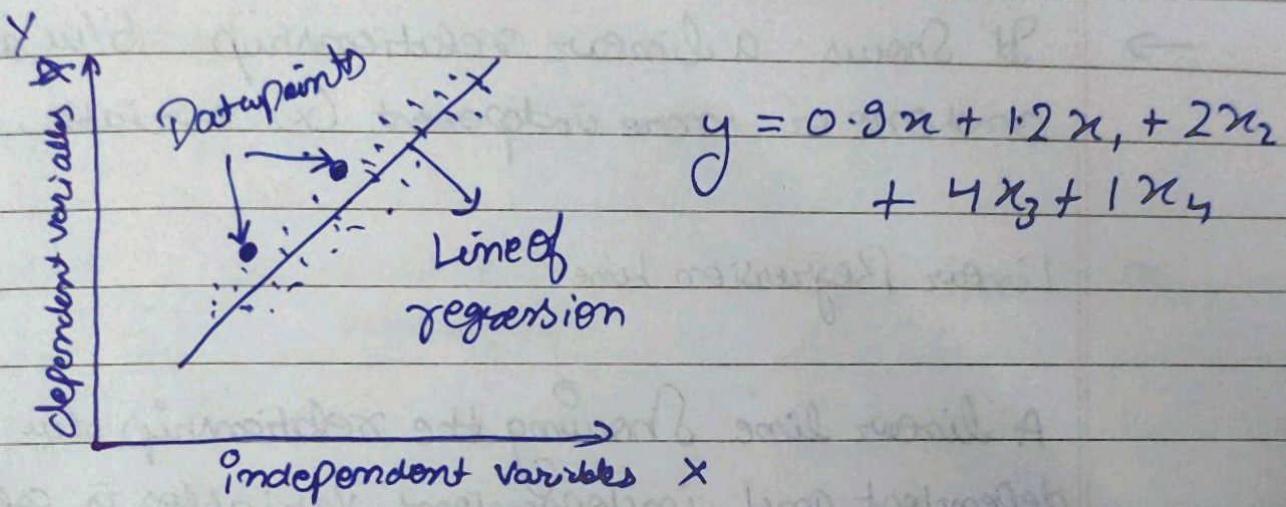


Unit-1Linear Regression

- Dependent Variable is continuous in nature.



Simple Linear equation

$$y = \alpha_0 + \alpha_1 x_1 \quad (y = c + mx)$$

Multiple Linear Equation

$$y = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_m x_m$$

α_i = Reg. Coeff.

x_i = Independent variable.

y = Dependent variable.

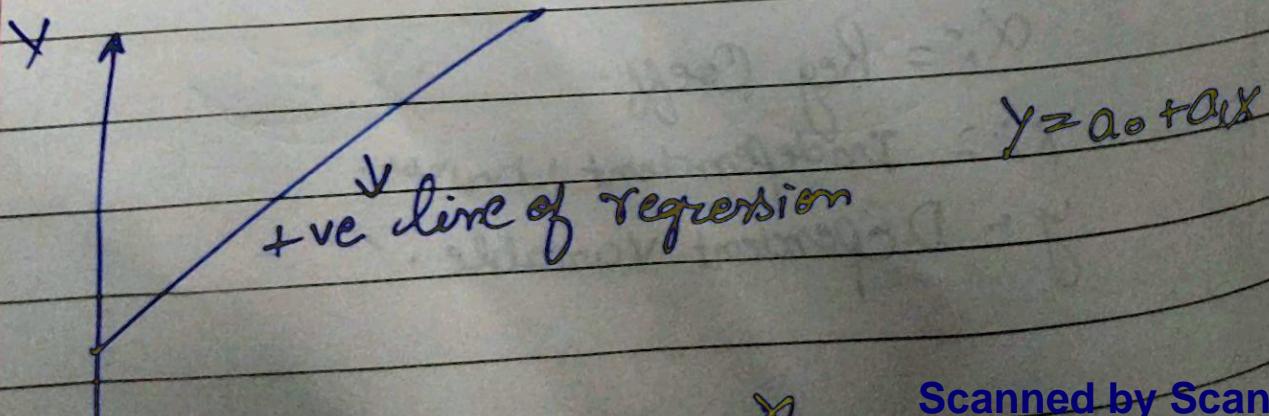
- Linear equation makes predictions for continuous real or numeric variables such as sales, salary, age, product price etc.
- It shows a linear relationship b/w a dependent (y) and one or more independent (x) variables, i.e.

