

# Statistics

## Data Types

Categorical

Normal

This is without  
order

e.g.

{ This is place,  
play, jump }

There is  
not  
order form  
in

Numerical

Discrete

{ 0, 1, 2, ... }

night  
193.2cm

ordinal

This is with

The format  
of order  
e.g. { child,  
young,  
old }

without  
decimal  
value

with  
decimal  
value

DATE 1/1

Statistics → Statistics is the discipline that concerns the collection, organization, analysis, interpretation and presentation of data

↓  
Descriptive

① Analyze data, summarizing data

organizing data that is the form of Number & Graph

② Bar plot, Histogram, pie chart  
pdf, cdf, normal distribution

③ Measure of Central tendency  
(Mean, median, mode)

④ Measure of variance  
standard deviation

↓  
Inferential Statistics

(Confidence interval)

{ - Z-test  
- T-test  
- Hypothesis test  
- Chi-square test  
Actually

taking the  
some sample  
doing some

inference  
Testing,  
(Take make  
a conclusion)

e.g. exit poll,

## Inferential Statistics →

Based on sample data of population

we will do some conclusion or decision  
that is called inferential statistics  
where there we made decision for this.

e.g. → Exit poll etc

some time called this to confidence  
interval

## Random Variable → (RV) →

every feature present in dataset is a  
Random variable

Random

Numerical Random

Discrete RV

Categorical Random

Counting

variable

R-variable

$$n = 24$$

↓

RV → RV is storing 24 (value)

RV is any type of data store, sometimes is called  
place holder.

Categorical RV → set is repeating  
The records in many times

discrete RV → gets ~~not~~ whole number

e.g. → family member, no. of people

continuous → loan amount, salary

Binarily decimal float value.

Age	Gender	
f <sub>1</sub>	f <sub>2</sub>	f <sub>3</sub>
23	M	
24	F	
25	M	

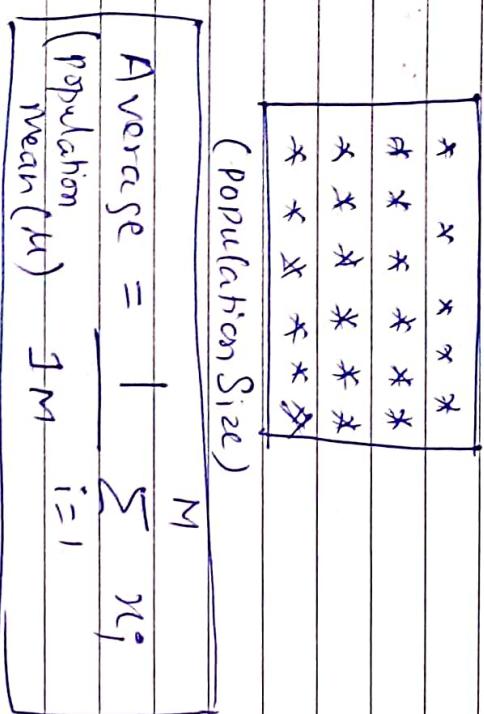
# Population $\rightarrow$ Sample Understand

## Sample

1 - Population (denoted by  $N$ )  
 2 - Sample (denoted by  $n$ )

1 - Population mean ( $\mu$ ) ( $\mu_N$ )  
 2 - Sample mean ( $\bar{x}$ ) ( $x_{\text{bar}}$ )  
 We calculate the Height of population  $\rightarrow$   
Population (Total No. of people denoted by  $N$ )

1 million



$$\text{Average} = \frac{\sum_{i=1}^{10000} (n_i \times \frac{1}{10000})}{n}$$

We can some sample of data from population in different different place  
 e.g exit poll (we can find the winner person is going to win)

## In ML Terms

All the records is the sample based on their we have the prediction  $\rightarrow$

Population Count =  $N$

Sample Count =  $n$

Note:- But we can not check each & every persons of height because it's very difficult to find it.

