

Introduction to Statistics in Data Science

we have two statistics

Descriptive stats

measure of Central Tendency
measure of Dispersion

Summarizing the data

Histograms, PDF, CDF....

probability

Permutations, mean,
median, mode, variance,
standard deviation.

Different types of
distribution

What is

STATISTICS

Statistics is the science of collecting, Organizing and
Analyzing data. { we are doing this for better Decision making }

Definition of DATA

Data → facts (or) pieces of information that can be

measured.

e.g.: The IQ of a class

{ 98, 97, 60, 55, 75, 65 }

Ages of students of a class

{ 30, 25, 24, 23, 21, 28 }

Types of Statistics

differential descriptive
differential

① DESCRIPTIVE STATS

It consists of organizing and summarizing of data

① measure of central tendency (mean, median, mode) b) measure of dispersion (var, σ)

② INFERENTIAL STATS

(using some experiments) Z test, t test ...

It is a technique where, we used the data that we have measured to form Conclusions (or) Inferences

we sample data $\xrightarrow{\text{perform experiments}}$ conclusions \rightarrow then conclude for other data {population data}

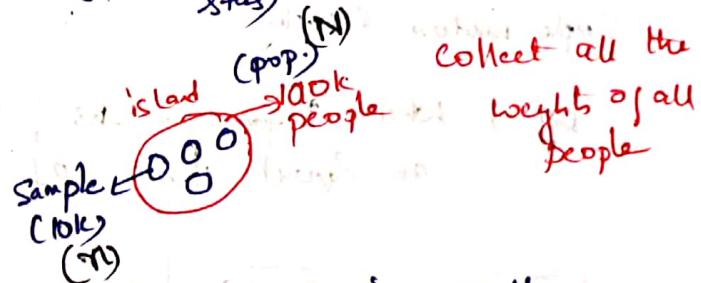
e.g.: let us say I have a classroom of Maths Students (20)
want to find marks of the 1st sem

marks $\Rightarrow 84, 86, 78, 72, 75, 65, 80, 81, 82, 95, 96, 97, \dots$

Q \Rightarrow What is the average ^{marks} of the students in the class (eg of descriptive class)

e.g.: Are the ^{marks} of the ^{students} of this classroom similar to the ^{marks} of the ^{college} (eg of inferential stats)

Population and Sample



e.g. Elections

which state are probably having the election recently

let (Goa, UP)

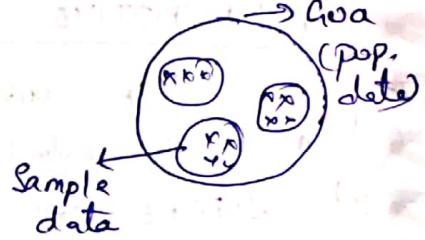
let election has finished and we really need to find out the exit poll

exit poll \Rightarrow they cannot go and ask each & every person that whom they have Voted. They report what they do is, they take up sample of population

from different different regions and they ask that whom do you vote, based on that maximum number of people whom they vote, they actually create their exit poll.



In this example, my population data is the entire population of Goa and we take some people in some regions, those are sample data. Population is denoted by N . Sample is " " in



Note \Rightarrow why you have selected samples randomly (or) is there any better ways to do sampling also (or) just we need to do the sample randomly.

Sampling Techniques

① Simple Random Sampling \Rightarrow

In this, we pick randomly some samples in the population. There is no any format, just we take randomly.

In exit poll, we can use random sampling,

Suppose, If we want to test for the medicines, then we cannot use simple random sampling.

Dg. \Rightarrow When performing SRS, every member of the population has an equal chance of being selected for your sample(n).

② Stratified Sampling \Rightarrow

Dg. \Rightarrow It is a technique where the population (N) is split into

\downarrow non-overlapping groups

Eg.: Gender

male

female

let us say I want to do survey for this I require some people, based on that gender, my samples will be divided

Eg. I want to do Survey on 0 to 10 years of kids Based on Age, I will take samples.

0-10 (10-20) (20-40) ...

Q \Rightarrow Can we do stratified Sampling based on profession?

No. Python, C, ...

we can apply this stratified sampling on doctors, engineers ...

③ Systematic Sampling →

In this, from population (N), we pick up every n th individual and will take as a Sample.

Eg.: I am outside the mall & I want to do a survey regarding Covid

I will do, every 7th or 8th person I see, for this person do the Survey. This is called as Systematic Sampling. There is no reason for selecting 8th or 9th person, you just said that take every 7th person & do survey.

④ Convenience Sampling →

lets consider that I am doing a survey only those people who are a domain expert in that particular survey, will be participating in this survey.

Let us say, I am doing the survey related to data science, I will say that any person who is probably interested in data science, if you consider only those people then it basically becomes a Convenient Sampling.

→ FBI do Survey with respect to household, for this household surveys, what kind of sampling they may use?

In household surveys, RBI make sure that they have to fill the survey from women, where probably they are trying to find out what is the cost expenditure owning a house etc.

Here we can do Stratified Sampling / Convenience Sampling

Note Sampling techniques are completely dependant on the use case that we are following.

Eg. A) A drug need to be tested, what kind of Sampling we use? Here I can think up multiple use cases, 1st of all to whom this drug need to be tested, GP: get that specific information

I will basically do the age groupings and then I apply... let us consider, the drug is for everyone, then I may consider up some samples, but I will put atleast put a condition that it should be > 15 years. So it is based on the care you do, based on the care you will select the sampling technique.

Variables

A Variable is a property that can take on any value.

e.g. height, weight

{188, 178, 168} {78, 99, 100, 60}

Here we have 2 kinds of Variables

- ① Quantitative Variables
- ② Qualitative / Categorical Variables

① Quantitative Variables

It can be measured numerically

{Add, subtract, multiply, divide}

e.g. Age, height, weight

Quantitative

Discrete Variables

Whole number

e.g. No. of Bank Accounts

Total children in a family
(no. of)

Continuous Variables

any values

e.g. height

(172.5, 162, 163.5 cm)

weight

(100kg, 99.5, 99.25)

amt of rainfall

(1.1 inches, 1.25 inches etc...)

② Qualitative / Categorical Variables

Based on some characteristics we can derive some categorical variables. In this we have some categories.

e.g. Gender

M → Male

F → Female

IQ → Less IQ

→ medium IQ

→ Good IQ

Blood group, T-shirt size ...

Age groups, ...

Gender, ...

Height, ...

- eg:
 1) what kind of Variable Gender is? Categorical Variable
 2) what kind of " Marital status? Categorical Variable
 3) what kind of " River length? Continuous Quantitative Variable
 4) what kind of " population of State is? Discrete "
 5) what " Song length? Continuous "
 6) what " Blood pressure? Continuous "
 7) " Pin code? Discrete "

Variable Measurement Scales

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

① Nominal

These are specifically categorical data.
 eg. gender, colour, type of flower

③ Interval

order matters, Value also matters.
 natural zero is not present
 eg. Temperatures
 -10-80F 80-90F 90-100F
 0F is not meaningful

distance

10-20 20-30 30-40

② Ordinal

order of the data matters but Value does not

④ Ratio

eg. student marks

	Rank
100	1
96	2
57	4
85	3
44	5

ordinal date
 Here we focus more on the order

Frequency Distribution

Suppose I have a sample dataset, which consists of 3 types of flowers - Rose, lily, Sunflower. Rose, lily, sunflower, Rose, lily, lily.

Based on the flower type, how much is the frequency?

Flower	Frequency
Rose	3
Lily	4
Sunflower	2

frequency distribution

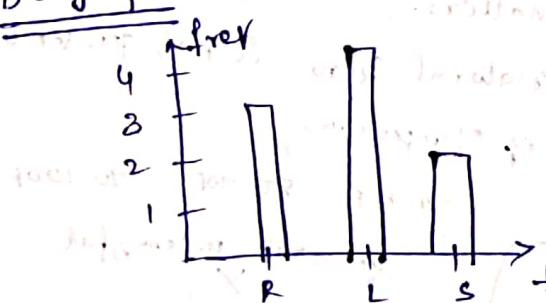
$$3 + 4 = 7$$

$$7 + 2 = 9$$

Total no of flowers

from this freq distribution, we can derive bar charts, pie charts etc

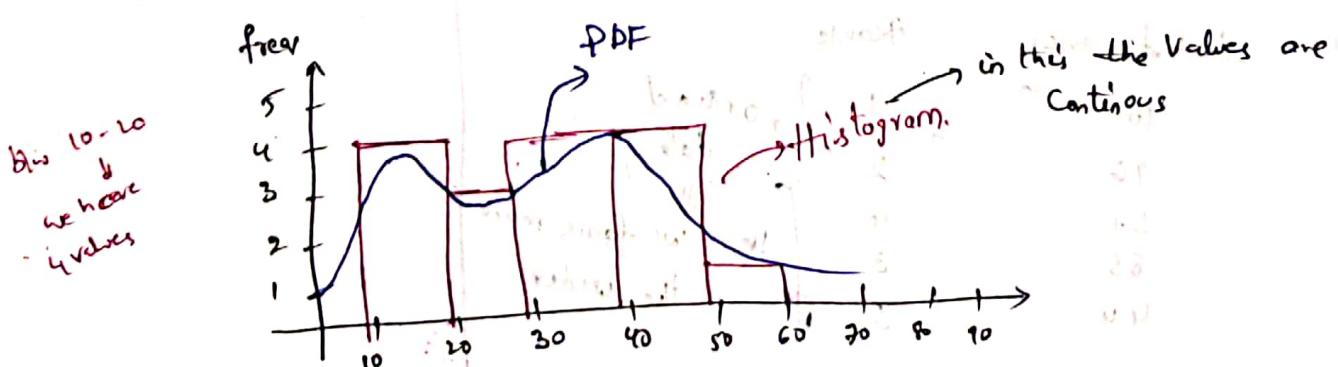
Bar graph (Data should be discrete)



Histogram (Data should be continuous)

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

In histogram, we make some kind of bins, by default the bin size = 10



PDF \Rightarrow Smoothing of histogram (it can be done using Kernel density estimator)
probability distribution function

BAR vs Histogram
used for discrete data
used for continuous data

Day 2 Basic to intermediate of statistics

Topics discussed ⇒

measure of central tendency

measure of dispersion

Gaussian Distribution

Z Score

Standard normal distribution

a It refers to the measure used to determine the centre of the distribution of data

Central Tendency

mean median mode

Arithmetic mean for population and sample.