Linear Regression Line

A linear line showing the relationship b/w the dependent and independent variables is called a regression line.

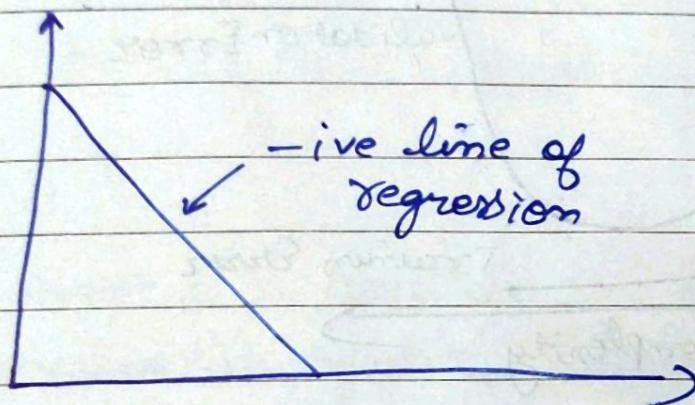
- Positive Linear Relationship

If the dependent variable increases on the Y axis and independent variable increases on X-axis, then such a relationship is termed as a positive linear relationship.



- Negative Linear Relationship

If the dependent Variable decreases on the Y-axis and independent Variable increases on the X-axis, then such a relationship is called a negative linear relationship.



