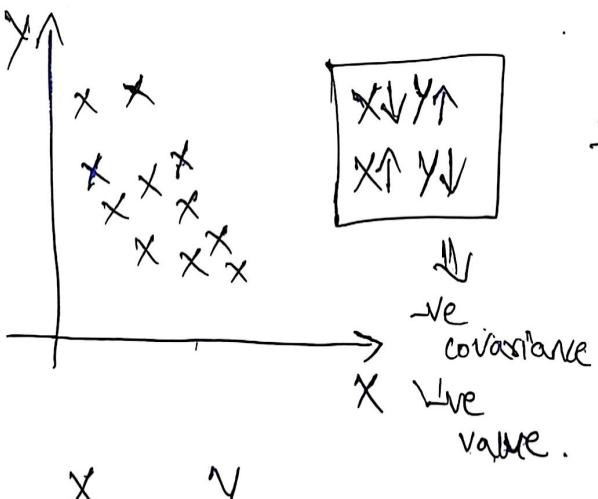
$x_i \rightarrow$ data points of random variable X

$\bar{x} \rightarrow$ sample mean of X

+ve covariance $y_i \rightarrow$ data points of random variable Y
+ve value

$\bar{y} \rightarrow$ sample mean of Y .

~~Ex 2 :-~~



X	Y
7	10
6	12
5	14
4	16

Example :-

Hour studied (X) Exam Score (Y)

2	50
3	60
4	70
5	80
6	90

Here $x↑ y↑$ $x↓ y↓$ \Rightarrow +ve covariance.

we will see with formula.

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

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

$$\textcircled{2} \bar{y} = \frac{50+60+70+80+90}{5} = 70$$

$$\text{Cov}(X,Y) = \frac{(2-4)(50-70) + (3-4)(60-70) + (6-4)(90-70)}{5-1}$$

④ The more the value towards +1, the more the correlated X & Y is.

$\text{Cov}(X,Y) \approx 20$, \rightarrow +ve value
+ve covariance.

④ The more the value towards -1, the more negative correlated X & Y is.

→ The positive covariance indicates the no. of hours studied increased the exam score also.

Advantages:

→ Quantify the relationship between X & Y.

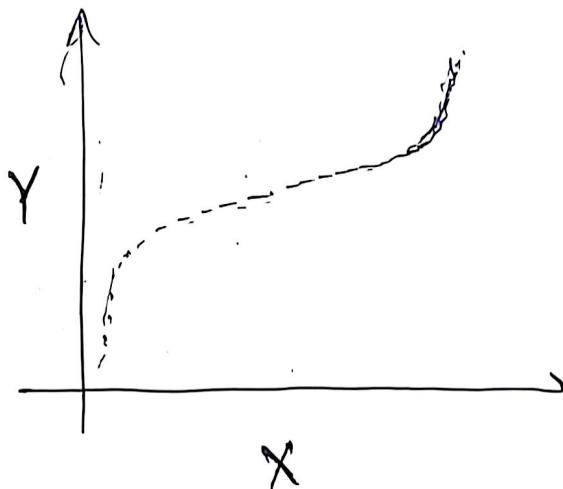
Disadvantage:

→ Covariance does not have a specific limit value.

(ii) Spearman Rank Correlation,

→ This is use for

non-linear data.



Spearman correlation = 1

Pearson correlation = 0.88

→ Here, a Spearman correlation of 1 results when the two variables being compared are monotonically related, even if their relationship is not linear.

② Correlation → Pearson Correlation Coefficient
; ; Spearman Rank Correlation.

Pearson Correlation Coefficient:

→ Value in b/w [-1, 1]

$$r_{xy} = \frac{\text{Cov}(X,Y)}{\sigma_x \sigma_y}$$

$$r_s = \frac{\text{Cov}(R(x), R(y))}{\sigma(R(x)) * \sigma(R(y))}$$

Distributions:

Probability Distribution functions

→ It describes how the probabilities are distributed over the values of a random variable.

Probability Mass function (PMF)

→ Used for discrete random variable.

Probability Density function (PDF)

→ Used for continuous random variable.

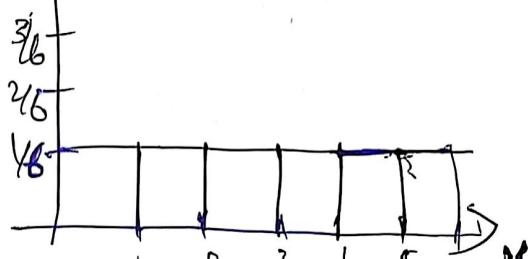
PMF

Rolling a dice.

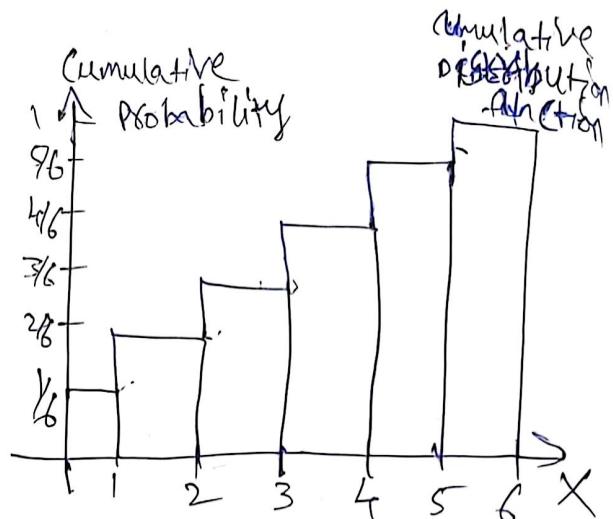
$$\{1, 2, 3, 4, 5, 6\}$$

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

PMF - Probability



Converting PMF to CDF



Cumulative Distribution function (cdf)

$$Pr(X \leq 2) = Pr(X=1) + Pr(X=2) \\ \approx \frac{1}{6} + \frac{1}{6} = \frac{2}{6} = \frac{1}{3}$$

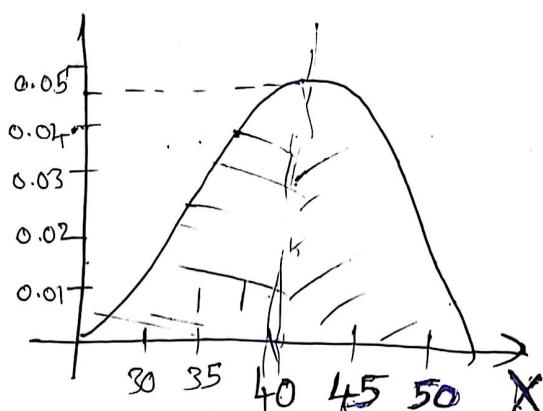
$$Pr(X \leq 6) = Pr(X=1) + Pr(X=2) + \dots$$

$$Pr(X \leq 6) = 1 - Pr(X=6)$$

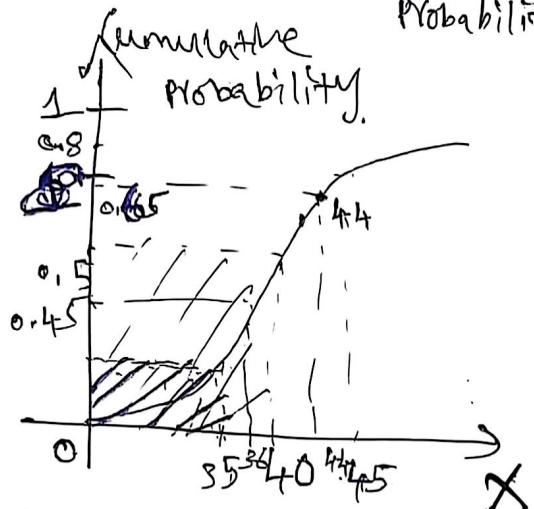
PDF

Ex:-

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



Converting PDF to CDF
Cumulative Probability



$$P(X \leq 40) = ?$$

→ Area under the Curve

$$P(X \leq 35) = ?$$

→ For getting probability density from cumulative density function we use slope.

$$\text{slope} = \frac{0.65 - 0.45}{44 - 38} = 0.05$$

Probability density.

$$\text{slope} = \frac{x_2 - x_1}{y_2 - y_1} \Rightarrow \text{Gradient}$$

* * * → Probability density is the gradient of cumulative density function.