Mean (Average)

Random variable population (N)

Let $X = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 7\}$

$$\mu = \frac{N}{P=1} \sum x(P)$$

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

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

Sample (n)

Same data

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

then we get same value

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

Central Tendency:

Median

will take same data set

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

$$\text{mean} = \frac{32+100}{11}$$

$$= \frac{132}{11} = 12$$

outlier
100
it is completely outside the values in dataset

Outlier has major impact

when 100 is not added

then mean = 3.2, if

100 is added in the same distribution then it becomes 12

so there is huge difference in mean because of this 100, so

this no. is called as

OUTLIER

So in this case we use median.

median \Rightarrow

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

Step 1 Set the no's

In this eg, already it is sorted

Step 2 Take the central element

If total count = odd no.

then we take the central element.

If total count = even no.

\Downarrow
we take middle nos.

& we take avg of the 2 nos.

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

median = 3 (As it is odd no.)

If 2 outliers are added

$$\{1, 1, 2, 2, 3, \textcircled{3}, 4, 5, 5, 6, 100, 112\}$$

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

median = 3.5 (It is even)

Note:- when we use median for one outlier, then the median = 3
when we have 2 outliers, then median = 3.5.

So there is less difference

* Median works well with the Outlier

Mode

Suppose dataset - {1, 2, 2, 3, 4, 5, 6, 6, 7, 8, 100, 200}

To calculate mode, we find out the most frequent element in D.

So for our eg \Rightarrow

mode = 6

Discd

Suppose I have many outliers

Dataset

order sales

12

Type of flower

PL PW

Rose

Lily

Sunflower

In this Dataset, suppose we have missing data (10%)

Here we use mode; because

the missing value is replaced by most frequent occurring element.

Mode is specially used for Categorical Value

eg: Age of students

25
26
27
—
—
32
34
36

for this what should I apply?
(mean / mode / median)

So I prefer mean in this example

because, I know students age will basically range from one value to another value.
so here Domain Knowledge will also be used.

Measure of Dispersion

Variance → Standard Deviation

Dispersion → Spread
(how well spread the data is)

Let $D_1 = \{1, 1, 2, 2, 4\} \Rightarrow M_1 = 2$

$D_2 = \{2, 2, 2, 2, 2\} \Rightarrow M_2 = 2$

for both distributions, we get same mean, so how these 2 distributions are different?

In order to find how two distributions are different, at that pt, we use

Variance

$$D_1 = \{2, 2, 4, 4\}$$

$$M_1 = 3$$

the 2 & 4 are closer to mean \Rightarrow less dispersion

$$D_2 = \{1, 1, 5, 5\}$$

$$M_2 = 3$$

here 5 & 1 are far away from mean \Rightarrow dispersion is more

* it is used to differentiate the data

Variance \rightarrow

pop. Variance \rightarrow find the diff b/w from mean

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

pop. size

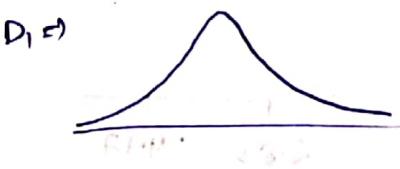
$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

why $n-1$? \downarrow
sample size

eg: $x \quad M \quad x - M \quad (x - M)^2$
1 2.83 -1.83 3.34
2 2.83 -0.83 0.6889
3 2.83 0.83 0.6889
4 2.83 +0.83 0.03
5 2.83 1.17 1.37
 $\bar{x} = 2.83$
 $s^2 = \frac{10.84}{4} = 2.71$

$$\sigma^2 = \frac{10.84}{5} = 2.168$$

let say if I have a dataset, which looks like



comparing these 2 data, where the variance is more?

In 2nd Distribution, variance is more because the spread is more.

Variance \uparrow when spread \uparrow

Elements present in the central region is more

Standard deviation \rightarrow It tells how far away from mean.

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

$$= \sqrt{1.81}$$

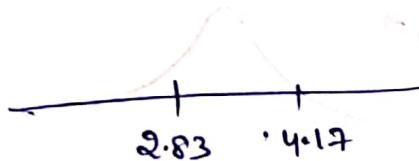
$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2$$

$$\boxed{\sigma = 1.345}$$

\rightarrow In this eg., mean = 2.83, from this mean, the data is distributed bcoz mean is basically specifying your measure of central tendency. It basically says that, where the center is there, for that particular distribution.

\rightarrow form the mean, if we go one standard deviation to the right, the next element that may probably found/fall b/w $2.83 + 1.345$
 $= 4.17$

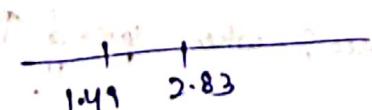
one SD



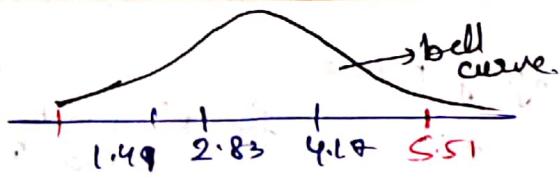
whatever elements are basically present b/w 2.83 & 4.17 will be falling within the 1st standard deviation.

If we consider something towards the left, we will subtract

$$2.83 - 1.34 = 1.49$$



Any element falls b/w 1.49 to 2.83 will be falling in this region i.e. one SD to the left



s.d is a very small number
graph looks like above.

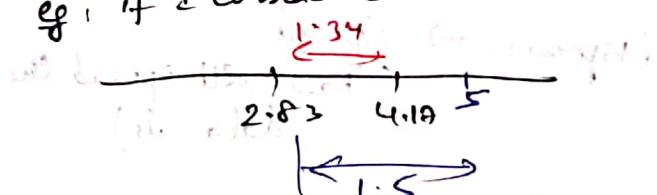
Based on SD and Variance we will be able to decide 2 imp things:

Variance \Rightarrow spread

SD \Rightarrow B/w one SD to the right and the left what may be the range of data that may be falling.

Standard deviation is ^{square} root of Variance ie from the mean how far a element can be

eg: if I consider 5



if you try to calculate, it may fall somewhere here
so SD is represented as 1.34σ from the mean.

So if we consider 1.34 σ to the left of the mean, we will get 1.49

Percentiles & Quantiles

This is the 1st step to find Outliers

percentage $\Rightarrow 1, 2, 3, 4, 5$

i.e. of the nos. that are odd?

$$q_1 = \frac{\# \text{ of numbers that are odd}}{\text{Total no's}}$$

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

8*

Percentile \Rightarrow

(CHAT, CAT, GMAT, SAT)

Definition:- It is a value below which a certain % of observation lie?

e.g. Dataset

2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 9, 9,
10, 11, 11, 12

Q? I want to find out, what is the percentile ranking of 10?

So percentile Ranking of 10

let $x = 10$

$$\text{percentile ranking of } x = \frac{\# \text{ of values below } x}{n} \times 100$$

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

= 80 percentile

i.e. 80 percentage of the entire distribution is less than 10
real meaning

a) what is the percentile ranking of 11?

$$\text{percentile ranking of } x = \frac{17}{20} \times 100 \\ = 85\%$$

b) what value exists at percentile ranking of 25%.

Formulae =

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

$$\text{Value} = \frac{25}{100} (20+1)$$

$$= \frac{25}{100} \times 21$$

$$\text{Value} = 5.25$$

it is the index position, not the Value.

in my dataset

2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 9, 9,
10, 11, 11, 12

5.25 index is not there so we take average of 5th & 6th index

$$\frac{5+6}{2} = \frac{11}{2} = 5.5$$

b) what value exists at percentile ranking of 75%.

$$\text{Value} = \frac{75}{100} (21)$$

$$= 15.75 (\text{index})$$

$$\begin{aligned} 15^{\text{th}} &\rightarrow 8 \\ 16^{\text{th}} &\rightarrow 9 \end{aligned} \quad \left\{ \frac{8+9}{2} = \frac{17}{2} \right.$$

Five Number Summary

- ① minimum value
- ② 1st Quartile (Q_1)
- ③ median
- ④ Third Quartile (Q_3)
- ⑤ Maximum value

Removing the Outlier

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

what is the outlier in this dataset

27

In order to remove outlier, we need to find lower fence & higher fence

nos < lower fence & nos > higher fence

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

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

$IQR = \text{Inter Quartile Range}$

$$= Q_3 - Q_1 \quad Q_3 = 75\% \quad Q_1 = 25\%$$

$$Q_1 \rightarrow 25\%$$

$$\text{Value} = \frac{25}{100} (19+1)$$

$$= \frac{25}{100} (20)$$

$$\text{Value} = 5 \text{ (index)} \quad Q_1 = 3$$

$$Q_3 \rightarrow 75\%$$

$$\text{Value} = \frac{25}{100} \times (20)$$

$$= 15 \text{ (index)}$$

$$Q_3 = 7$$

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

$$= 7 - 3$$

$$\text{IQR} = 4$$

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

$$= 3 - 1.5(4)$$

$$> -3$$

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

$$= 7 + 1.5(4)$$

$$= 19$$

[lower fence \longleftrightarrow higher fence]

$[-3 \longleftrightarrow 19]$

so any thing < -3 and anything > 19 is considered as outlier in the dataset D.