## Normal Gaussian Distribution

$$Y \sim \text{CD}(\mu, \sigma^2)$$

belong

to

mean

Standard deviation

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

Random variable

of  $n$ 

Gaussian or Normal

Distribution

Probability

When  $\rightarrow$  Data is Normal  
Centre point is median

Center point at

Distribution

50%

Bell curve

Mean, Median, Mode

is equal to 0

50% (Probability)

Because

why saying  $\mu \approx \text{CD}/\text{ND}$ 

and standard

deviation is 1

Belongs to

$$\Omega = P[Y \in [\mu - 2\sigma, \mu + 2\sigma]] \approx 95\%$$

$$\mu - 2\sigma, \mu, \mu + 2\sigma$$

Range of

First Standard

deviation

Second Range of Standard deviation

Third range of Standard deviation

If  $X$  (Random Variable) follow the Normal distribution, Then there are three properties → These properties is called Empirical formula

$$3 - \{ \mu - 3\sigma \leq x \leq \mu + 3\sigma \} \approx 99.7\%$$

Mean / Median / Mode

Lies to Third Range of Standard deviation

below

$$n = [1, 2, 3, 7, 6]$$

$$\bar{n} \text{ or } \mu = \frac{1+2+3+7+6}{5}$$

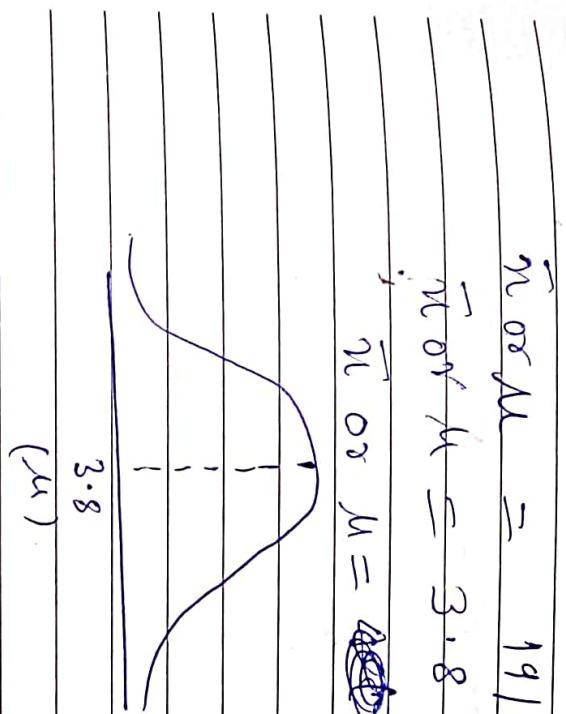
In MC terms

When we are doing EDA

to plot by some of the feature if it is having follow the same type of the curve

Internally follow of Empirical formula of AI

three properties.



Mean is ~~resistant~~ measure of centre tendency.

Mean / median / mode is help to Add some missing value.

$$\text{Mode} = \frac{1}{n} \sum_{i=1}^n n_i$$

Mode = No. of times the particular value occur

## Range

The Range is the difference between the lowest and highest values.