1. Normal Distribution

(Gaussian distribution)

→ Symmetrical, bell-shaped

curve characterised by its mean (center) μ standard deviation (spread)

→ Empirical rule (68-95-99.7)

→ If a normal distribution has a mean of 0 e.g. a standard deviation 1, it is referred to as a standard normal distribution.

→ Z-score can be used to standardize values from any normal distribution to the standard normal distribution.



② Bernoulli Distribution /

→ It is a simple & fundamental discrete probability distribution that models a random experiment with two possible outcomes

"success" $\rightarrow 1$

"failure" $\rightarrow 0$

Ex! flipping a coin
(heads or tails)

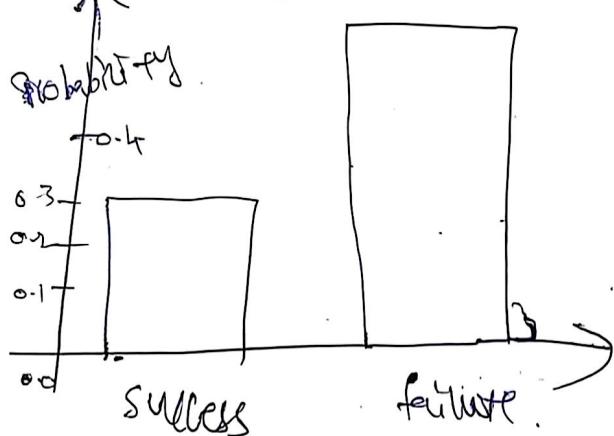
Ent of a single question

Yes / no

mean (expected value) $E[X] = p$

Variance $\text{Var}[X] = p(1-p)$

Bernoulli distribution.

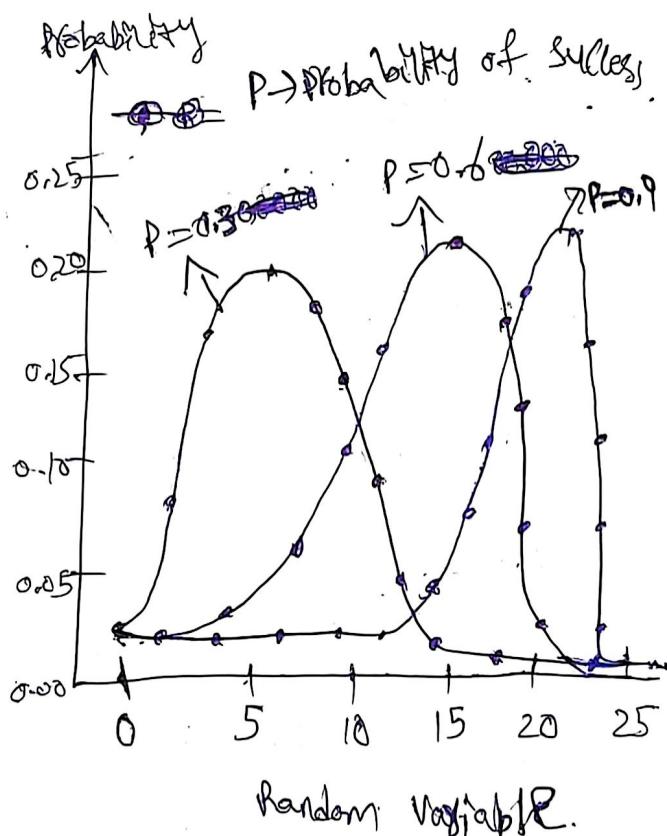


Prob of getting success

$$= 0.3$$

$E[X]$ in failure = 0.7

- Ex-
→ no. of heads in a fixed no. of coin flips.
→ defective/non-defective items in a production process.



Binomial Distribution

→ It is a discrete probability distribution.

→ It describes the no. of successes in a fixed no. of independent and identical Bernoulli trials.

mean $\rightarrow E[X] = np$

$\text{Var}[X] = np(1-p)$

Poisson Distribution

→ It is a discrete probability distribution.

→ It models the no. of events occurring in a fixed interval of time or space.

→ It is particularly useful for situations where events are rare & random.

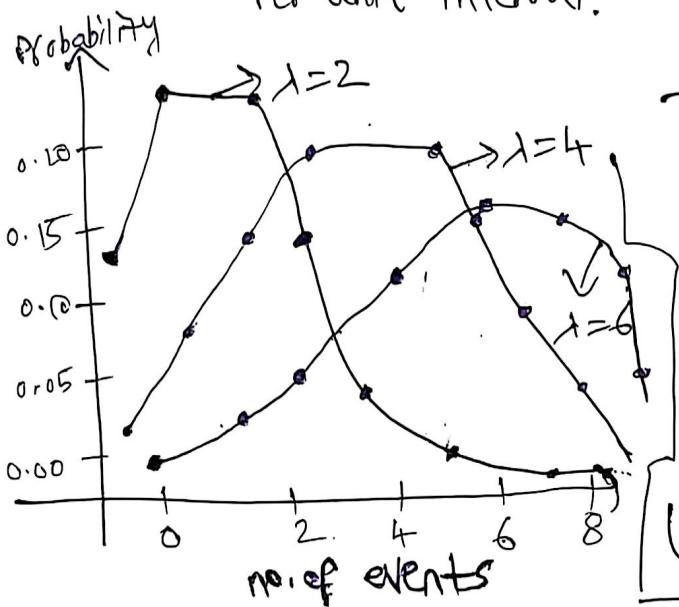
Ex:- no. of phone calls at a call center in a given minute.

→ No. of arrivals at a service point.

$$E[X] = \lambda$$

$$\text{Var}[X] = \lambda$$

$\lambda \rightarrow$ Avg rate of events per unit interval.



Poisson distribution varying λ .

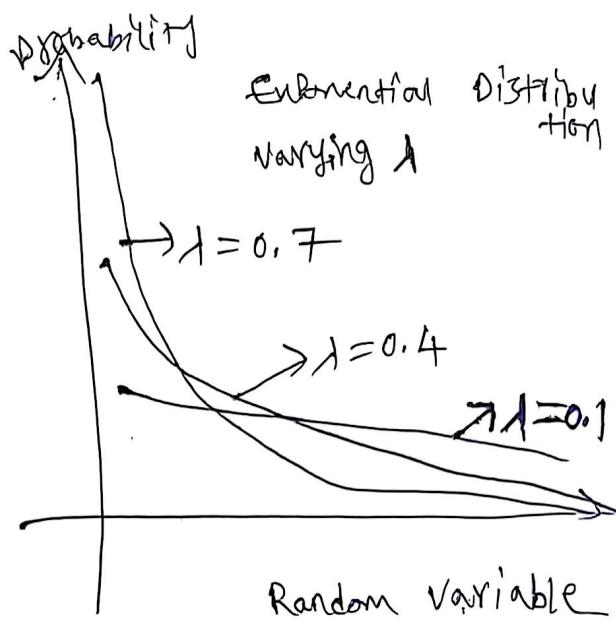
Exponential Distribution

→ It is a continuous probability distribution that models the time b/w events in a Poisson process.

Ex:- time b/w arrivals

Ex:- at a service point.
→ Time until a system failure.
mean $E[T] = \lambda$

$$\text{Var}[T] = \frac{1}{\lambda^2}$$

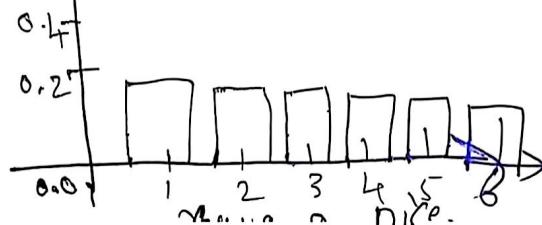


Uniform Distribution

→ All values in distribution have an equal probability of occurring.

→ There are continuous & discrete versions of the uniform distribution.

Ex:- generating a random number.