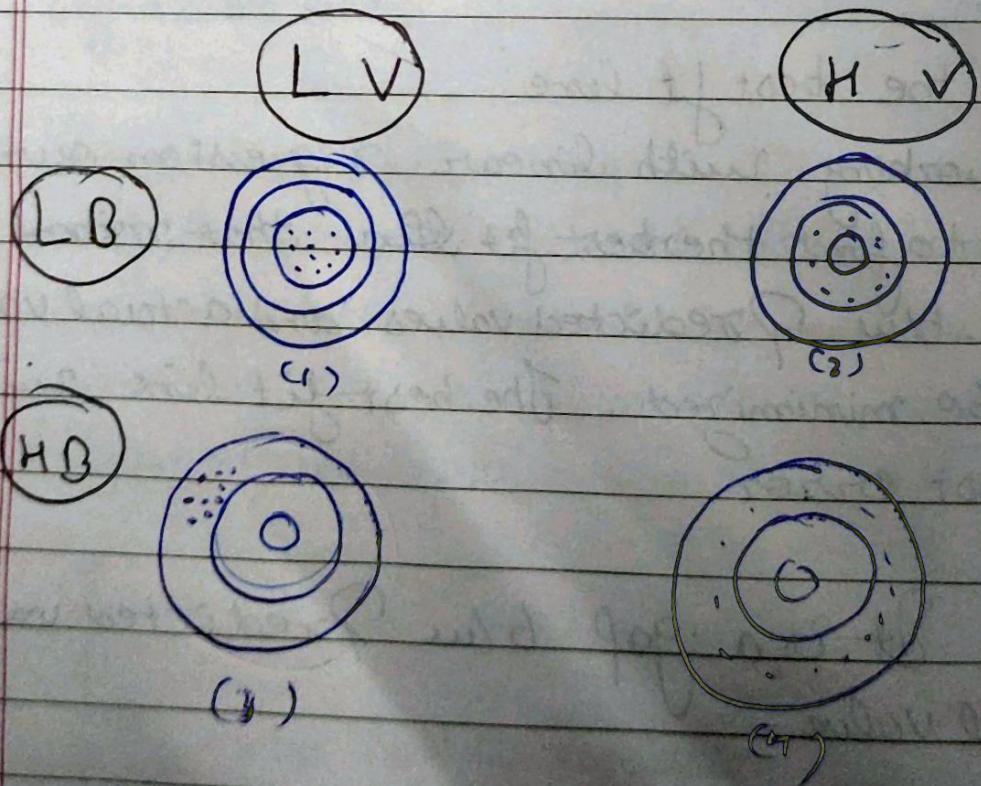
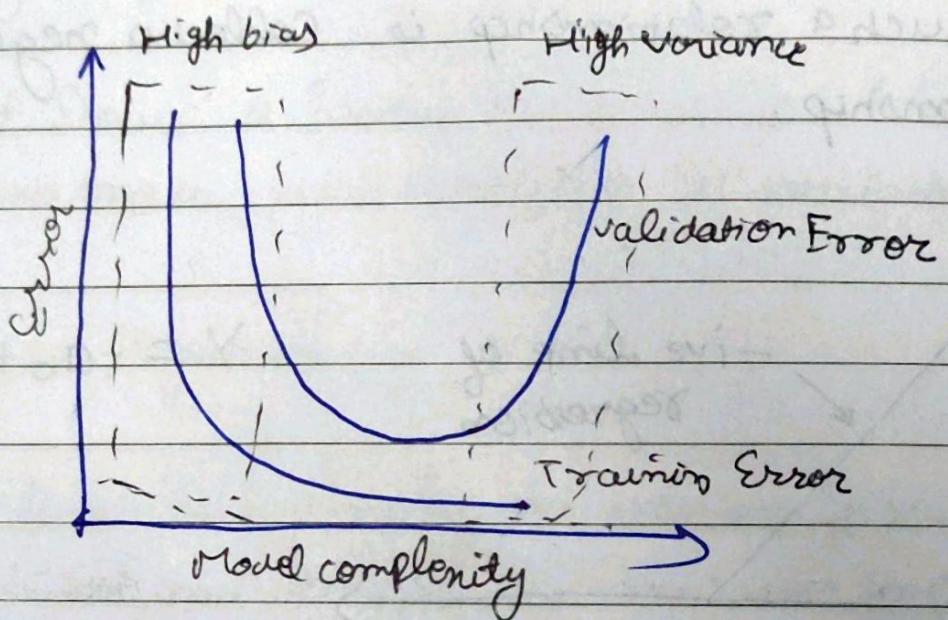
$$Y = -a_0 + a_1 x$$

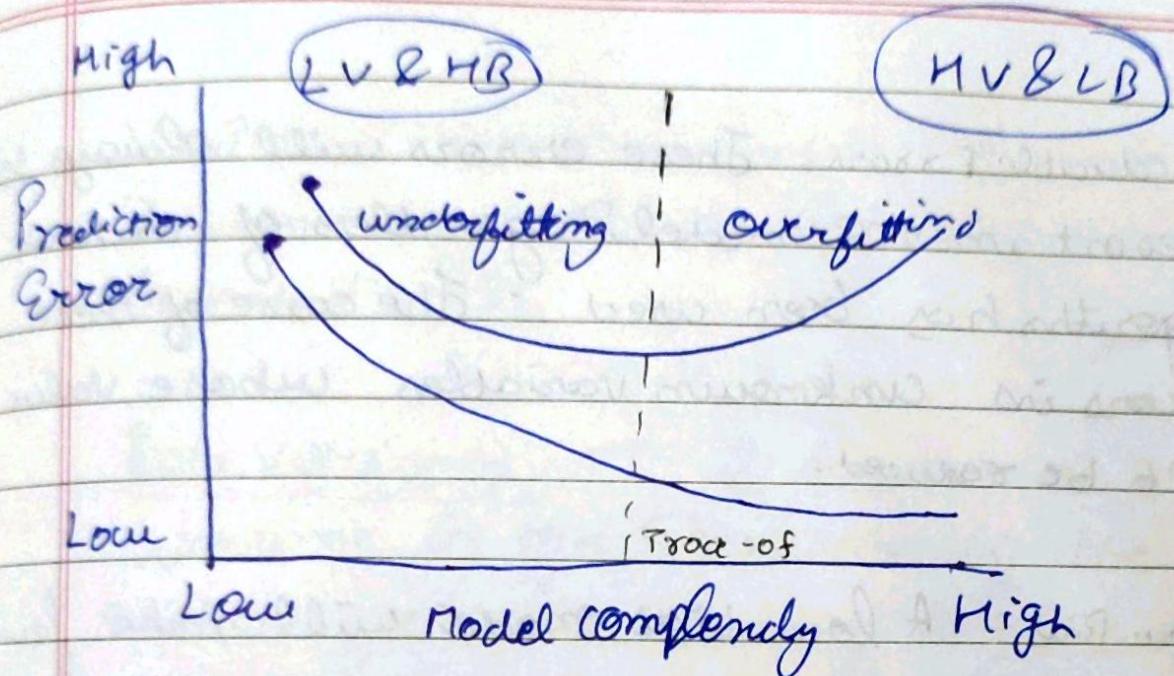
Finding the best fit line:

When working with linear regression our main goal is to find the best fit line that means the error b/w Predicted values and actual values should be minimized. The best fit line will have the least error.

Ans \rightarrow It is a gap b/w Predicted value & actual value.

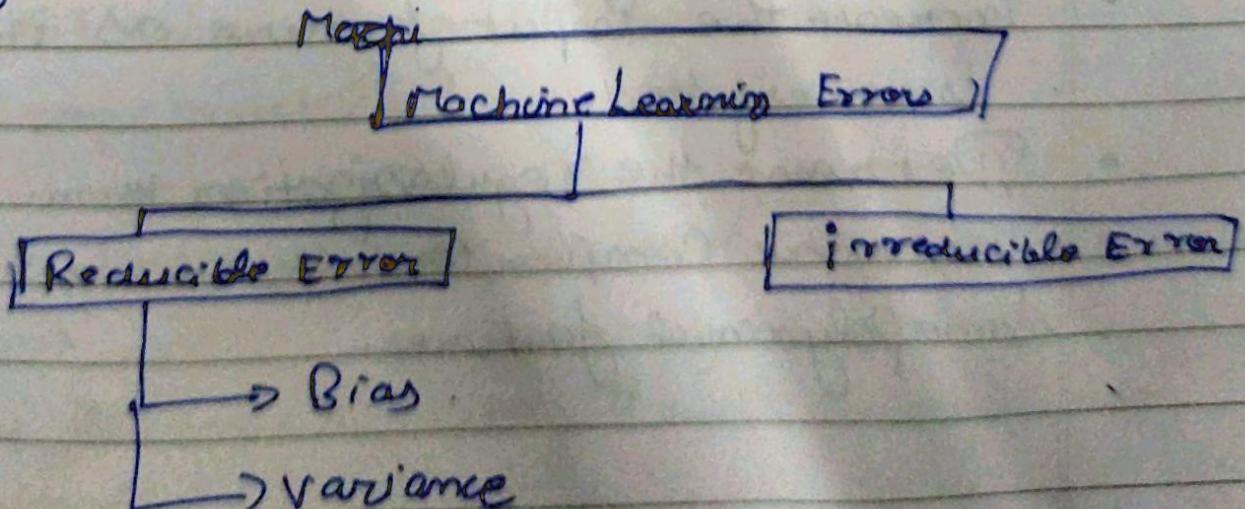
Bias is a Phenomenon that Skews the result of an algorithm in favour or against an idea.





An error is a measure of how accurately an algorithm can make predictions for the previously unknown dataset.

Reducible errors: These errors can be reduced to improve the model accuracy. Such errors can further be classified into bias & variance.



irreducible Errors: These errors will always be present in the model regardless of which algorithm has been used. The cause of these errors is unknown variables whose value can't be reduced.

Low Bias: A low bias model will make fewer assumptions about the form of the target function.

High Bias: A model with a high bias makes more assumptions, and the model becomes unable to capture the important features of our dataset.