So 27 becomes as

outlier and we should remove it.

Remaining data is

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

minimum value = 1

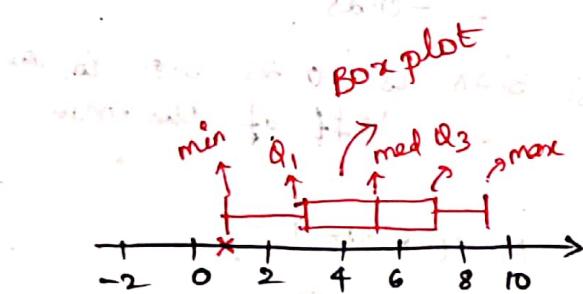
1st quartile = 3
(Q1)

median = 5

3rd quartile = 7
(Q3)

Max value = 9

By this specific data we can draw a Boxplot



This technique of removing an outlier we basically say with respect to lower fence & higher fence & we use IQR.

Day 3

Distributions

Let we have a dataset

Ages = {24, 26, 27, 28, 30, 32, ...}

Here we have a dataset, so how do we basically see this dataset in a visualized way:
If I want to start analysis I really need to see lot of visualized diagrams of that & where, when consider this entire distribution. Here are multiple ways to visualize this data through various graphs.

As we have ages dataset, I want to plot this data and the best & easy way that can probably think about is Histogram

So we will go to distributions

① Gaussian Distribution / Normal distribution

If I have a distribution

if this distribution

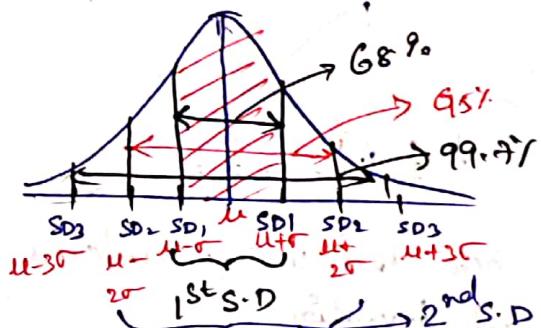
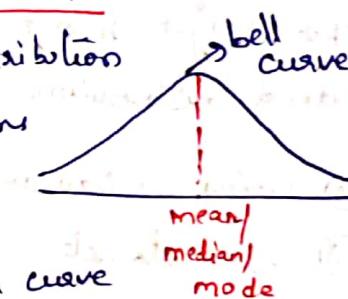
is like a bell

curve, main

thing in this bell curve

is symmetrical on b.s.

The right part data present = left part data



Empirical Formulae

(68-95-99.7% rule)

68 → ~~within~~ B/w the 1st standard deviation region, around 68% of the distribution is present

let $D = 100$ pts

so out of 100 pts, 68 pts present in the 1st S.D. region (so it is in form of bell curve)

95 with in 2nd S.D. region, around 95% of the entire data lie in this region

99.7% \rightarrow Within 3rd S.D region, around 99.7% of the entire distribution will fall in this region.

Q1) Height

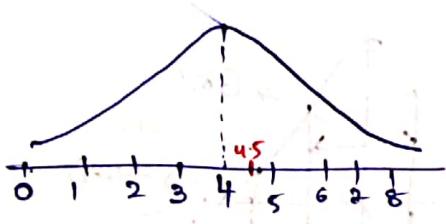
It is normally distributed

So Domain expert (doctors) says that height belongs to Gaussian/ normal distribution.

Q2) Weight \rightarrow Gaussiandest

Q3) Fridge set

Ex: Suppose I have $\mu = 4$, $\sigma = 1$
then can I find distribution?



Q1) Ques where does 4.5 falls in terms of standard deviation.

Ans: it is 0.5 SD right

Q2) $4.25 ?$

the $S.D = 1$, in case of 4.25
it is very much difficult to do
do the calculation

so we can use Zscore,
PE used to find out how much
SD away from the mean

$$Zscore = \frac{4.25 - 4}{1}$$

$$= 0.25$$

$\therefore 4.25$ is 0.25 S.D to the right

As it is +ve \Rightarrow So it is right to mean

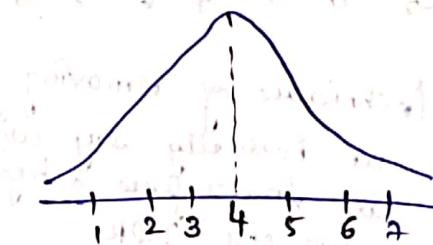
Q1) 3.75 ?

$$Zscore = \frac{3.75 - 4}{1}$$

$$= -0.25$$

$\therefore 3.75$ is 0.25 S.D to the left of the mean

$\rightarrow \mu = 4$ & $S.D = 1$



lets apply Zscore to every values

$$Zscore = \frac{x_i - \mu}{\sigma}$$

initially my distribution is

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

after applying Zscore my distribution is

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

$$Z(1) = \frac{1-4}{1} = -3 \quad Z(3) = 1$$

$$Z(2) = \frac{2-4}{1} = -2 \quad Z(6) = 2$$

$$Z(3) = \frac{3-4}{1} = -1 \quad Z(7) = 3$$

$$Z(4) = 0$$

$\{1, 2, 3, 4, 5, 6, 7\}$, normal/gaussian distribution

zScore

\downarrow

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

This distribution is

called as

STANDARD NORMAL DISTRIBUTION

$\mu = 0 \text{ & } \sigma = 1$

$y \in SND(\mu=0, \sigma=1)$

Standard normal distribution

Q2:- Why do we do this?

we do this in ML algorithms.

Suppose Dataset is

Age (years)	Salary (Rs)	Weight (kg)
24	40K	70
25	90K	80
26	60K	55
27	70K	45

One main target is, that we should bring up in a form where $\mu=0$ & $\sigma=1$ at that pt of time we can apply standard normal distribution. I can scale up this entire data and apply zscore & convert this into SND.

This process is called as

STANDARDIZATION

$\mu=0 \text{ & } \sigma=1$

whenever we talk about Standardization, internally there is a Zscore formulae applying on it.

NORMALIZATION

→ let us say I want to change my all values in dataset b/w 0 to 1 then we use Normalization.

This normalization can be done by MinMax Scaler.

In this you have to provide 0 to 1 & automatically this kind of normalization happens.

Normalization gives you a process where you can basically define the lower bound and upper bound and you can convert your data b/w them.

Where do we use normalization?

In CNN, whenever you are doing image Training / Object detection, every image has a pixels, each image has a pixels, each pixel ranges b/w 0-255 and every pixel ranges b/w 0-255.

Before we start training this can be applied with min-max scaler and it gets converted b/w 0 to 1 where the min value '0' is assigned to 0 & the max value 255 is assigned to 1.
So here we are doing Normalization.

Another type of normalization $\{0 \text{ to } 255\} \rightarrow \{255 \text{ to } 0\}$

$\{0 \text{ to } 1\}$

Practical ex :-

Cricket match
{ India vs South Africa }

ODI series { 2021 }

The series average $2021 = 250$,

The S.D of the Score = 10
(Standard deviation)

Rishabh final score = 240
(Team) avg { 2020 }

Series average 2020 = 260

S.D = 12

Rishabh final score = 245
(Team) avg

Question -> Compare to both the series, in which year Rishabh's final score was better?

Ans

2021

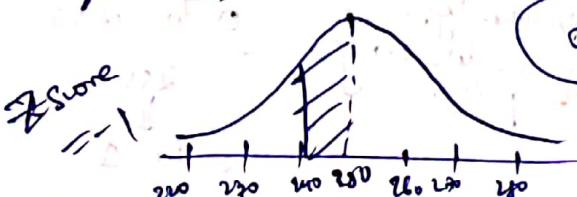
$$\text{Zscore} = \frac{x_0 - \mu}{\sigma}$$
$$= \frac{240 - 250}{10} = -1$$

2020

$$\text{Zscore} = \frac{x_0 - \mu}{\sigma}$$
$$= \frac{245 - 260}{12} = -1.25$$

In 2021, Rishabh

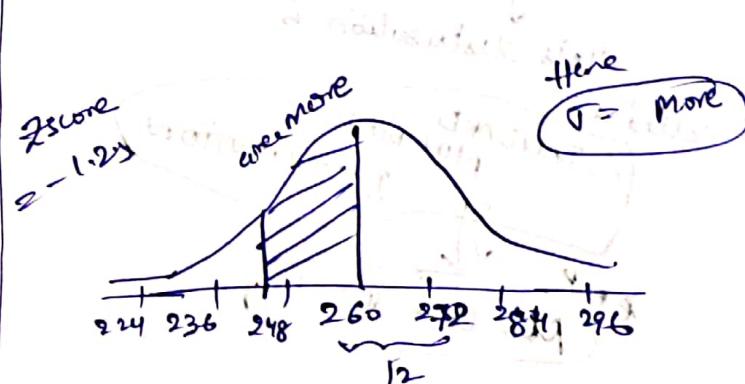
$$\mu = 250, x_0 = 240, \sigma = 10$$



240 is falling into -1 standard deviation

In 2020

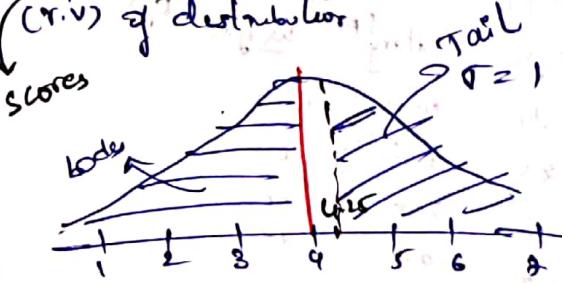
$$\mu = 260, x_0 = 245, \sigma = 12$$



India performed well in 2021
becoz Zscore is more for 2021

Practical ex 2

X = it has the below kind (R.V) of distribution



Q1 what % of Scores fall above 4.25 ?

meaning is, what is the % of scores that are greater than 4.25

Ans here we use Zscore

$$\text{Zscore} = \frac{x_0 - \mu}{\sigma}$$
$$= \frac{4.25 - 4}{1}$$

$$\text{Zscore} = 0.25$$

-4.25 falls 0.25 standard deviation from the mean.

$$Z\text{-score} = -0.25$$

This is a symmetric bell curve.
i.e. the entire area can consider it as one.