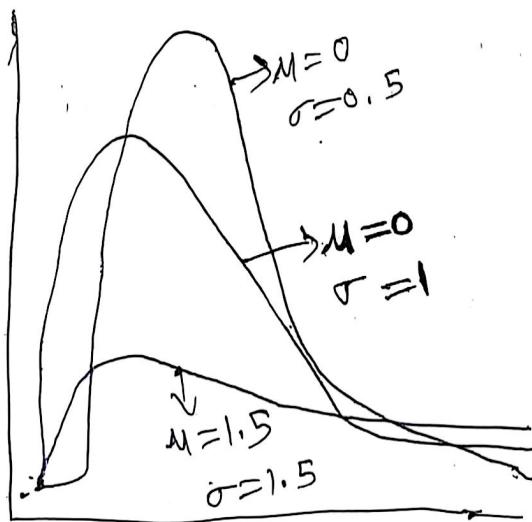
Log-normal Distribution | Power-Law Distribution

→ It is a continuous probability distribution of a random variable whose logarithm is normally distributed.

→ If you takes a logarithm for a Gaussian variable

~~stays~~

→ By taking gaussian variable and exponentiate it, you get the lognormal distribution.



Ex! - The length of ~~Comments~~ posted in Internet discussion forums follows a log normal distribution.

→ It describes a functional relationship between two variables, where one variable is proportional to a power of the other,

$$f(x) = Kx^\alpha$$

$x, y \rightarrow$ variables of interest

$\alpha \rightarrow$ law exponent

$K \rightarrow$ constant

Ex! If side of a square increased by two units, the area increases by four units.

Pareto Distribution

→ A specific type of distribution that follows power law is called Pareto distribution.

→ Pareto principle states that 80% of the

effects come from 20% of the cause.

→ for Ex) 80% of the world's wealth is earned by 20% of the people.

→ The idea sometimes expressed as "80-20" rule".

→ The CLT says that the sampling distribution of 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 mean will be normal.

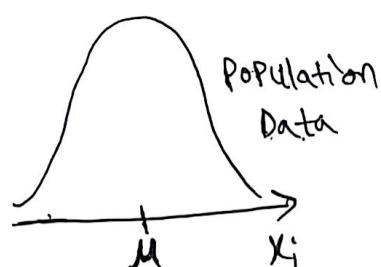
Note!

Need to study formulae (derivations) from Krish Naik course.

Central Limit Theorem (CLT)

→ The Central Limit Theorem relies on the concept of a sampling distribution, which is the probability distribution of a statistic for a large no. of samples taken from a Population.

$$① X \sim N(\mu, \sigma)$$



n = sample size \Rightarrow any value

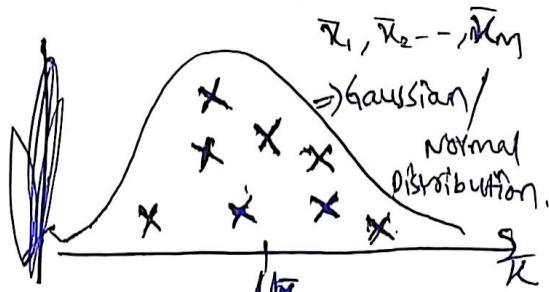
$$S_1 = \{x_1, x_2, \dots, x_n\} = \bar{x}_1$$

$$S_2 = \{x_2, x_3, \dots, x_n\} = \bar{x}_2$$

$$\vdots \qquad \qquad \qquad = \bar{x}_3$$

$$S_m \qquad \qquad \qquad = \bar{x}_m$$

Sampling distribution of means



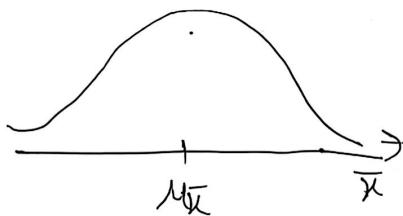
② $X \not\sim N(\mu_0)$ \rightarrow not a normal distribution.
 \downarrow random variable

here sample size
must be $n \geq 30$

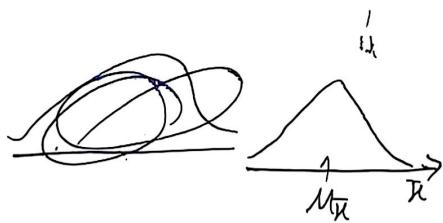
$$S_1 = \{x_1, \dots, x_{30}\} = \bar{x}_1$$

$$S_2 = \{x_1, \dots, x_{30}\} = \bar{x}_2$$

$$\begin{array}{ll} S_3 & = \dots = \bar{x}_3 \\ \vdots & \\ S_m & = \dots = \bar{x}_m \end{array}$$



Sampling distribution of mean



σ = population std deviation

μ = population mean

n = sample size.

$$X \sim N(\mu, \frac{\sigma^2}{\sqrt{n}})$$

AFTER DISTRIBUTION

~~for~~ for large sample sizes, the sampling distribution of means will approximate to normal distribution even if the population distribution is not normal.

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.

A Hypothesis and Hypothesis Testing Mechanism

① Null hypothesis (H_0)

↳ It is the assumption we are beginning with

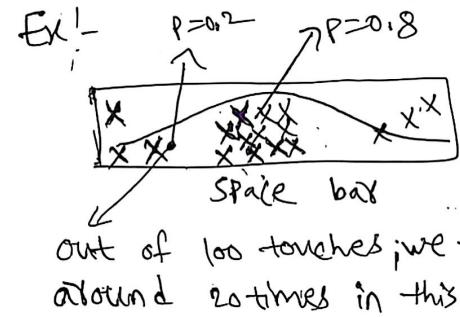
② Alternate hypothesis (H_1)

↳ opposite of null hypothesis.

③ Experiments → statistical Analysis
 {P-value, significance value}

④ Accept the null hypothesis

or Reject " " "



Hypothesis Testing

Eg) coin is fair or 100 times

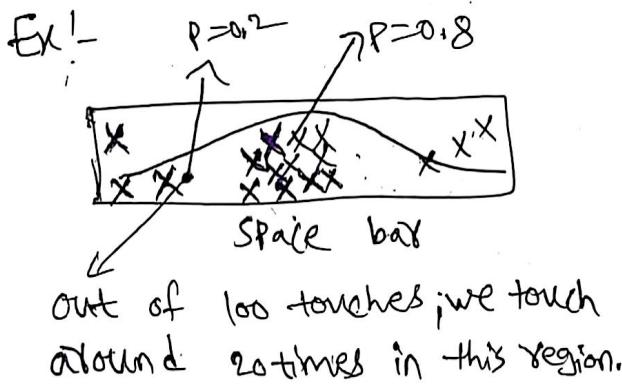
① Null hypothesis:

$H_0 \rightarrow$ coin is fair

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 times?}

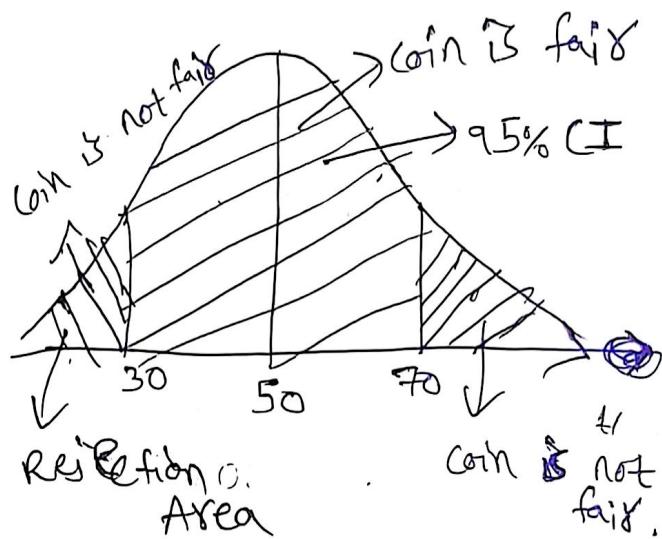
① Null Hypothesis:

$H_0 \rightarrow$ coin is fair

② Alternate hypothesis

$H_1 \rightarrow$ coin is not fair

③ Experiment: {100 times?}



④ Significance value

$$\alpha = 0.05$$

$$C.I = 1 - 0.05 = 0.95$$

⑤ Conclusion

$p <$ significance value

Reject the null hypothesis.

else

Fail to reject the null hypothesis.