Example

$$4, 5, 6, 9, 3, 7$$

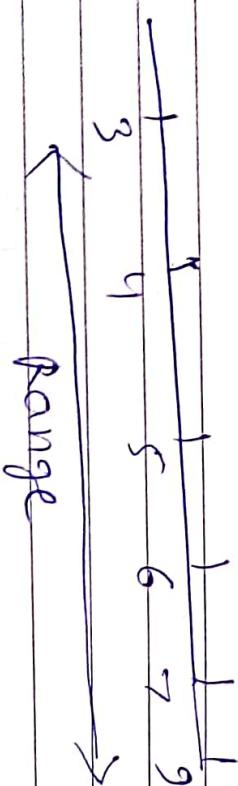
lowest value is = 3

Highest value is = 9

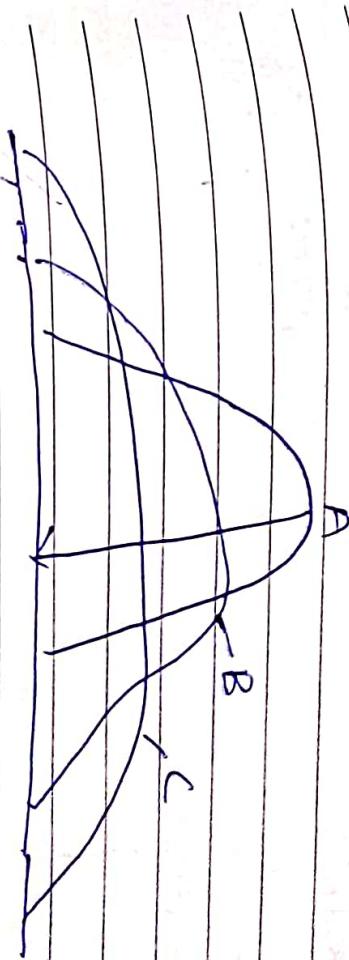
$$\text{Range}(x) = \text{Max}(x) - \text{Min}(x)$$

$$\text{Range}(x) = 9 - 3.$$

$$\text{Range}(x) = 6$$



Dispersion helps to understand the distribution of data.



Measure of dispersion help to interpret the variability of data.

In measurement, it describes

the variation of "observed variables"