Since I am interested in above 4.25 region, ^{I say as} this region has tail. The remaining region is said as Body.

Based on Z-score, I need to find the area of the body curve.

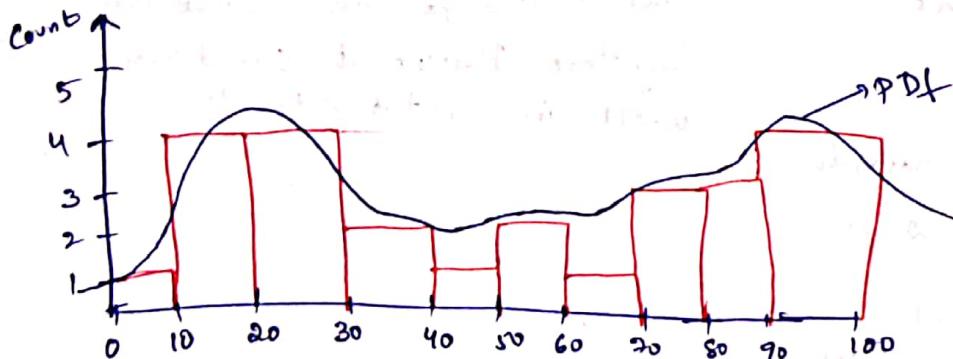
Histograms

Dataset =

10, 12, 13, 14, 20, 22, 24, 25, 26, 35, 38, 42, 45, 56, 68, 82, 84, 86, 92, 93, 94, 100

These nos. are Quantitative nos.

bin = 10



If we smoothen this histogram, it gets converted into PDF

Smoothed Version of Histogram
= PDF

This histogram will also useful for Qualitative data

Standard Normal distribution
(continuation of Gaussian dist)

$$y \approx SND(\mu=0, \sigma=1)$$

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



$$y \approx SND(\mu=0, \sigma=1)$$

To convert Normal distribution to SND, we use Z-score.

$$z = \frac{x - \mu}{\sigma}$$

$$\text{eg: } \theta = 1, 2, 3, 4, 5$$

$$\mu = 3$$

$$\sigma = 1.1 \quad (\text{let } \sigma = 1)$$

$$z = \frac{x - \mu}{\sigma} = \frac{1 - 3}{1} = -2$$

$$z_2 = \frac{4 - 3}{1} = 1$$

$$z_3 = \frac{3 - 3}{1} = 0$$

$$z_4 = \frac{4 - 3}{1} = 1$$

$$z_5 = -1$$

$$\theta = \{-2, -1, 0, 1, 2\}$$

$$\mu = 0 \quad \sigma = 1$$

Q) Why we are doing this?
(converting any dist. into SND)

e.g. Suppose

Dataset

	Age (years)	wt (kg)	distance (km)
24	100	100	100
25	120	200	200
26	75	500	500
27	80	600	600

In this dataset, all the units to measure the features are different.

In one ML model, we should try to bring all the features in the same scale.

In this kind of dataset, in order to bring into same scale, we apply z-score to each & every feature & convert into SND.

Application of Z Score

Z scores → These are standardized values that can be used to compare scores in different distributions.

e.g. Cricket (India vs Australia)

2020 2021

2020 (Tests)

$$\text{Test avg} = 181 \\ SD = 12$$

$$\text{Rishabh pant} = 187$$

2021

$$\text{Test avg} = 182 \\ SD = 5$$

$$\text{Rishabh pant} = 185 \\ \text{find z score}$$

Compared to the rest of the entire test series in which year was Rishabh pant's score in first game better?

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

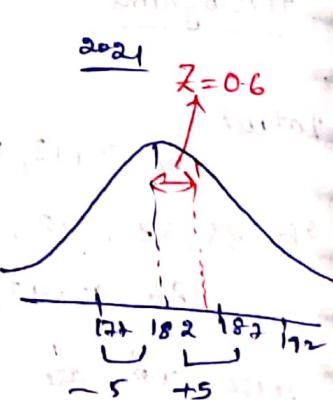
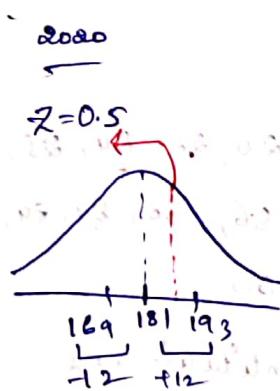
2020

$$x = 187, \mu = 181 \\ Z = \frac{187 - 181}{12}$$

$$Z = 0.5$$

2021

$$Z = \frac{185 - 182}{5} = \frac{3}{5} = 0.6$$



Since $Z = 0.6$ for 2021, so we can say that Rishabh Pant scored very well in 2021 test series. ($0.6 > 0.5$)

Eg of Zscore

a) In India the avg IQ = 100, with a standard deviation of 15. what % of the population would you expect to have an IQ lower than 85?

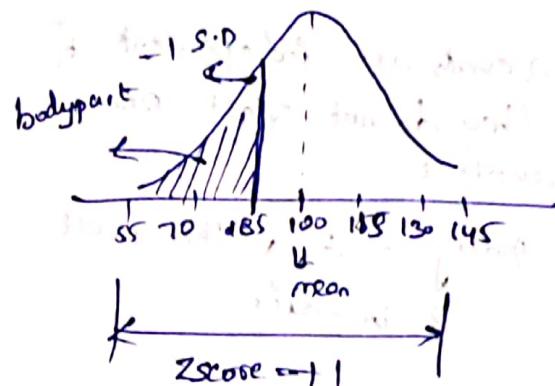
$$\text{Eq: } Z\text{ score} = \frac{x - \mu}{\sigma}$$

$$\text{solution: } z = \frac{85 - 100}{15} = \frac{-15}{15} = -1$$

from Z-table, the area of the body curve is 0.84134 \Rightarrow for +1

$$1 - 0.84134 = 1.5866.$$

for -1 \Rightarrow 1.5866%.



Basics of probability

probability \Rightarrow it is a measure of the likelihood of an Event/experiment

e.g. flipping of a coin.

What is the prob. of getting head = ?

by sample space = {H, T}

probability = $\frac{\text{No. of ways an event can occur}}{\text{No. of possible outcomes}}$

$$= \frac{1}{2} = 0.5\%$$

Addition Rule

Mutual exclusive events \Rightarrow 2 events are mutual exclusive if they cannot occur at the same time.

e.g. when flipping a coin, the 2 events Head & tail are mutual exclusive.

Non-mutual exclusive events \Rightarrow 2 events are non-mutually exclusive if they can occur at the same time.

e.g. Deck of Cards \Rightarrow when picking a card randomly, the 2 events Queen of Hearts.

Multiplication Rule

Independent Event

2 events are independent if they do not effect one another.

e.g. rolling a '6' & then roll '3' is a dice

Dependent Event

2 events are dependent, if they effect one another.

e.g. consider we have Bag of marbles which has 4 red & 2 black marbles.

$$P_r(Red) = \frac{4}{6} = \frac{2}{3}$$

$$P_r(B) = \frac{2}{5}$$

multiplication rule for indep. events

Q: what is the prob. of rolling 6 & then '3' with a normal 6-faced die?

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

$$P(6 \text{ and } 3) = P(6) * P(3)$$

$$= \frac{1}{6} * \frac{1}{6} = \frac{1}{36} = \frac{1}{3}$$

Q) multiplication rule for dependent events
what is the prob. of drawing a Queen, then drawing a King from a deck of cards?

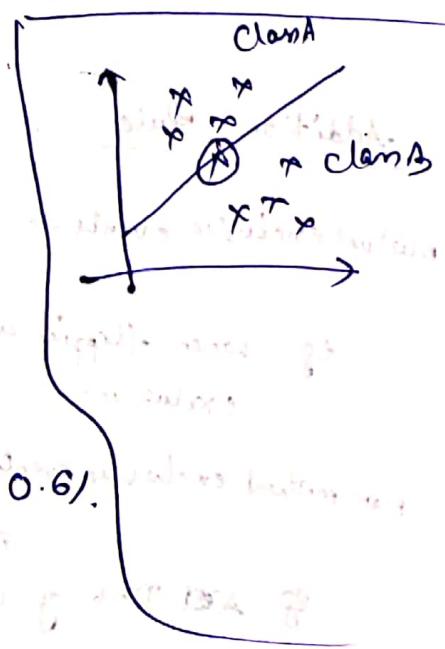
$$P(A \text{ and } B) = P(A) * P(B|A)$$

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

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

$$P(K|Q) = \frac{4}{51}$$

$$P(Q \text{ and } K) = \frac{4}{52} * \frac{4}{51} = 0.06 = 0.6\%$$

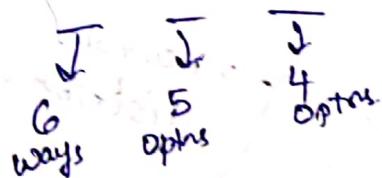


Permutations and Combinations

Permutations \Rightarrow

e.g.: In a school trip, there are 50 students, we go to a chocolate factory, I have given a task that note down the 1st three chocolates that you probably see?