Hypothesis Testing and

Statistical Analysis:

- ① Z Test
 - ② t Test
 - ③ CHI SQUARE \Rightarrow Categorical Data
 - ④ ANNOVA \Rightarrow Variance.
- Average (mean) \Rightarrow
 z-table \Rightarrow p-value \Rightarrow t-table

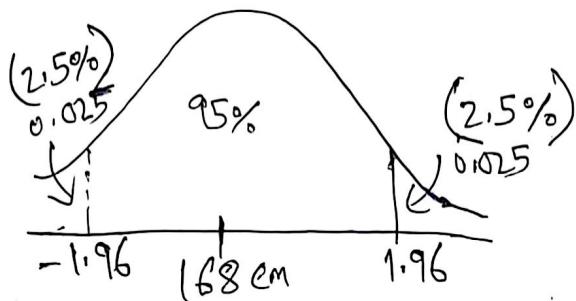
① Null Hypothesis

$$H_0 = \mu = 168\text{cm}$$

② Alternate Hypothesis

$$H_1 = \mu \neq 168\text{cm}$$

③ Based on C.I. we will draw decision boundary.



We use Z-test when you know
(i) population standard deviation
(ii) $n \geq 30$

① Z-Test

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

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

Ans)

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

$$n = 36 \quad \bar{x} = 169.5$$

$$C.I = 0.95 \quad \alpha = 1 - 0.95$$

* This is 2 tail test, because we are trying to find, the avg height is greater than or less than 168cm. We are checking in both the directions, so it is 2-tail.

Area under curve
 $1 - 0.025 = 0.9750$

Check Z-table &
Compute Z-score.
For 0.9750 the
Z-score is 1.96

If Z is less than -1.96 or greater than +1.96, Reject the null hypothesis.

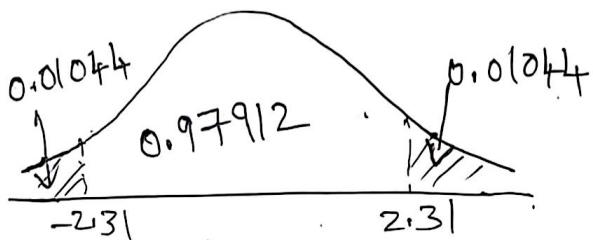
Z-test

$$Z_d = \frac{\bar{x} - \mu}{\frac{3.9}{\sqrt{n}}}$$
$$= \frac{169.5 - 168}{\frac{3.9}{\sqrt{36}}}$$
$$\boxed{Z_d = 2.31}$$

Conclusion :-

$2.31 > 1.96$, Reject null hypothesis

Based on calculation.



from Z-table

The area corresponds to ' 2.31 ' is 0.98956

Getting the tail Area

$$1 - 0.98956 \\ = 0.01044$$

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

$$P < 0.05$$

$$0.02088 < 0.05$$

\Rightarrow Reject the null hypothesis.

Final Conclusion :-

Avg Height $\neq 168\text{cm's}$

\rightarrow The avg height seems to be increasing & based on sample heights.

② A factory manufactures bulbs with a avg warranty of 5 years with standard deviation of 0.50. A worker believes that the bulb will malfunction is less than 5 years. He tests a sample of 40 bulbs and find the average time to be 4.8 years

(a) State null and Alternative hypothesis

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

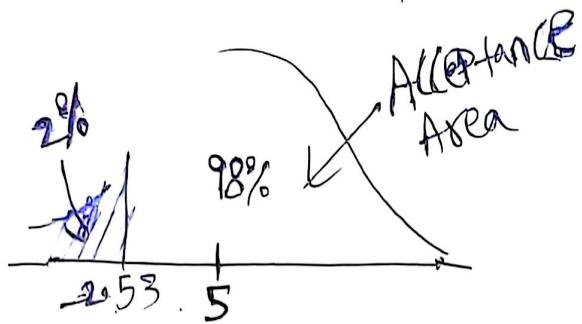
Ans:- $\mu = 5$ $\sigma = 0.5$ $n = 40$
 $\bar{x} = 4.8$

③ Null hypothesis $H_0 = \mu = 5$

Alternative " $H_1 = \mu < 5$

{ 1 tail Test }

(b) Decision Boundary



Z-test

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

From Z-table

Area under curve with Z-score -2.53 is 0.0570

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

Compare P-value with α

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

we accept the null hypothesis

we fail to reject the

null hypothesis.

Student t distribution

~~→~~

→ 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 don't know the population standard deviation?



for this scenario we use Student's t distribution.

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

s = sample standard deviation,

T-stats & T-test

↳ 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.

Ans:-

$$\mu = 100 \quad n = 30$$

$$\bar{x} = 140 \quad s = 20$$

$$CI = 95\%$$

$$\alpha = 0.05$$

i) Null Hypothesis $H_0: \mu = 100$

Alternate ii) $H_1: \mu \neq 100$

{Two tail test}

① Calculate Test statistic /

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

$$t = 10.96$$

$$10.96 > 2.045$$

Conclusion Reject null hypothesis

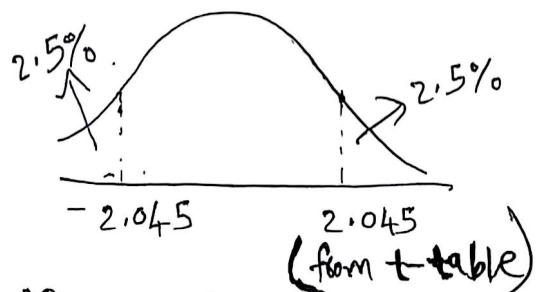
* Medication used has affected the intelligence
medication has increased the intelligence

iii) $\alpha = 0.05$

iv) Degree of freedom

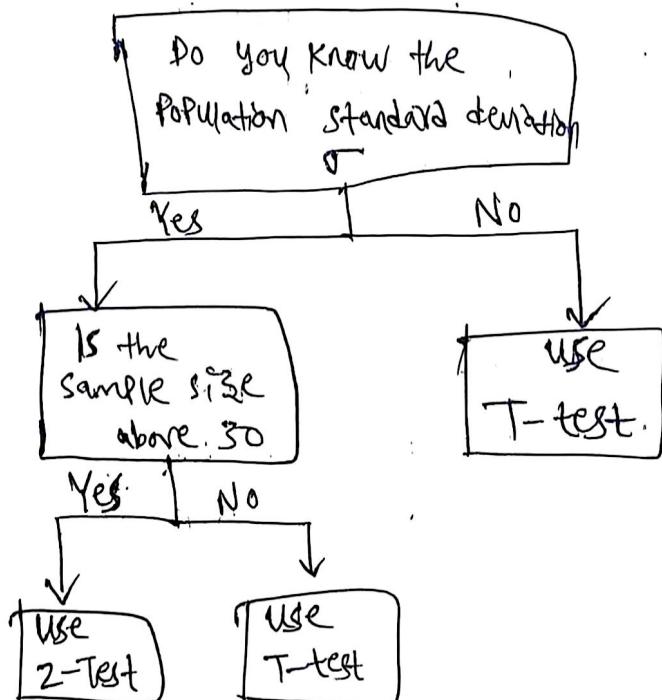
$$df = n - 1 = 30 - 1 = 29$$

v) Decision Rule)



→ if t-test is less than -2.045 or greater than 2.045 , reject the null hypothesis.

When to use T-tests vs Z-tests



Type 1 and Type 2 Errors) Confidence Interval & Margin of Error

Ques

Reality : Null Hypothesis is True

Point Estimate \pm Margin of Error

or " " " false If it is Z-test

Decision : Null Hypothesis is True

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

or " " " false

Outcome 1 : we reject the null hypothesis when in reality it is false \rightarrow good.

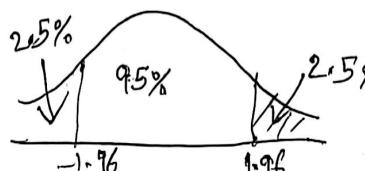
Ex: On the Verbal section of CAT exam, the standard deviation is known to be 100. A sample of 30 test takers has a mean of 520. Construct 95% CI about the mean.