↓

To find how much homogenous or heterogeneous the data

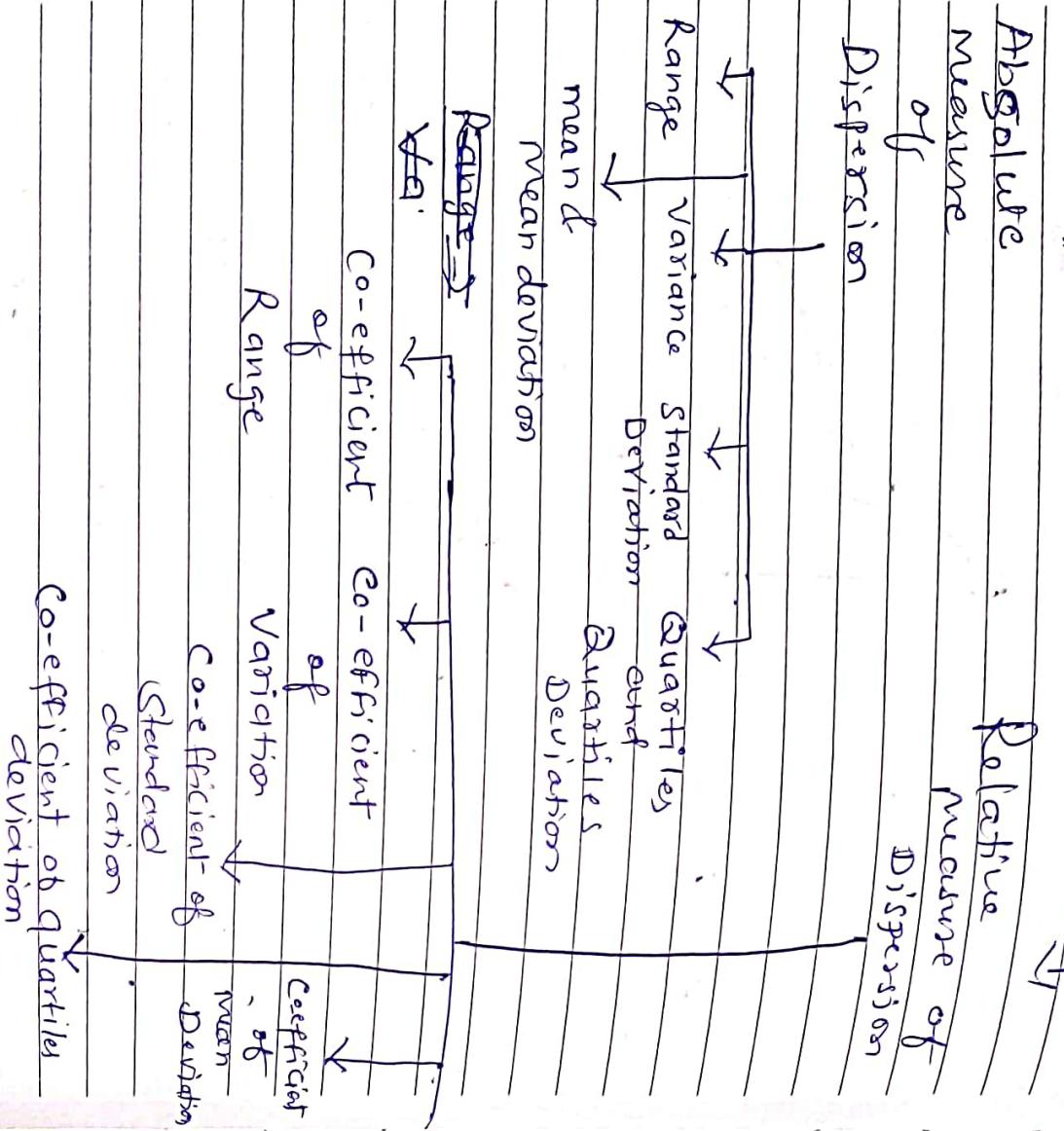
↓  
a simple term

It shows how squeezed the variable is.

## Measures of Dispersion

## Type of measure of dispersion

DATE — 1 / 1



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

where  $\mu = \text{mean}$ ,  $N = \text{total number of population}$

Variance =  $\sigma^2 = (\text{standard deviation})^2$

$$\text{Standard Deviation} = \sqrt{\sigma}$$

$$\text{Standard deviation} = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}}$$

**When  $n = \text{no. of observations in Sample set}$**

$\bar{x} = \mu = \text{mean}$

$n_i = \text{no. of records of the each data}$

### Adding of $(x_i - \bar{m})^2$

e.g. → Find the variance of the  
Number

$$3, 8, 6, 10, 12, 9, 11, 10, 12, 7$$

Sol →

$$\text{Mean} = \frac{3+8+6+10+12+9+11+10+12+7}{10}$$

$$= 73.6$$

$$\sigma^2 = \frac{\sum (x_i - \bar{m})^2}{N}$$

$$\mu = 8.8$$

$$x_i - \mu = 3 - 8.8 = -5.8$$

$$8 - 8.8 = -0.8$$

$$\sigma^2 = \frac{1}{10} (\text{Total No. of Sample})$$

$$6^2 = 7.36$$

$$7 - 8.8 = -1.8$$

$$(x_i - \mu)^2 = (-5.8)^2 = 33.64$$

$$= (-0.8)^2 = 0.64$$

$$= (-2.8)^2 = 7.84$$

$$(-1.8)^2 = 3.24$$

## Quartile Formula

$$\text{Range of Co-efficient} = \frac{x_{\max} - x_{\min}}{x_{\max} + x_{\min}}$$

$$\text{A quartile divide the set of observations into 4 equal parts.}$$

~~25% 75% 25%~~

$$\text{Quartile Deviation} = \frac{(Q_3 - Q_1)}{(Q_3 + Q_1)}$$

$Q.D \rightarrow$  Mean the semi variation between the upper

quartiles ( $Q_3$ ) and lower quartiles

$$S.D = \frac{\sum D}{\text{Mean}}$$

Mean  $\equiv$  Mean deviation

Deviation  $\equiv$  Average

Lower half      Upper half

12, 13, 15, 16, 19, 21, 22, 23, 32, 35, 38, 64



Median = 21

Interquartile Range (IQR)  $\rightarrow$

$Q_1$	$Q_2$	$Q_3$
20.5%	25%	25% 25%

Median =  $Q_3 - Q_1$

$$Q_1 = \left( \frac{n+1}{4} \right)^{\text{th}} \text{ term}$$

(25 percentile)