Chocolate factory (6 chocolates)



$\Rightarrow 6 \times 5 \times 4 = 120 \Rightarrow$ i.e., they have 120 options (all) permutations to be write the chocolate name.

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

$n =$ total no. of objects
 $r =$ # of objects you are picking

$$6_{P_3} = \frac{6!}{(6-3)!} = \frac{6!}{3!} = 120$$

Combinations \Rightarrow

b) Write down the 1st three chocolates that you see, but they should not repeat

$$n_{Cr} = \frac{n!}{(n-r)!r!} = \frac{6!}{3!3!}$$

$$6_{C_3} = 20 \text{ unique combinations}$$

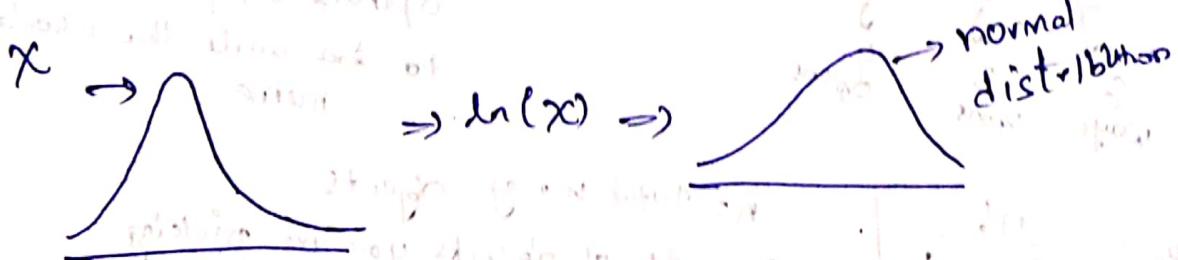
\therefore combinations < permutation

Log Normal Distribution

w.k.t. $X \approx N(\mu, \sigma)$

$X \approx \text{log normal distribution}$
If then

$y = \ln(X) \Rightarrow$ it has a normal distribution



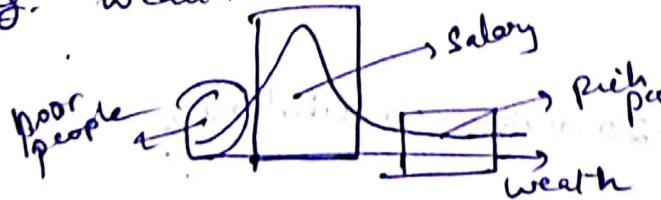
then $X \approx \text{log Normal distribution}$

if the R.V X is a log-normally distributed : then $y = \ln(X)$ is a normal distribution.

$$\boxed{y = \ln(x)}$$

$$x = e^{y}$$

e.g. ① wealth distribution

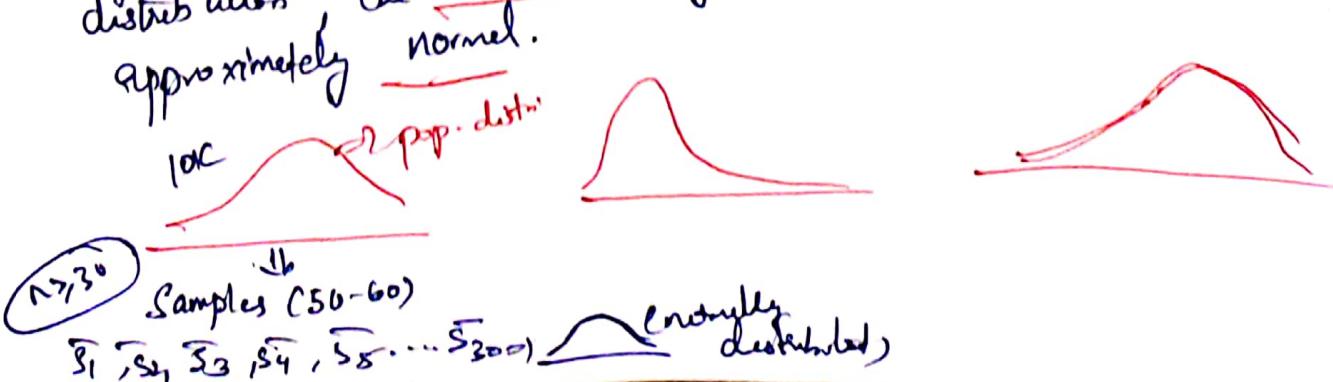


② length of the comment
in any channel.

③ salaries of people in companies.

Central Limit Theorem

The CLT states that regardless of the shape of population distribution, the distribution of sample means will be approximately normal.

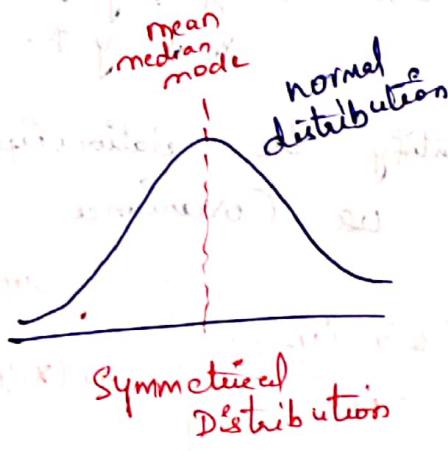
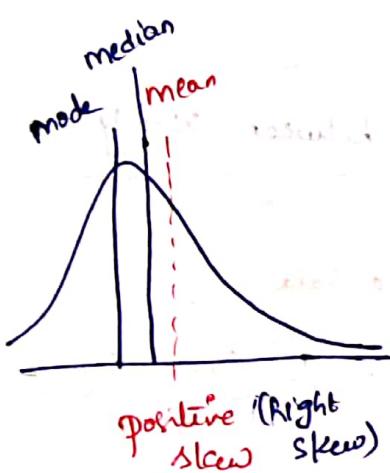


Here $n \geq 30$

2 imp points in CLT,

- ① The distribution of Sample means will become normal as its sample size increases.
- ② Good rule thumb:- Sample distribution will be approximately if their sample size is $n \geq 30$

Left skewed & right skewed distribution



e.g.: wealth distribution
length of comments

Mean > median > mode

In this most of the scores occur at the lower end of distribution & few scores at higher end.



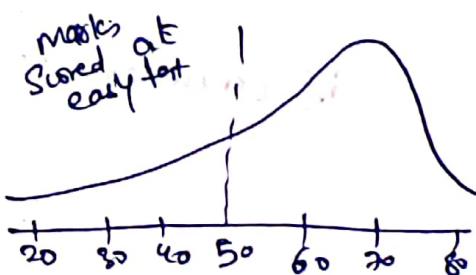
e.g.: Age
weight
height
IPIC data set

mean = median = mode

e.g.: Life span of human being

mode > median > mean

In this most of the scores occur at the higher end of distribution & few scores at lower end.



Co-Variance \rightarrow it helps to quantify the relationship b/w features, R+V
in a dataset

Covariance is used to find
the direction of the Relationship

e.g. Suppose I have a dataset

(x) Years ^{Age}	weight(kg)
20	75
18	63
15	45
14	40
25	78

As the $x \uparrow \Rightarrow y \uparrow$ from
dataset

$x \downarrow y \downarrow$

In some cases of dataset $\Rightarrow x \uparrow y \downarrow$
 $x \downarrow y \uparrow$.

In order to Quantify the relationship between X & Y,

(i.e features) then we use Covariance

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

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

↓
Becoz of
Bessel's correction

e.g. Economic Growth(Y)

(x)

2.1

2.5

4.0

3.6

Nifty 50 index(Y)

8.7

12

14

10

Age (X)

20

18

15

14

25

wt (kg)

75

63

45

40

78

covariance calculation

if covariance is +ve

$$\begin{cases} x \uparrow y \uparrow \\ x \downarrow y \downarrow \end{cases}$$

$$\begin{cases} x \uparrow y \downarrow \\ x \downarrow y \uparrow \end{cases}$$

X	y	\bar{x}	\bar{y}	$x-\bar{x}$	$y-\bar{y}$
2.1	8	3.1	11	-1	-3
2.5	12	3.1	11	-0.6	1
4.0	14	3.1	11	0.9	3
3.6	10	3.1	11	0.5	-1

$x \uparrow \Rightarrow y \uparrow \Rightarrow \text{cov} = +\text{ve}$
 $x \uparrow y \downarrow \Rightarrow \text{cov} = -\text{ve}$

$$\text{cov}(x,y) = \frac{(-1)(-3) + (-0.6)(1) + (0.9)(3) + (0.5)(-1)}{4-1}$$

$$= \frac{3 - 0.6 + 2.7 - 0.5}{3} = \frac{4.6}{3}$$

$\boxed{\text{cov}(x,y) = 1.533}$

+ve value (less than zero)

you are satisfying all conditions

① when economic growth $\uparrow \Rightarrow$ Nifty 50 index growth \uparrow

② " " " $\downarrow \Rightarrow$ " " "

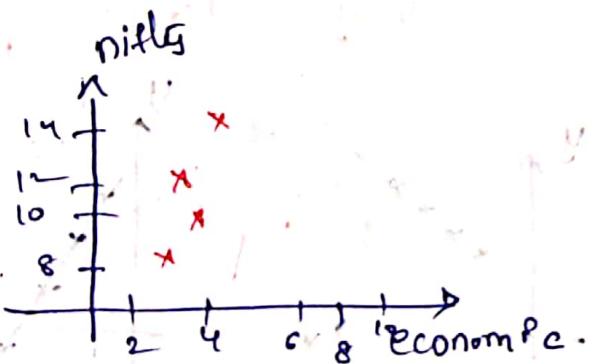
covariance can be positive / negative value, there is no restriction at all. but we should restrict. In order to prevent this we use Pearson & Spearman Rank correlation.

use Pearson Correlation coefficient (it is used for feature selection)

Pearson Correlation coefficient

Covariance is positive / negative but how strong it is +ve (or) how strong it is -ve, it is difficult to find (or) to restrict the value. In order to overcome we use Pearson Correlation Coefficient.