Ways to reduce high bias:

- increase the input feature as the model is underfitted.
- Decrease the regularization term.
- use more complex models, such as including some polynomial features.

Variance:- Variance tells that how much a random variable is different from its expected value.

Low variance → It means there is a small variance in the prediction of the target func with changes in the training data set.

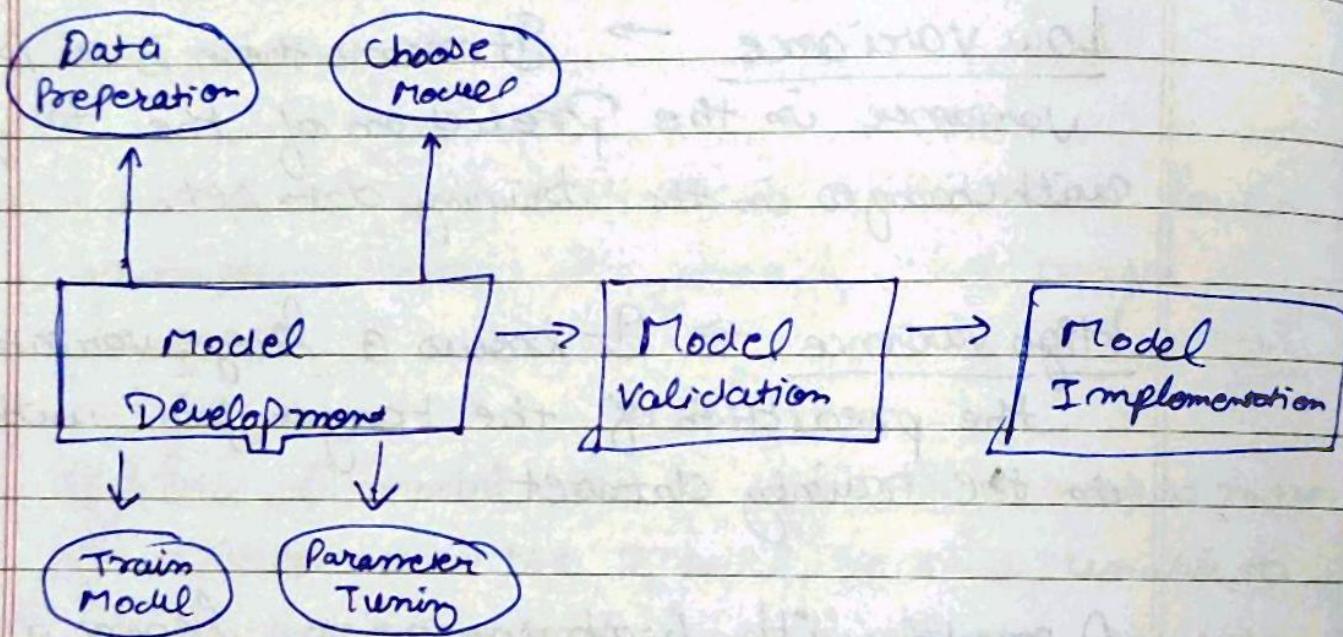
High variance → It shows a large variation in the prediction of the target func with changes in the training dataset.

A model with high variance learns a lot and perform well with the training dataset, and does not generalize well with the unseen dataset.

- A high variance model leads to overfitting
- Increase model complexities.
- If we decrease the variance, it will increase the bias
- If we decrease the bias, it will increase the variance.