$$Q_2 = \left( \frac{3(n+1)}{4} \right)^{\text{th}} \text{ term}$$

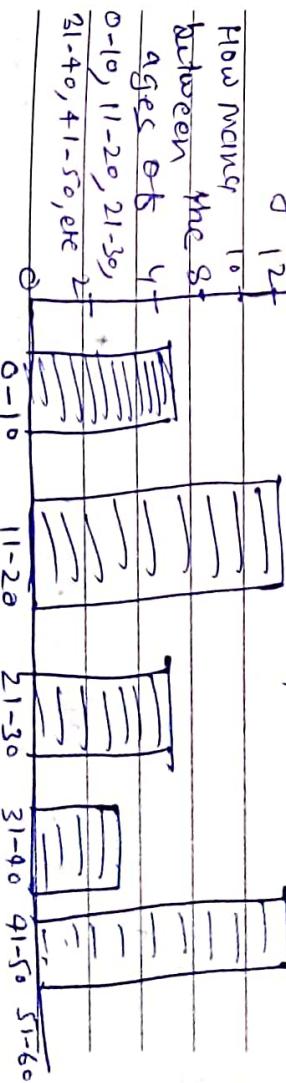
(50th percentile)

$$Q_3 = \left( \frac{3(n+1)}{2} \right)^{\text{th}} \text{ term}$$

$$IQR = Q_3 - Q_1$$

$IQR = \text{Upper Quartile} - \text{Lower Quartile}$

e.g. Age (one feature)



## Histogram

A histogram provide a visual representation of distribution of a dataset

Skewness of the data  
It helps to visualize whether the distribution is symmetric or skewed left or right

Histogram is used to summarize discrete or continuous data

Histogram is represent mostly by Bar chart.

Statistics to demonstrate how many of a certain type of variable occurs within a specific range

Covariance  $\rightarrow$  is a measure

of the relationship between two random variables.

— Histogram, it provide a visual

representation of numerical

data by showing the number of data

points fall within a "specified

range of values" (bin is called)

It is similar to vertical graph

CENTRAL LIMIT THEOREM

$$\mathcal{N}(CD(\mu, \sigma^2))$$

maybe / maybe

belong to

where  $n \geq 30$  $n$  is data point is grain

then it equal to 30 or more.

Suppose we take 30 random sample =

$$S_1 = n_1 - \dots - n_{30} = \bar{n}_1$$

$$S_2 = n_1 - \dots - n_{30} = \bar{n}_2$$

$$\vdots = \vdots$$

Correlation equation

$$S_{100} = n_1 - \dots - n_{100} = \bar{n}_{100}$$

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

$$N \sim \text{GD}\left(\mu, \frac{\sigma^2}{n}\right)$$

that is a  
derivation of  
100 Random sample

Suppose we have a example of house pricing  $\rightarrow$ 

Size	Price
1200 Sqft	10000 \$
1300 Sqft	15000 \$
1400 Sqft	20000 \$

See

Size	Price
1200	10000
1300	15000
1400	20000

Quantify relationship between them

(size increase) Size  $\uparrow$  Price  $\uparrow$  (price increase)(size decrease) Size  $\downarrow$  Price (price decrease)(size decrease) Size  $\downarrow$  Price (price decrease)