Outcome 2 :

we reject the null hypothesis when in reality it is true

\hookrightarrow Type 1 Error

$$\alpha = 0.05 \quad \bar{x} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$



from z-table

Outcome 3 :

we retain the null hypothesis, when in reality it is false

\hookrightarrow Type 2 Error.

$$\text{Lower CI} = 520 - (1.96) \frac{100}{\sqrt{25}} = 480.8$$

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

Outcome 4 :-

we retain the null hypothesis when in reality it is TRUE \rightarrow good.

Conclusion :-

I am 95% confident about the mean CAT score is between 480.8 and 559.2.

CHI SQUARE TEST

- The Chi square Test for Goodness of fit test claims about population proportions.
- It is a non-parametric test that is performed on Categorical [Ordinal & Nominal] data.

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

	Theory	Sample
Yellow Bike	$\frac{1}{3}$	22
Red	"	17
Orange	"	59
	↓ Theory categorical distribution	↓ observed categorical distribution

Goodness of fit test

- In a science class of 75 students, 9 are left handed.

Does this class fit the theory that 12% of people are left handed

	O	E
left handed	9	$\rightarrow \frac{12}{100} \times 75$
Right handed	66	

Ex-1

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

<50kg	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-1

2010	<50kg	50-75	>75
Expected	20%	30%	50%

2020	<50	50-75	>75
Observed $n=500$	140	160	200

Expected	<50	50-75	>75
	$0.2 \times 500 = 100$	$0.3 \times 500 = 150$	$0.5 \times 500 = 250$

① Null Hypothesis : H_0

The data meets the expectation.

Alternate Hypothesis : H_1

The data doesn't meet the expectation.

⑩ $\alpha = 0.05$ ($C.I = 95\%$)

② Calculate Chi Square Test Statistic

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

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

$$\chi^2 = 26.66$$

③ Dof

$$3 - 1 = 2$$

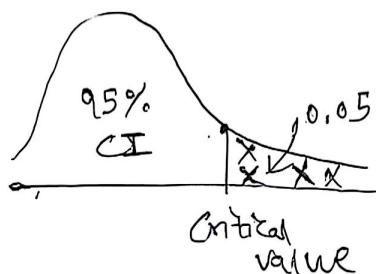
$$\chi^2 > 5.991$$

Reject Null Hypothesis.

④ Decision Boundary

Conclusion :

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



$$dof = 2$$

$$\text{significance} = 0.05$$

see CHI SQUARE Table for value \rightarrow ANOVA is a statistical method used to compare the means of three or more groups.

If χ^2 is greater than 5.99,

Reject H_0

else
we fail to reject the null hypothesis

Analysis of Variance (ANOVA)

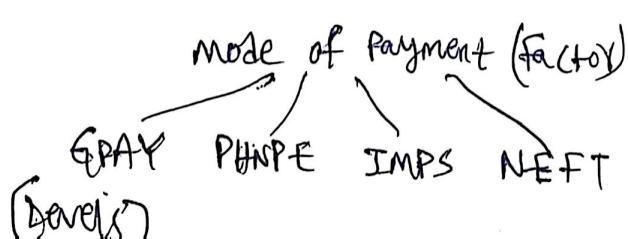
ANOVA

① Factors (variables)

② Levels

e.g. medicine (Factor)

[Dosage] 5mg 10mg 15mg \rightarrow levels



Assumptions in ANOVA

① Normality of Sampling

Distribution of Mean.

→ the distribution of sample mean is normally distributed.

② Absence of outliers.

→ outliers need to be removed from the dataset.

③ Homogeneity of Variance

→ population variance in different levels of each independent variable are equal.

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

④ Samples are independent & Random.

Mathematics - And statistics for data science

by Code Basic

Complete Statistics for Data Science By Krish

① Descriptive statistics

- (i) Measure of Central Tendency
(ii) " " Dispersion

summarizing the data.

- Histograms, Pdf, Cdf, Probability
- Permutation, Mean, median, mode,
Variance, Standard deviation.

→ Gaussian distribution

→ LogNormal distribution

→ Binomial "

→ Bernoulli's "

→ Poisson distribution

→ Standard Normal distribution

→ Transformation and standardization

→ Q-Q plot

Naik |

② Inferential statistics

→ Z-test

→ t-test

→ ANOVA (also called F-test)

→ CHISQUARE

→ HYPOTHESIS TESTING {P value
Confidence Intervals}

What is Statistics?

→ It is the science of collecting, organizing and analyzing data.

→ Statistics helps in taking better decision making.

What is Data?

→ Facts or pieces of information that can be measured.

Eg: The IQ of a class

{98, 97, 60, 55, 75, 65} → Data

Age of students of a class

{30, 25, 24, 23, 27, 28} → Data

Types of statistics

① DESCRIPTIVE STATS

→ It consists of organizing and summarizing data

② INFERENTIAL STATS

→ Techniques where in we used the data that we have measured to form conclusions.

Eg:-

Classroom of Maths student(20)

Marks of the 1st sem

84, 86, 78, 72, 75, 65, 80, 81,
92, —

Eg:- Descriptive stats

① What is the average marks of the students in the class?

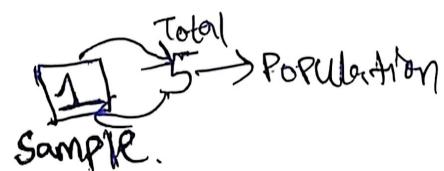
Eg: Inferential stats question
be like

→ Are the marks of the students of this classroom similar to the marks of the maths classroom in the college?

→ If there are 5 classrooms, then I take only 1 classroom to analyse.

→ Population

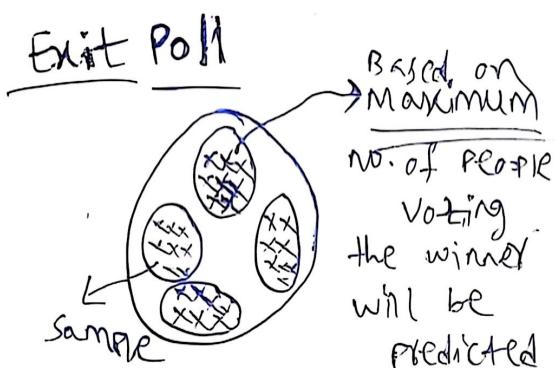
✓ Sample



Population and Sample: ~~Eg based on profession~~

Elections \rightarrow Goa, UP

~~Don't mix up Python~~

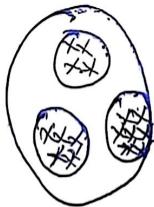


Population (N) Sample (n)

Sampling Techniques

① Simple Random Sampling

\rightarrow When performing simple random sampling, every member of the population (N) has an equal chance of being selected for your sample (n)



② Stratified Sampling

\rightarrow In this the population (N) is split into non-overlapping groups (strata)

Eg: Gender \rightarrow Male
 Female

Eg: Survey

Age groups

Gender (Male/Female) \rightarrow In this there is no change of distribution

③ Systematic Sampling

from population (N) \rightarrow n th individual can be examined

Eg: Mall \rightarrow Survey (Covid)

for every 8th person you do survey

It's not fixed, you can do survey to every 1st person or 3rd person, whatever person you need you can do survey to the same n th person in a systematic way

④ Convenience Sampling

\rightarrow For example we are doing survey, we can take only expertise persons in that domain.

Eg: Data Science

\rightarrow If we are doing survey on data science, we can consider persons who have knowledge on data science only.

Variables

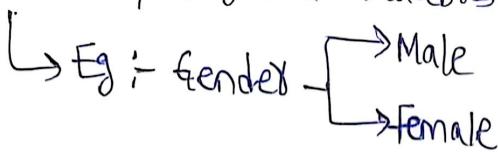
\rightarrow A variable is a property that can take on any value.

Eg: Height = {182, 178, 168, 150, 160, 170}
Weight = {78, 99, 100, 60, 50}

TWO Kinds of Variables

① Quantitative Variable → Measured Numerically, {Add, subtract, multiply, divide}

② Qualitative / Categorical Variables



{Based on some characteristics }
We can derive categorical variables

Eg :- IQ

0 - 10	10 - 50	50 - 100
↓	↓	↓
less IQ	medium IQ	Good IQ

→ It is two types

Quantitative

~~Qualitative~~

Discrete Variables

Eg :- Whole number

① No. of Bank Accounts

2, 3, 4, ...