Model validation

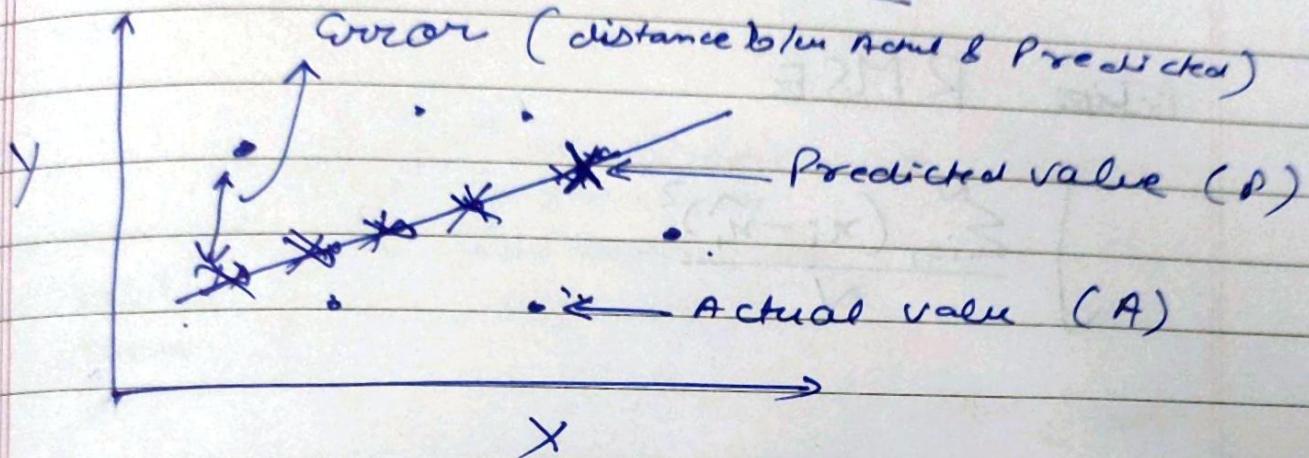
- A model Should Perform well in the new data



Model validation is the process that is Carried out after Model Training where the trained model is evaluated with a testing data set.

The testing data may or may not be a chunk of the same data set from which the training set is procured

Mean Squared Error



$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Model A

Actual	Predicted	Error	Error Square
10	9	1	1
13	14	-1	1
14	12	2	4
9	10	-1	1
17	15	2	4