Covariance = the change in one variable is equal to change in another variable

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

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

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

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

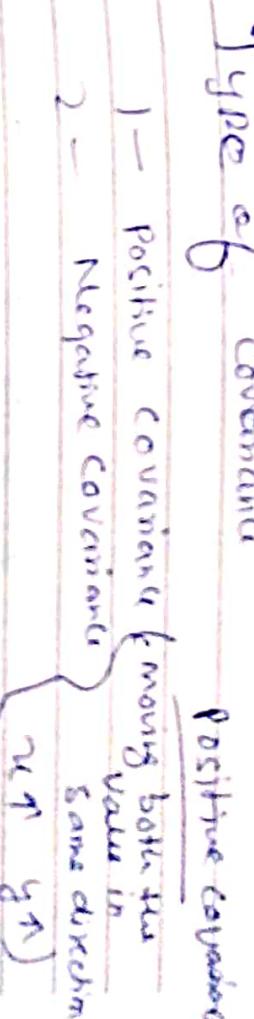
where  $\bar{y}_x = \text{mean of } x$

$$x \uparrow y \uparrow = \boxed{\quad} \text{ +ve values}$$

$$x \uparrow y \downarrow = \boxed{\quad} \text{ -ve values}$$

### Type of covariance

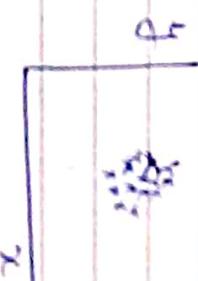
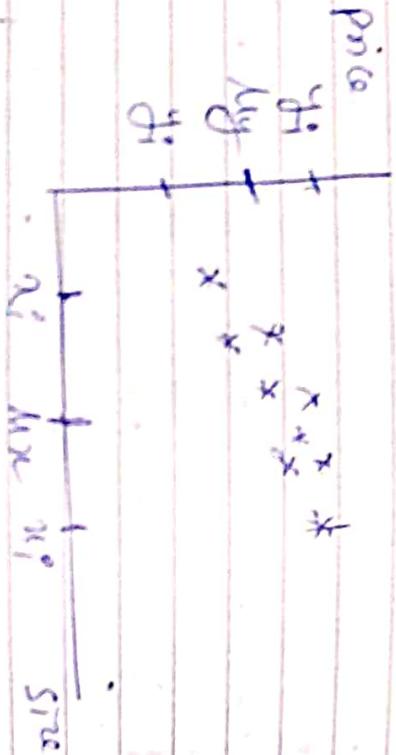
- 1 - Positive covariance (moving both the value in same direction)
- 2 - Negative covariance



$\text{cov}(x, y) > 0$   $\rightarrow$   $\text{cov}(x, y) < 0$  (less than zero)

positive

Negative



$\text{cov}(x, y) = 0$  (is zero)

new no relation

## Pearson Correlation Coefficient

it only happen for linearly data

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

$$2 - \text{Pearson CC} = \rho(x, y) = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y}$$

$\sigma_x = \text{Standard deviation of } x$      $\sigma_y = \text{Standard deviation of } y$

This is use in feature selection process.

$$3 - \text{Spearman Rank Correlation Coefficient}$$

Covariance Revise

$$\text{height}(x) \quad \text{weight}(y)$$

$$\text{cov}(x, y) = n \uparrow \quad y \uparrow \quad (\text{increasing})$$

$$\mu \uparrow \quad y \downarrow \quad (\text{decreasing})$$

What is the relationship between  $x, y$

is say the direction of relationship

positive, bland, How Negative

## Covariance is just find

The direction of relationship of How Positive, Bland, How Negative

## Pearson Correlation Coefficient

Activity help to find the

1 - Strength

2 - Direction of Relationship

$$n \uparrow \quad y \uparrow$$

\* 1 -  $\text{cov}(x, y)$  says only the relationship of two random variable, how much positive, but never said how much correlated?

So this problem is solved by Pearson correlation Coefficient.