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



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

$$= \sqrt{(-1)^2 + (-0.6)^2 + (1.0)^2 + (0.5)^2}$$

4-1

$$= \sqrt{1.0 + 0.36 + 1.0 + 0.25} = \sqrt{2.58} = 2.58$$

= 0.8981

$$\sigma_y = \sqrt{\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}}$$

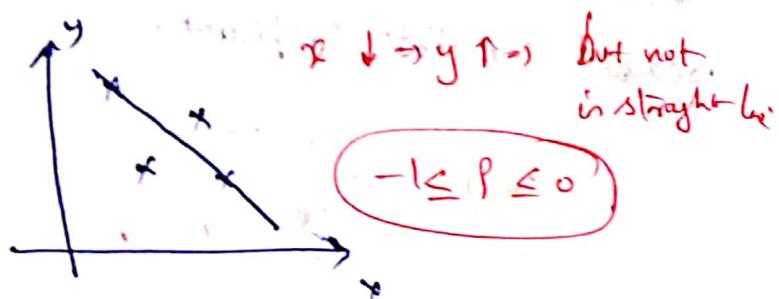
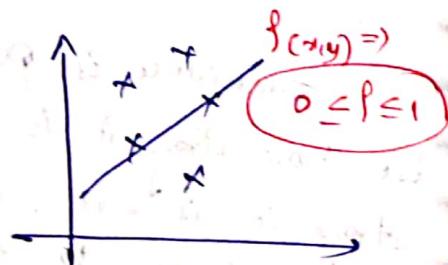
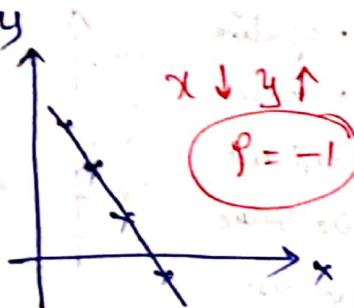
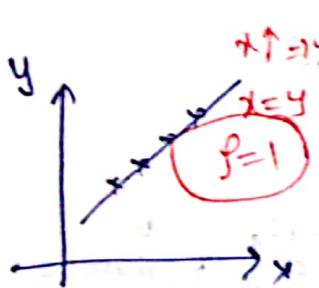
$$= \frac{(-3)^2 + (1)^2 + (2)^2 + (-1)^2}{3}$$

$$r_{(x,y)} = \frac{1.533}{(0.898)(2.58)} = 0.6611 \Rightarrow +ve \text{ correlation.}$$

\therefore x & y is truly correlated.

Based on Variance of x & y , Pearson Correlation Coefficient
will tell you about both Strength and Direction of Relationship

always $r_{(x,y)} = \boxed{-1 \leq r \leq 1}$



when we have 2 features x & x' , suppose if $1 \mapsto x' \mapsto$ then we say that $f=1 \Rightarrow$ 2 features are same \rightarrow so in this case we can drop one feature and apply ML models.

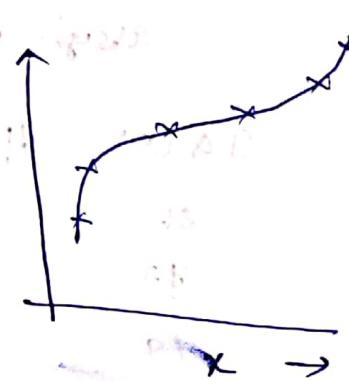
Spearman's rank correlation Coefficient

$$r_s = \frac{\text{Cov}(r_{gx}, r_{gy})}{\sigma_{rg_x} \cdot \sigma_{rg_y}}$$

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

$\text{Cov}(r_{gx}, r_{gy})$ = covariance of rank variables

$\sigma_{rg_x}, \sigma_{rg_y}$ = standard deviations of rank variables



eg The table below shows to calculate the correlation between IQ of a person with the no. of hours spent in front of TV per week

IQ (X) Hours b/Tv / week (Y)

106	7
86	2
101	50
99	28
103	29
97	20
113	12
112	6
110	17

Sol If we want to apply Spearman's rank correlation, steps are

Step 1 \Rightarrow Sort the data by the 1st col. (X), Create a new column X_p & assign the ranked values 1, 2, ...

2 \Rightarrow sort the data by 2nd col (y), Create a col. y_p & assign ranked values

$X_p (X)$	Hours of TV (Y)	Rank x_p	Rank y_p	d_i	d_i^2
81	0	1	1	0	0
97	20	2	6	-4	16
99	28	3	8	-5	25
100	27	4	7	-3	9
101	50	5	10	-5	25
103	89	6	9	-3	9
106	7	7	3	4	16
110	17	8	5	3	9
112	6	9	2	7	49
113	12	10	4	6	36

Step 3 \Rightarrow Create a 5th col. d_p to find the differences b/w 2 ranks (x_p & y_p) & square the d_i value.

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

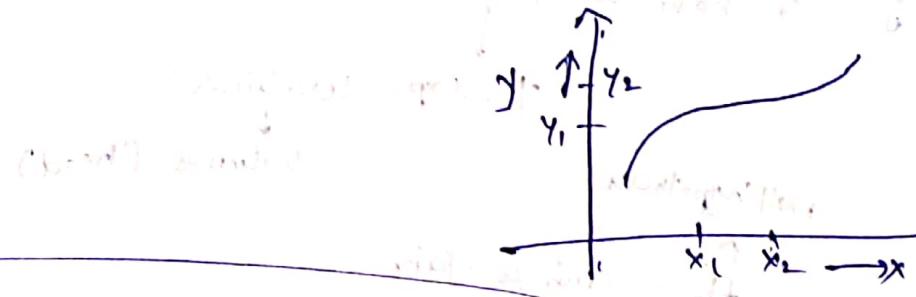
$$= 1 - \frac{6 \times 194}{10(10^2 - 1)}$$

$$r = -0.12575$$

Correlation is very very less.

\Rightarrow very near to zero \Rightarrow hours spent on TV is very less

Even though the data is non-linear, we can find the correlation / we can clearly do the correlation b/w 2 features by using Spearman's rank correlation.



e.g:

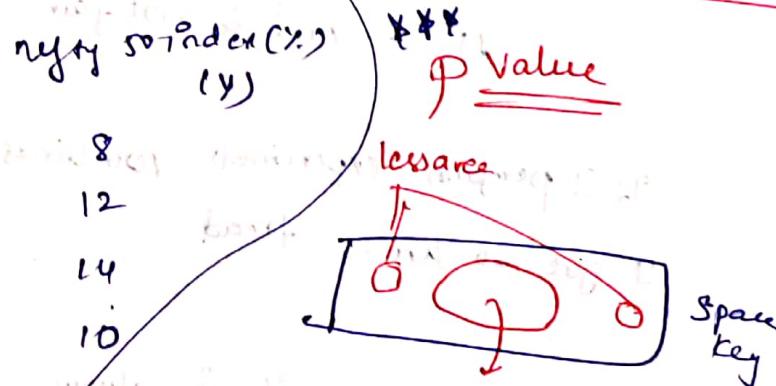
Economic growth (Y.)

	1	2	3	4	5
1. Economic growth	2.1	8	3.1		
2. Inflation rate	2.5	12	3.1		
3. Interest rates	4.0	14	3.1		
4. Unemployment	3.6	10	3.1		

Corr(x,y)

x	y	\bar{x}
2.1	8	3.1
2.5	12	3.1
4.0	14	3.1
3.6	10	3.1

Spearman's rank index (Y.)

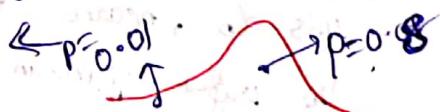


P Value

less than

Space key

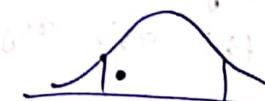
probability



where there is maximum no. of touchings, there we have a bulb part, where there is less no. of touchings, there it is tail part. Here we have 2 tailed test.

Let we touch at the

dotted point, -then let $p=0.01$



it indicates, if we repeat this experiment 100 times, the no. of times we will be touching this area is 1 time

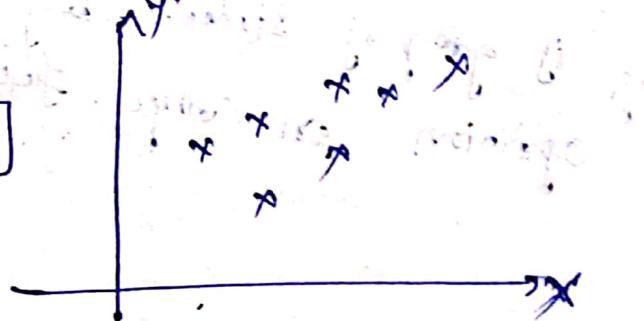
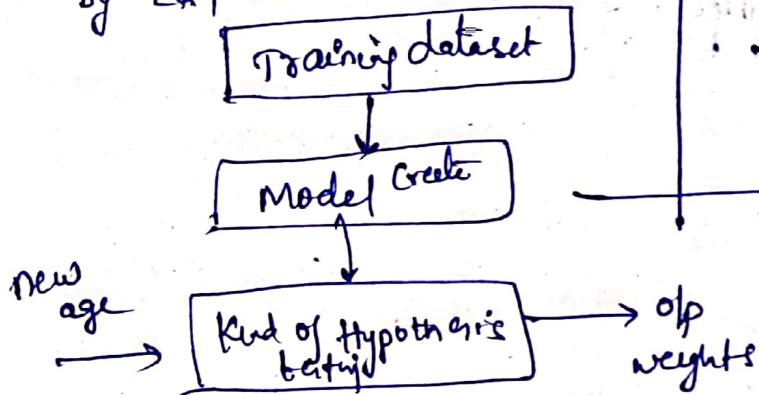
\rightarrow if $p=0.8 \Rightarrow$ it touches 80 times

p value \Rightarrow it is the probability for the null hypothesis to be true.