② Total no. of child in family

2, 3, 4

Continuous Variables

Eg :-

Height = {172.5, 162, 163.7, ...}

Weight = {100, 99.5, 99.75, ...}

Rainfall = {1, 1.25, 1.3, ...}

Examples

④ What kind of variable Gender is? Categorical

⑤ " " " Marital status? Categorical

⑥ " " " River length? Continuous

⑦ " " " Population of state? Discrete

⑧ " " " Song length? Continuous
Blood pressure? Continuous

PIN CODE? {Discrete or Categorical}

X
As classes are going
we will know about it

Variable Measurement Scales | Ratio data

4 types of measured variable

- ① Nominal ② Ordinal
- ③ Interval ④ Ratio

→ This have all the properties of interval variables plus a real absolute zero.

→ Value of zero represents the total absence of the variable.

Ex: 100 feet is twice as long as 50 feet.

Nominal data

→ This is categorical / qualitative
→ This comes from names. data

Eg: Colors, Gender, Type of flower

Eg: Married or not married, Yes or No, blood type (A_B, AB, O) etc.

Ordinal

→ Order of data matters, value does not

→ Ordinal means order. Order means ranking

Eg: Students (marks)	Rank
100	1
96	2
57	4
85	3
44	5

→ Ordinal scales both name and order.

Ex: finish a race (first, second, third...)

Interval data

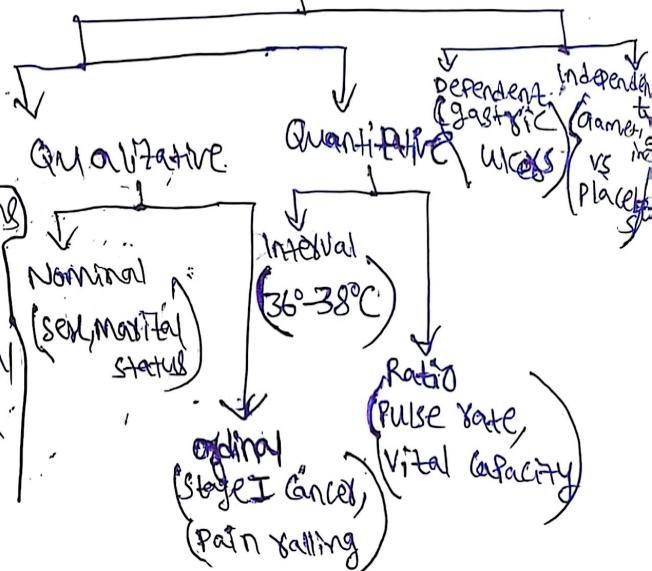
→ In this order matters, values also matters but natural zero is not present

Eg: Temperature

Fahrenheit

-70-80, 80-90, 90-100, 0
 ↑
 ↓
 This is not present in this. That's why natural zero is not present

Types of Variables



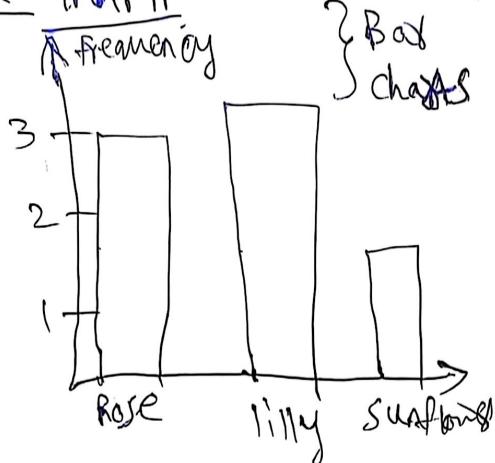
Frequency Distribution

Sample data: Rose, lily, sunflower, Rose, lily, sunflower, Rose, lily, lily

Flower	Frequency	Cumulative Frequency
Rose	3	3
Lily	4	7
Sunflower	2	9

This is natural zero.

① BAR GRAPH



BAR IS HISTOGRAM }

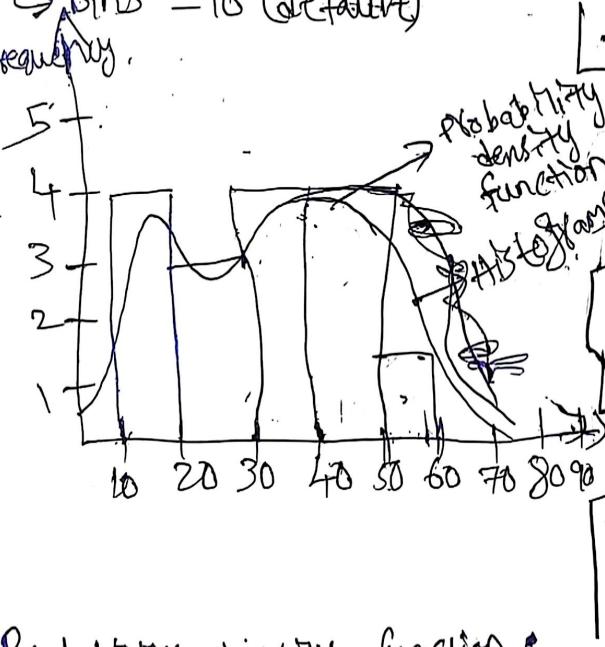
→ Bar is used for discrete values

→ Histogram is used for continuous values.

② Histograms: continuous

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

→ BINS = 10 (Default)
frequency.



Basics to Intermediate stats:

→ Measure of Central Tendency

→ Measure of dispersion

→ Gaussian distribution

→ Z-score

→ Standard normal distribution

③ Arithmetric mean for Population & Sample:

Mean (Average)

Population (N) | Sample (n)

Probability density function:-

→ Smoothening of histograms.

How the smoothening is done?

Ans: Based on {kernel: density}
estimator

we will learn about
it in upcoming
classes.

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

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

$$\bar{M} = \frac{\sum_{i=1}^N X_i}{N}$$

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

$$= \frac{32}{10} = 3.2$$

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

$$= 3.2$$

~~1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100, 100~~

even number = 12

$$\text{median} = \frac{12}{2} = 6$$

= ~~avg~~ of 6th & 7th element

$$= \frac{3+4}{2} = \frac{7}{2} = 3.5$$

$$\text{median} = 3.5$$

Measure of Central Tendency!

What is Central Tendency?

→ Refers to the measure used to determine the centre of the distribution of data.

→ median works well ~~ever~~ if outliers present.

① Mean ② Median ③ Mode

② Median }

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

$$\text{Mean} = \frac{1+1+2+2+3+3+4+4+5+5+6}{11}$$

$$= \frac{32}{11} = 2.9$$

→ outliers can effect the "mean".

→ outliers can make adverse effects on the distribution (mean)

$$\text{median! } \{1, 1, 2, 2, 3, 3, 4, 4, 5, 5, 6\}$$

→ first sort the numbers

$$\text{odd number} = 11$$

$$\text{median} = \frac{11}{2} = 5.5$$

= 6th element
is median

~~even number = 12~~

③ Mode }

$$\{1, 2, 2, 3, 4, 4, 5, 5, 6, 6, 6, 6, 7, 8, 100, 100, 100, 100\}$$

Mode = {most frequent Element}

$$\text{Mode} = 3$$

→ Mode works well with "categorical" where mode is used? ~~total~~

Type of flower | petal length | width | DATASET
flavor | DATASET

Rose

Lily

Sunflower

Missing data

{ missing value → filled with most recent element }

mode works well with categorical variable

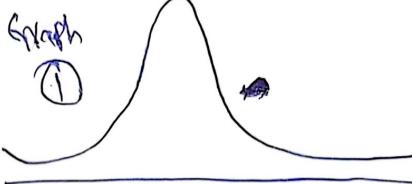
Measure of Dispersion

① Variance

② Standard Deviation

↓
it means how well the data is spread.

$$\sigma^2 = \frac{10.84}{6} = 1.81$$

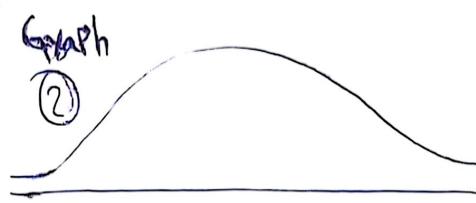


$$\{1, 1, 2, 2, 2\} \rightarrow \mu = \frac{8}{5} = 2$$

$$\{2, 2, 2, 2, 2\} \rightarrow \mu = \frac{10}{5} = 2$$

How these two distributions is different?