$$\rho \text{ is lies } -1 \leq \rho \leq 1$$

$\rho$  (Pearson CC) is range from

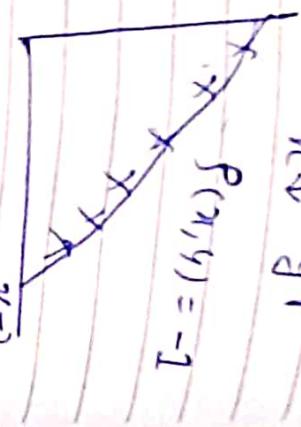
$$-1 \text{ to } +1$$

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

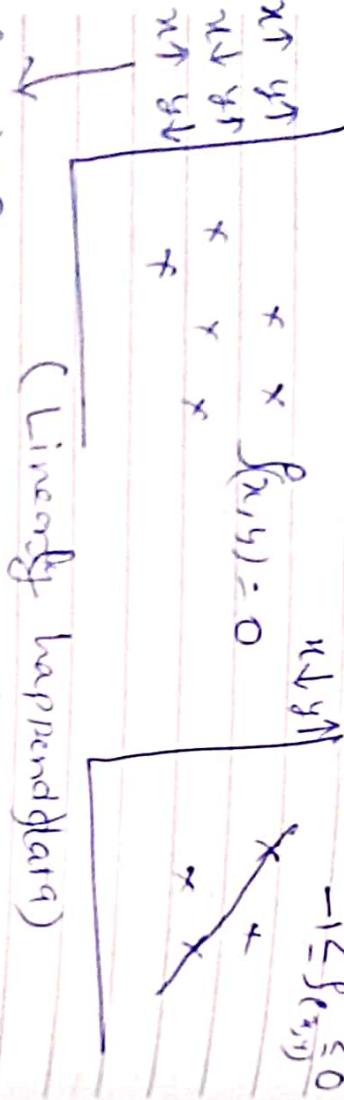
Range -1 to 1

$$n \downarrow y \uparrow$$

$$\rho_{x,y} = -1$$



$$-1 \leq \rho_{x,y} \leq 0$$



(Linearly independent)

$$\rho_{x,y} = 0$$

These techniques are best for above Selection.

Suppose we have three feature columns if

$x_1, x_2$  is  $x_1 \uparrow x_2 \uparrow$   
then Remove one feature. Now keep

Spearman's Rank Correlation Coefficient

This is useful on non-linear Structure Data.  
Pearson CC result = 0.88



→ its Non-linear → in data

formula

$$\text{Cov}(X_{\text{rank}}, Y_{\text{rank}})$$

$$\rho_{(X_{\text{rank}}, Y_{\text{rank}})} = \frac{\text{Cov}(X_{\text{rank}}, Y_{\text{rank}})}{\sigma_{\text{rank}} \times \sigma_{\text{rank}}}$$

using famous Spearman Rank the  $X_{\text{rank}}$ ,  $Y_{\text{rank}}$ ,

DATE — / /

$$\rho_{\text{Spear}} = 1 - \frac{6 \sum d_i^2}{n(n^2-1)}$$

Check the wiki

Spearman Rank CC is help us  
find spee non-linearly data  
to how correlated to each  
if very close then result  
will be 0.96 or 0.85 kind  
and very low - 0.0076 from  
Nesgatral

We Can - Alcohol is

Very High So Scale  
down

Normalisation

Normalisation  $\rightarrow$  In this approach we

Will Scale down the values of features between 0 to 1

$$X_{\text{norm}} = \frac{X - X_{\text{min}}}{X_{\text{max}} - X_{\text{min}}}$$

Using MinMaxScaler in sklearn  
few data sample

Class	Alcohol	Malic
1	14.23	1.71
1	13.20	1.78
1	13.16	2.26
1	14.37	1.95
1	13.24	2.59

Output after applied

MinMaxScaler()

$$\left[ \begin{array}{l} 0.869532, 0.233201 \\ 0.783261, 0.21834 \end{array} \right]$$

It Mean both the column Alcohol & Malic is scaled down to 0 to 1 form.

Standardization

(Z-Score Normalization)

Here all the feature will be transformed in such a way that it will have the properties of standard normal distribution with mean ( $\mu = 0$ ) and Standard deviation ( $\sigma = 1$ )

$$\bar{Z} = \frac{N - \mu}{\sigma}$$

Using Sclar Using StandardScalar  
Value Considering the mean=0

$$\alpha = 1$$

of both column (Alcohol & Matrix)

$$\left[ \begin{array}{c} 0.06732, -0.542316 \\ \cancel{1.03467}, -0.326476 \end{array} \right]$$

Basically  $\rightarrow$

When we use lot in machine learning term is StandardScaler because perform well.

But MinMaxScalar is bit not much equal to StandardScaler so in lot term StandardScaler is use more.

DL forming they use ~~Normal~~ CNN which Max Scalar because it's converted 0 to 1 form.

But DL is use MinMaxScalar because it's convert 0 to 1 to good to if open scale is 0 to 1 to good to go with.

~~Because~~ Basically,

its scale down the feature, it

bring happen time & the

gradient descent. Actually

reduce the weight & gets the

new weights.