Linear Regression

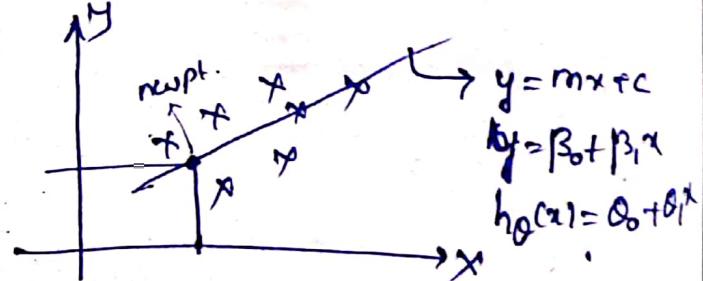
Suppose I have 2 features (Age, wt). based on these features we have data pts.

By LR, we



By LR, we will try to find out a best fit line, which will help us to do prediction

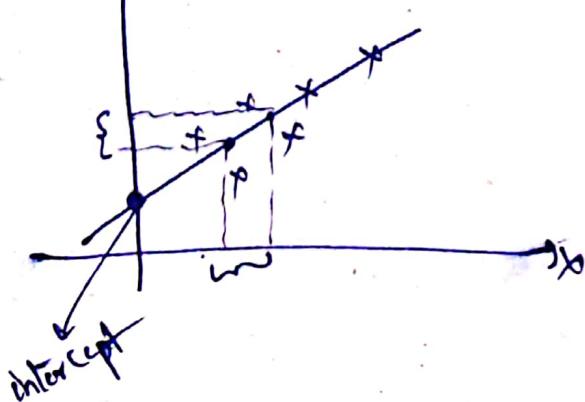
$$y = \text{linear function}(x)$$



In order to create a straight line, we use an equation

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

$$y = c + \theta_1 x$$



θ_0 = intercept

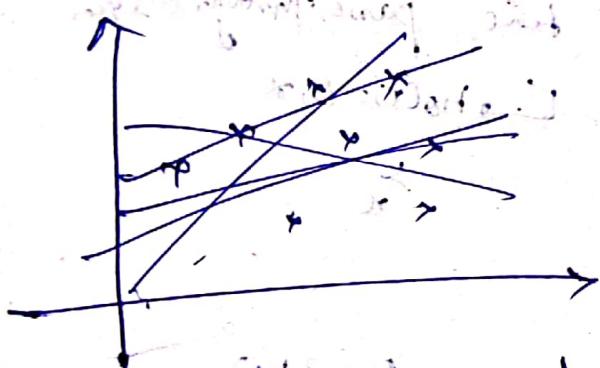
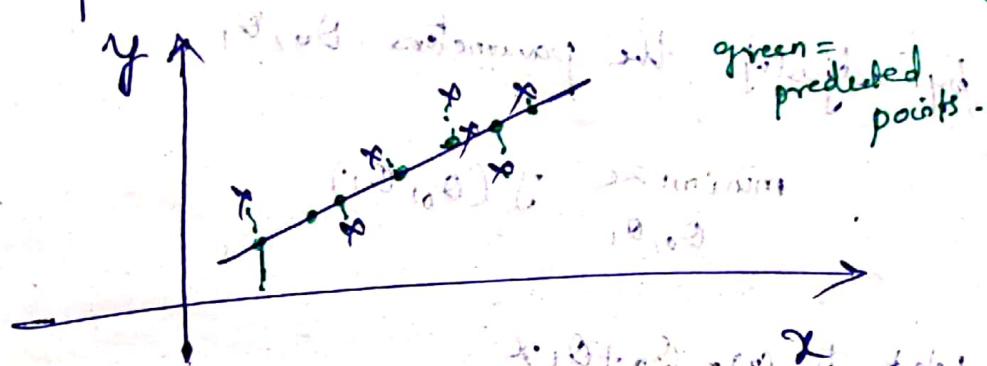
when $x=0 \Rightarrow h_{\theta}(x)=\theta_0$

at what pt, you are metting the y axis

θ_1 = Slope / Coefficient

in 1-unit moment
on x-axis, what is
the 1-unit moment is
in y-axis

x_p = data points
 Main aim \Rightarrow to create a best fit line in such a way that the distance b/w the data pts & have & predicted pts should be very less.



we should do many iterations and how do you know which line is best fit line?

for this we should start at one point & then go towards finding the best fit line.

Following this, we use a cost function.

Cost function

formulae for finding distance

$$C.f = \frac{1}{m} \sum_{i=1}^m (h_0(x_i) - y_i)^2$$

$$J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_0(x_i) - y_i)^2$$

entire eq called as squared error function

y = original pt

$h_0(x)$ = predicted pt.

\rightarrow remove -ve value

m = no. of data pts.

\sum = add all the distances.

$\frac{1}{m}$ = avg of all

$\frac{1}{m} \rightarrow$ for derivative purpose

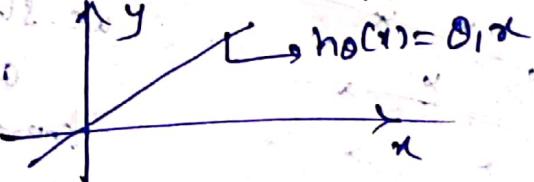
What we need to solve

I need to minimize $\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$
by adjusting the parameters θ_0, θ_1 .

minimize $J(\theta_0, \theta_1)$.
 θ_0, θ_1

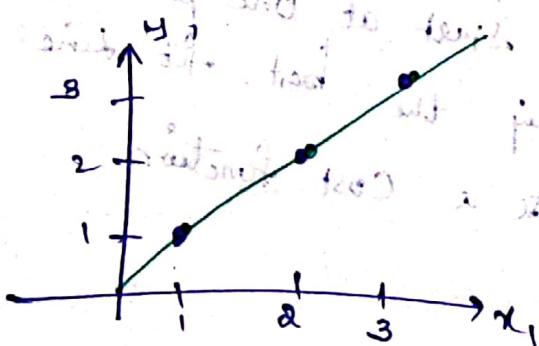
Let $h_\theta(x) = \theta_0 + \theta_1 x$

if $\theta_0 = 0 \Rightarrow$ the best fit line passes through origin



e.g.

① $h_\theta(x) = \theta_1 x, \theta_0 = 0, \theta_1 = 1$



$$D_1 = (1, 1)$$

$$D_2 = (2, 2)$$

$$D_3 = (3, 3)$$

If $\theta_1 = 1 \Rightarrow$ then

straight line passes through
all the pts.

When $\theta_1 = 1 \Rightarrow$ as predicted pts & original data pts are
same then the str. line passes through all the pts

Cost function $J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m [(h_\theta(x^{(i)}) - y^{(i)})^2]$

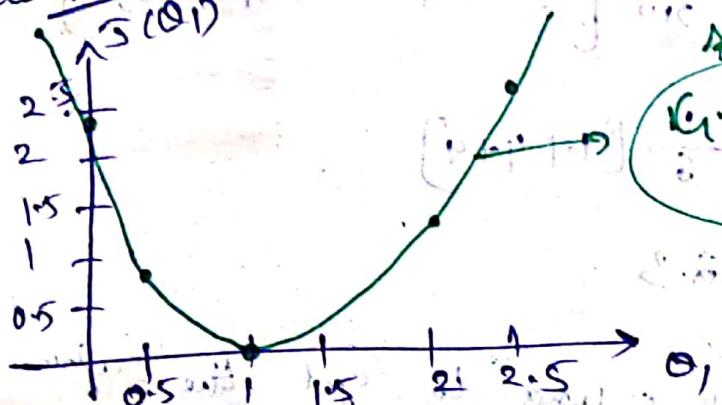
$$J(\theta_1) = \frac{1}{2m} \cdot [(1-1)^2 + (2-2)^2 + (3-3)^2]$$