→ So, to say difference we use variance and standard deviation.



→ The variance is more in Graph 2. Because the spread is more in Graph 2 compared to Graph 1

① Variance (σ^2)

→ It is nothing but spread.

Population Variance

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

~~standard deviation~~

Sample Variance s^2 sample mean

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

Ex:-

$$x \quad \mu \quad x-\mu \quad (x-\mu)^2$$

$$1 \quad 2.83 \quad -1.83 \quad 3.34$$

$$2 \quad 2.83 \quad -0.83 \quad 0.6889$$

$$2 \quad 2.83 \quad +0.83 \quad 0.6889$$

$$3 \quad 2.83 \quad 0.17 \quad 0.03$$

$$4 \quad 2.83 \quad 1.17 \quad 1.37$$

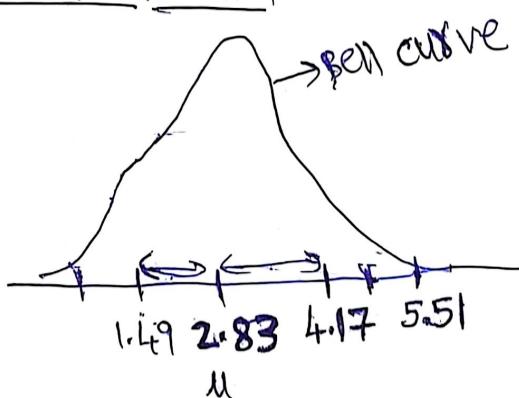
$$5 \quad 2.83 \quad 2.17 \quad 4.71$$

$$\overline{x} = 2.83$$

$$\frac{10.84}{10}$$

$$\sigma = \sqrt{\text{Variance}} = \sqrt{1.81} = 1.345$$

~~standard deviation~~



$$\mu = 2.83$$

$$\sigma = 1.34$$

$2.83 - 1.34 = 1.49$ → first standard deviation towards left
 $2.83 + 1.34 = 4.17$ → first standard deviation towards right

$2.83 - 1.34 = 1.49$ → first standard deviation towards left
 $2.83 + 1.34 = 4.17$ → first standard deviation towards right

$2.83 + 1.34 = 4.17$ → second standard deviation towards right

If we want to find '5'.

→ It falls between 4.17 & 5.51

So, it is 1.5% from the M

How many are below '11'?

$$= \frac{17}{20} \times \frac{5}{100} = 85\%$$

Percentiles and Quartiles

→ These are used to find the outliers.

percentage!

1, 2, 3, 4, 5

% of the numbers that are odd?

$$\% = \frac{\text{# of numbers that are odd}}{\text{Total numbers}}$$

$$= \frac{3}{5} = 0.6 = 60\%$$

Q2 what value exists at percentile ranking of 25%?

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

$$= \frac{25}{100} \times (2) = 5.25$$

← Index position
So take 5th & 6th element and do average.

In the example the 25th percentile value is '5'.

Percentile:

→ A percentile is a value below which a certain percentage of observation lie.

Ex?

Dataset: 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8,

8, 9, 9, 10, 11, 11, 12

What is the percentile ranking of 10?

x = 10

Quartiles

Five Number Summary

- ① Minimum
- ② First quartile (Q_1)
- ③ Median
- ④ Third " (" Q_3)
- ⑤ Maximum

Removing the outliers

{1, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 7, 8, 8, 9, 27}

Lower fence = $Q_1 - 1.5(IQR)$

Percentile rank of $x = \frac{\text{# of values below } x}{n} \times 100$ unless fence = $Q_3 + 1.5(IQR)$

$$= \frac{16}{20} \times \frac{100}{100} = 80\%$$

Q_1 - 25% percentile

Q_3 - 75% percentile

As per example

$$Q_1 = \frac{25}{100} \times (19+1) = .5 \rightarrow \text{index}$$

$$= 3$$

$$Q_1 = 3$$

$$Q_3 = \frac{75}{100} \times (19+1) = 15 \rightarrow \text{index}$$

$$Q_3 = 7$$

Interquartile Range (IQR)

$$\text{IQR} = Q_3 - Q_1$$

$$= 7 - 3$$

$$= 4$$

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

$$= 3 - 1.5(4)$$

$$= -3$$

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

$$= 7 + 1.5(4)$$

$$= 7 + 6 = 13$$

→ so, anything less than lower fence and anything higher than upper fence are outliers

Here in example less than ~~-3~~ & higher than 13 are outliers.

→ After Removing outliers
the data is

$$\text{minimum} = 1$$

$$Q_1 = 3$$

$$\text{Median} = 5$$

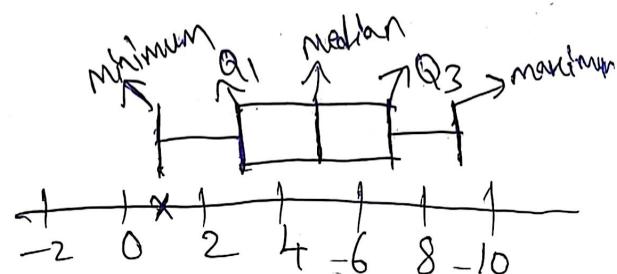
$$Q_3 = 7$$

$$\text{Max} = 9$$

→ 5 number summary

Let us draw Box Plot by using that data.

Box Plot



→ To find outliers we use

Box Plot

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

Bessel's correction
degree of freedom

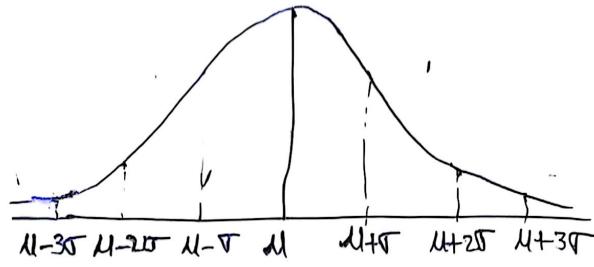
Why sample variance is $n-1$?

And the answer to this question.

Advance Statistics

Distributions

- Normal Distribution
- Standard Normal Distribution
- Z-score Distribution
- Log Normal Distribution
- Bernoulli Distribution
- Binomial Distribution



Empirical Formula:

→ 68% - 95% - 99.7%

Practical Programming with Python

- ① Mean, median, mode
- ② Variance, standard deviation
- ③ Histogram, PDF, Box Plot, Violin plot
- ④ IQR
- ⑤ Log normal distribution.

Example

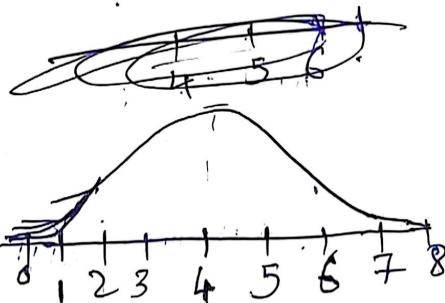
Examples of Gaussian Distribution

Ex. Height → Normal distribution

↳ domain expert will say this.
② weight ③ Iris dataset

↳ Normal distribution

e.g. $\mu = 4 \sigma = 1$



where 4.5 → standard deviation is
we can say easily, it 0.5σ to
the right

But when 4.75 or something like
that comes, we can use Z-score

$$\begin{aligned} \text{Z-score} &= \frac{x_i - \mu}{\sigma} \\ &= \frac{4.75 - 4}{1} = 0.75 \end{aligned}$$

→ Then we can say, it is 0.75σ
to the right

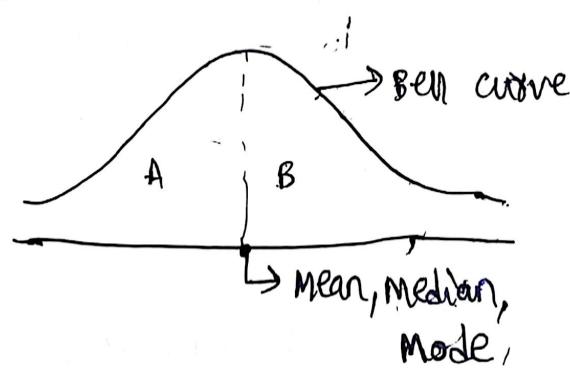
① Distributions

$$\text{Ages} = \{24, 26, 27, 28, 30, \dots\}$$

- Distributions are used to
know the data how it is
distributed.