CNN

Linear Regression

KNN

K-means

ANN

K-nearest neighbour

No Scaling

Decision tree

Random forest

XGBoost

Boosting All technique

Bagging tree

Decision

branch create

# Right Skewed

DATE \_\_\_\_\_

## Skewed Distribution

Q1 - What is the Right skew distribution?

Q2 - What is the Left skew distribution?

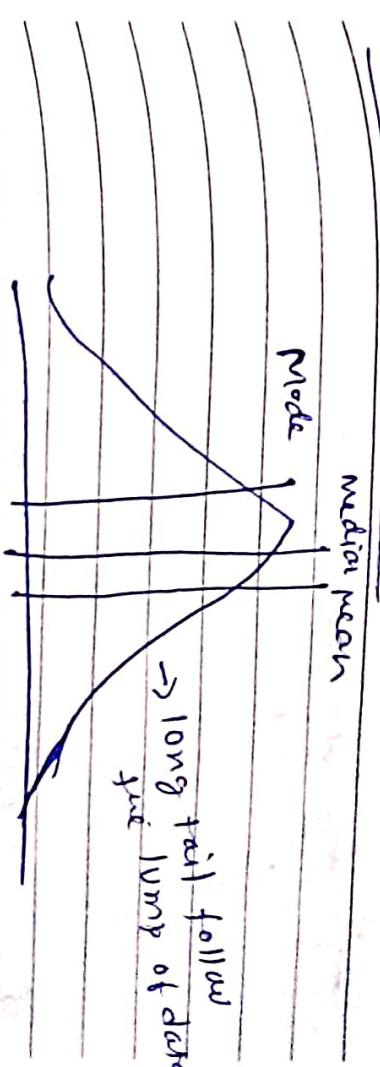
Q3 - Classical example of

Right Skewed distribution & Left skewed distribution?

Q4 - What is Relationship bet.

Mean, Median, Mode with Right skewed & left skewed distribution?

~~Right Skew~~



Positive

Right skew mean > median > mode

Right hand Side is tail aligned to it

Right - wealth distribution

- it mean lot of people

wealth is more like Amancio  
Tita to normal one

→ Cricket Score (few score good pattern)

→ Exam Results (many low, less high  
median result)

~~Left Skew~~

mode

R.G - life  
span

ob

Human

② Human life  
cycle

(some die early  
before the age)

Left skew mode > median > mean → Relationship

## Sampling techniques 1

mean  
traction

Normal distribution

Q - age distribution mode

weight "

height "

TRSS "

In childhood, in cricket one person randomly choose number of batting position of every player, while another person pointing out the number by hand backside of 1st person, this is called simple random sampling

- Q - What is Sampling ?  
Q - What are different sampling technique ?  
Q - What are types of sampling

## Q - Stratified random Sampling

In school prayer in the morning there is particular line for every class and different line for boys and girls, so here there are n groups of particular class student. And here if teacher choose one boy from every class or one girl from every class this is Stratified Random Sampling

## 3 - Cluster random Sampling

length

DATE - 11

Assume in ~~stretten~~ school suddenly

one guest come and principal ordered every teacher to all students to come at play ground, This is an immediate order so teacher and student can't make proper line according to class wise, Here there can be or can't be students of different class are stood in one line, so here every line is called as cluster, How teacher order one student from every line, This is called cluster Random Sampling.

4 → Systematic Random Sampling →

Assume everyday in school assembly, teacher order students by their roll number from every line here it can be stratified or cluster Random Sampling), so take

randomly from group, after

random position, This type of sampling is a systematic Random Sampling.

What's average ~~length~~ ~~of~~ shark in a sea?

(1) Hypothesis testing → Let's assume theoretical value of shark length around 2 feet

H0 (Null Hypothesis): The population mean (2 feet) is same as sample mean.

H1 (Alternate hypothesis): The population mean is different

Set the P Value as 0.05 and collect the shark length data by using Z-test / T-test we will conclude whether we can accept the Null or Reject the Null.