$J(\theta_1) = 0 \rightarrow$ when $\theta_1 = 1$

$$\theta_1 = 1 \Rightarrow J(\theta_1) = 0$$

$$J(1) = 0$$

cost function graph



Gradient
Descent

② $\theta_1 = 0.5$

$$h_{\theta}(1) = 0.5(1) = 0.5$$

$$h_{\theta}(2) = 0.5(2) = 1$$

$$h_{\theta}(3) = 0.5(3) = 1.5$$

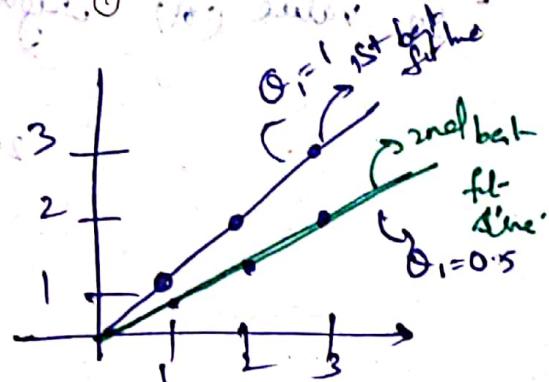
~~keep~~

$$J(\theta_1) = \frac{1}{2m} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2]$$

$$= \frac{1}{2(3)} [(0.5)^2 + 1 + 2.25]$$

$$J(\theta_1) \approx 0.58$$

$$J(0.5) = 0.58$$

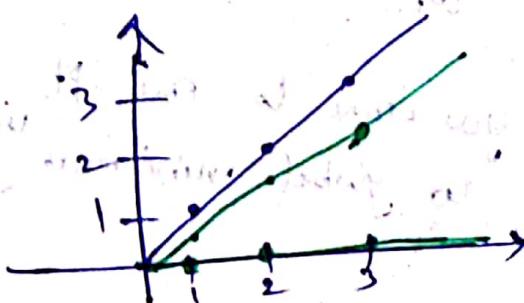


③ $\theta_1 = 0 = \cancel{\theta_1(0)}$

$$h_{\theta}(1) = 0$$

$$h_{\theta}(2) = 0$$

$$h_{\theta}(3) = 0$$

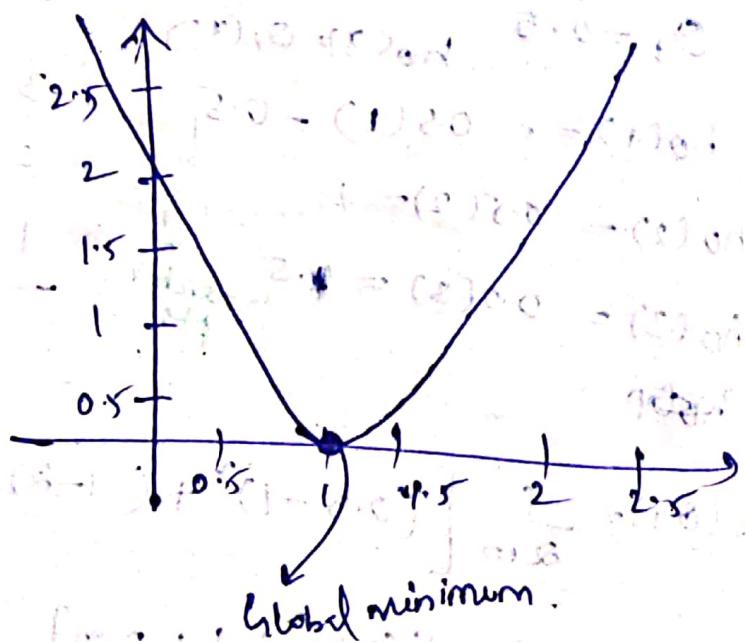


$$J(\theta_0) = J(0) = \frac{1}{2m} \left[(0-1)^2 + (0-2)^2 + (0-3)^2 \right]$$

$$= \frac{1}{6} [1+4+9]$$

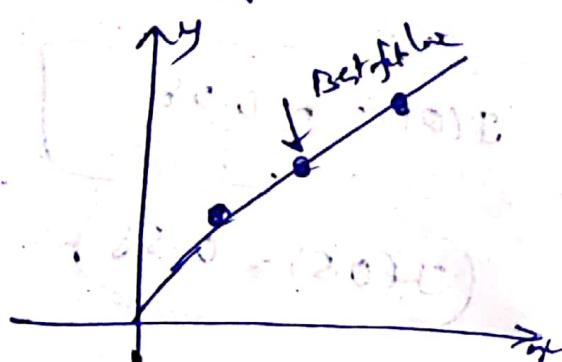
$$J(0) \approx 2.3$$

Gradient descent is used to get the right θ_0 value (or) right slope value.



we got 3 less, which is the best fit line

when θ_0 is at global minimum, the distance b/w actual pt & predicted pt is very very less.



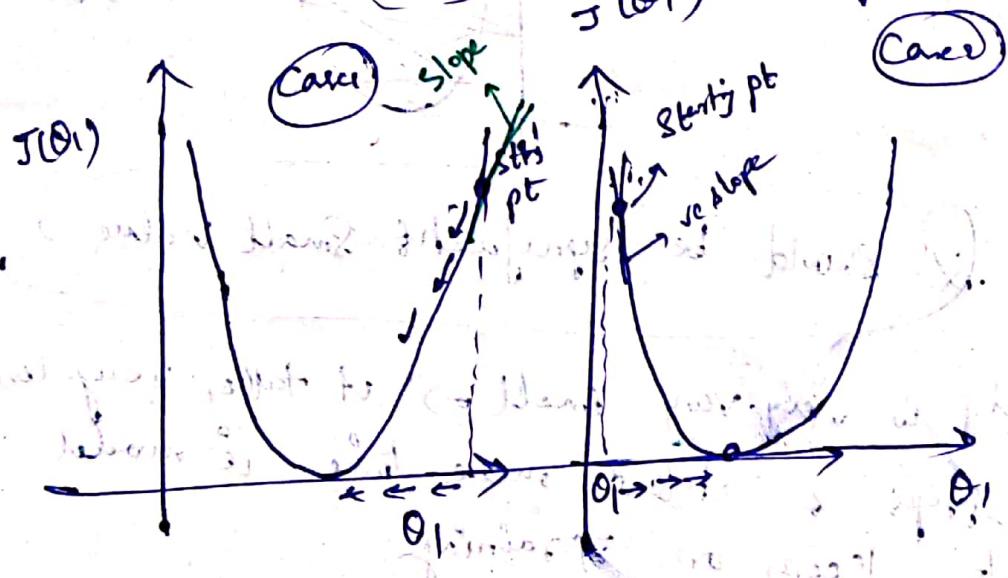
→ you come to one pt and then we reach to global minimum instead of using θ_0 value. So, how do you this?

for this we use Convergence Algorithm

Repeat until Convergence
Condition update

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j} \quad \text{derivative}$$

θ_j & slope.



after reaching to one pt, we will apply derivative i.e.
we are trying to find out the slope.

In the 1st case, at starting pt, we calculate
slope, it is +ve slope.
update the weights now we should

$$\theta_1 := \theta_1 - \alpha (+ve \text{ value})$$

as we have +ve slope.

$$\theta_1 := \theta_1 - \alpha (-ve \text{ value})$$

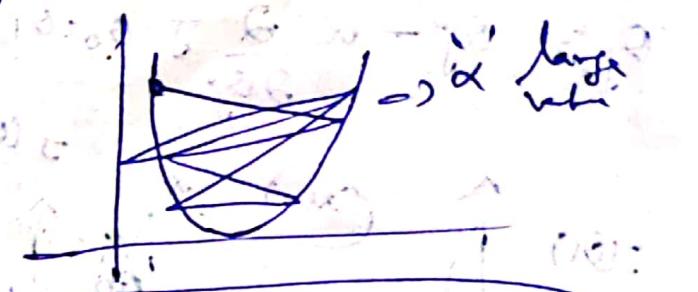
-ve slope
for corr

α = learning rate

by what speed we
should reach to global
minimum

If $\alpha = 0.01 \Rightarrow$ low value \Rightarrow it vis small steps to reach global min

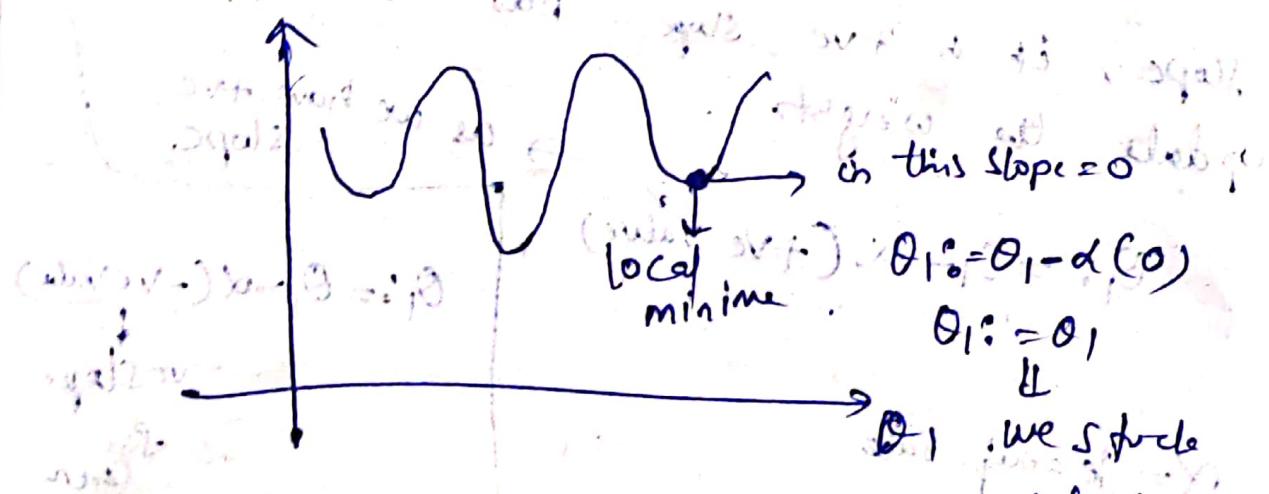
$\Rightarrow \alpha$ very large value \Rightarrow update of θ_1 keeps on jumping, it will never reach c.m.



Should be significant small value

α is very very small \Rightarrow it takes very less steps & takes more time i.e model gets keeps on training.

\rightarrow If my cost function has local minima.



Local minima is not in U.R, but is Deep local.

Gradient Descent Algorithm

repeat until convergence

{

$$\theta_j^* = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$

$$\frac{1}{2m} (x^2) - \frac{\mu x}{2m}$$

if $j=0$ =

$$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$= \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$j=1 \Rightarrow \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

$$h_\theta(x) = \theta_0 + \theta_1 x$$

repeat until convergence

{

$$\theta_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$= \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) (x^{(i)})$$