$$\mu = \dots$$

(i) Gaussian / Normal Distribution



→ Part A < Part B are symmetrical

→ If you want to find how many standard deviation to the right of how many to the left you can use 'Z-score'

Normalization

→ To bring data into (0 to 1 range) you can apply Normalization. For this you can use MinMaxScaler (0 to 1)

Ex:-

$$\mu = 4 \quad \sigma = 1$$

$$\{1, 2, 3, 4, 5, 6, 7\} \rightarrow \text{Normal}$$

Distribution

↓ After applying
Z-score

$$\{-3, -2, -1, 0, 1, 2, 3\}$$

Standard Normal
Distribution

$$(\mu = 0, \sigma = 1)$$

~~+ 100~~
~~- 100~~
~~100~~

Practical example on Z-score

Get your local pizza shop claims their large is atleast 16 inches or free. Their pizza is normally distributed with $\mu = 16.3$ inches & $\sigma = 0.2$ inches.

(i) What's the probability of getting free pizza

(ii) What's the probability of getting lucky with pizza over 16.5 inches?

(iii) What's the probability of getting pizza between 15.95 & 16.63 inches?

→ If the condition is satisfied

$(\mu = 0, \sigma = 1)$, then it is the (i) so,

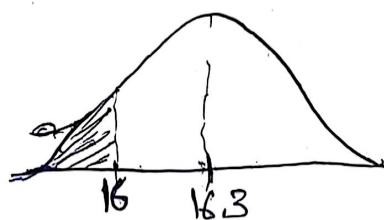
"standard normal distribution"

Practical Application

Dataset

units

Age (years)	Salary (krs)	Weight (kg)
24	40k	70
25	80k	80
26	60k	55
27	70k	45



$$x_i = 16$$

First step is to standardize the distribution

$$z = \frac{x_i - \mu}{\sigma} = \frac{16 - 16.3}{0.2} = -1.5$$

$$P(X < 16) \Rightarrow P(z < -1.5)$$

Look ~~at~~ Left Z-table

$$P(z < -1.5) = 0.0668$$

→ we have 6.68% of chances of getting free pizza.

→ These all are having different units, to bring them into one range you have to apply Z-score ($\mu = 0, \sigma = 1$)

→ This process is called Standardization.

$$(i) P(X > 16.5)$$



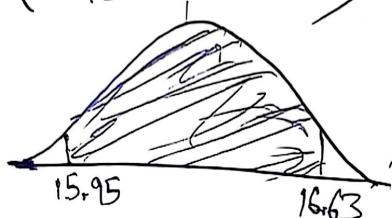
$$Z = \frac{x_i - \mu}{\sigma} = \frac{16.5 - 16.3}{0.2} = 1$$

$$P(X > 16.5) \Rightarrow P(Z > 1) = 0.1587$$

see right z-table

Probability of getting lucky with P329 over 16.5 inches is 15.87%

$$(ii) P(15.95 < X < 16.63)$$



$$Z = \frac{x_i - \mu}{\sigma} = \frac{15.95 - 16.3}{0.2} = -1.75$$

$$= \frac{16.63 - 16.3}{0.2} = 1.65$$

$$P(-1.75 < Z < 1.65)$$

$$= 0.9505 - 0.0401$$

$$= 0.9104.$$

91.04% people will get P329 of length in between 15.95 & 16.63 inches

Day 4 → statistics

- ① IQR - Python
- ② Probability
- ③ Permutation & Combination
- ④ Confidence Intervals
- ⑤ P-value
- ⑥ Hypothesis Testing

① Inter quartile range (IQR)

→ this is used to detect outliers

→ use sns.boxplot to see outliers

→ find q_1 & q_3 quartiles using

$$q_1, q_3 = \text{np.percentile}(\text{dataset}, [25, 75])$$

$$\text{IQR} = q_3 - q_1$$

$$\text{lower fence} = q_1 - (1.5 * \text{IQR})$$

$$\text{higher fence} = q_3 + (1.5 * \text{IQR})$$

② Probability

→ Probability is a measure of the likelihood of an event

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

$$P(\text{getting } 6) = \frac{1}{6}$$

(i) Addition Rule (Probability, or)

Mutually Exclusive Event

If two events are mutually exclusive, they cannot occur at the same time.

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

either 1 or 2 or ... or 6, only one digit can happen at once.

Ans: Non-mutual exclusive

$$P(Q) = \frac{4}{52} \quad P(\heartsuit) = \frac{13}{52}$$

$$P(Q \text{ and } \heartsuit) = \frac{1}{52}$$

Non mutual Exclusive:

multiple events can occur at the same time.

Eg:- Deck of cards $\{Q\}$

↓
It can be heart, diamond, shade etc.

Inclusive Example

Q: If I toss a coin, what is the probability of the coin landing on heads or tails?

Ans: mutual exclusive

Addition Rule

$$P(A \text{ or } B) = P(A) + P(B)$$

$$= \frac{1}{2} + \frac{1}{2}$$

$$= 1$$

$$P(A \text{ or } B) = 1$$

Rolling a dice

$$P(1 \text{ or } 3 \text{ or } 6) = P(1) + P(3) + P(6)$$

$$= \frac{1}{6} + \frac{1}{6} + \frac{1}{6}$$

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

Non-mutual exclusive example

If you are picking a card randomly from deck. What is the probability of choosing a card that is queen **or** a heart?

Addition rule for non-mutual exclusive events.

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

$$P(Q \text{ or } \heartsuit) = P(Q) + P(\heartsuit) - P(Q \text{ and } \heartsuit)$$

$$= \frac{4}{52} + \frac{13}{52} - \frac{1}{52}$$

$$= \frac{16}{52} = \frac{4}{13}$$

Note:

→ Addition rule for mutual exclusive events

$$\boxed{P(A \text{ or } B) = P(A) + P(B)}$$

→ Addition rule for non-mutual exclusive events

$$\boxed{P(A \text{ or } B) = P(A) + P(B) - P(A \cap B)}$$

(ii) Multiplication Rule:

Independent Events

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

one event cannot depend on another. Each & every event are independent.

When you rolled a dice, coming of 1 or anything, cannot depend on another

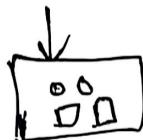
Dependent Event:

Ex:-



$$P(\text{circle}) = \frac{3}{5}$$

If we take out one circle



$$P(\text{square}) = \frac{2}{4}$$

→ Naive Bayes algorithm is based on Dependent Event.

Example of Independent Events:

Q: What is the probability of rolling a "5" and then a "4" in dice?

Ans: Independent Event

Multiplication Rule

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

$$P(4 \text{ and } 5) = \frac{1}{6} * \frac{1}{6} = \frac{1}{36}$$

Dependent Event Example

What is the probability of drawing a Queen & then a Ace from a deck of cards?

Ans: Dependent

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

Conditional Probability

$$P(Q \text{ and Aces}) = P(Q) * P(\text{Aces}|Q)$$

$$= \frac{4}{52} * \frac{4}{51}$$

③ Permutation & Combination:

Permutation:

Ex:-

School trip (Chocolate factory)

Dairy, 5 star, Milky bar, Eclairs, Fren, Kit Kat.

Student {Assignments}

When he visit how many 3 types of chocolates he can see.

$$\underline{6} \times \underline{5} \times \underline{4} = 120$$

formula

$$nP_8 = \frac{n!}{(n-8)!} = \frac{6!}{(6-3)!}$$

$$= \frac{6 \times 5 \times 4 \times 3!}{3!} \\ = 120.$$

Combinations

$$nC_r = \frac{n!}{r!(n-r)!}$$

$$6C_3 = \frac{6!}{3!(6-3)!} = \frac{6 \times 5 \times 4 \times 3!}{3! \times 3 \times 2} \\ = 20.$$

20 unique combinations