ΣE^2

11

MSE

2.2

Model B

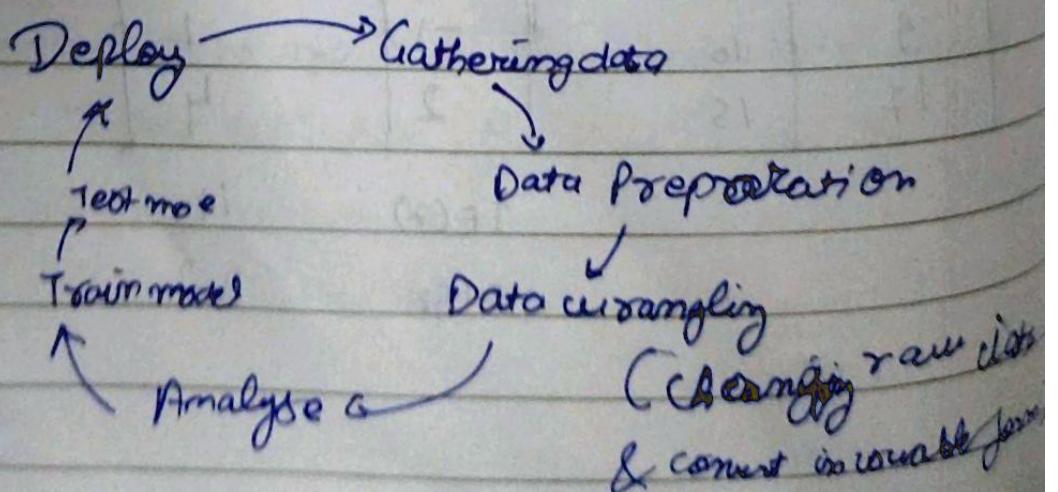
Actual RMSE

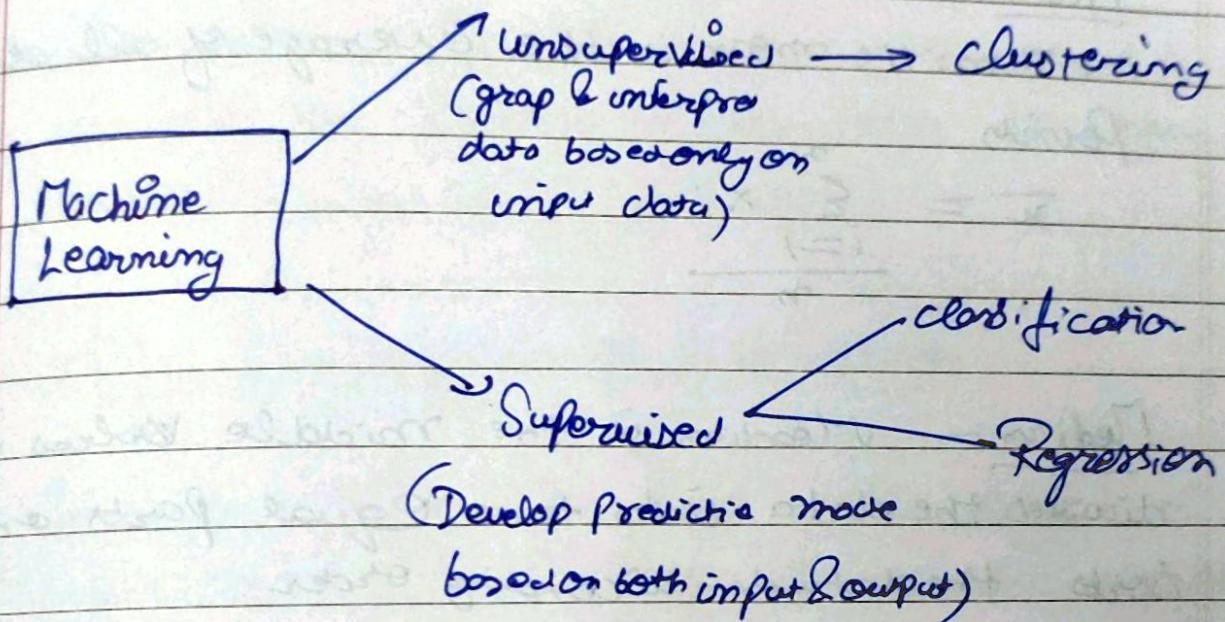
$$\sqrt{\frac{\sum_{i=1}^N (x_i - \hat{x}_i)^2}{N}}$$

Types of ML:-

- i) Supervised - labelled data for training
- ii) unsupervised - unlabelled data for training
- iii) Reinforcement - Algo discovered data through a process of trial & error then decides what action results in higher rewards

ML Life cycle:



ML techniquesRegression

Help Predictor interpret a Particular numerical Value based on prior data.

Classification

Predict or explain a class value.

Clustering → Aim to group or group Observation with similar characteristics

Statistics & Exploratory Data Analysis

- Mean -

Arithmetic mean is the average of all data points

points

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

- Median - Median is the middle value that divides the data into two equal parts and sorts the data in ascending order.

If total number of data points (n) is odd, the median is the value at position $(\frac{n+1}{2})$.

- Mode - Mode is value that occurs most frequently in the data set.

- Correlation - Explains how one or more variables are related to each other.

- i) Positive Correlation - When value of one variable increases then the value of other variable also increases.

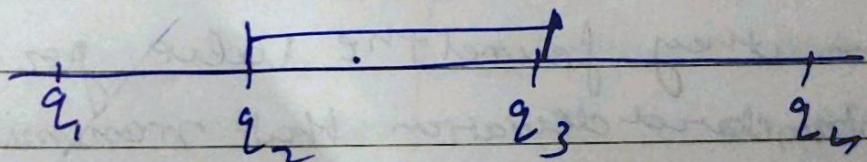
* ii) - Ve Correlation - when value of one variable increase the value of other variable decrease.

* Covariance - Is a measure for how two variables are related to each other i.e. how two variables vary with each other.

Diff. b/w covariance & correlation -

Covariance only tells about the directional relation b/w two variable but correlation along with direction also tells its strength of that relation.

* Quartile: Quartile are position indicators that divides a sequence into 4 equal parts.



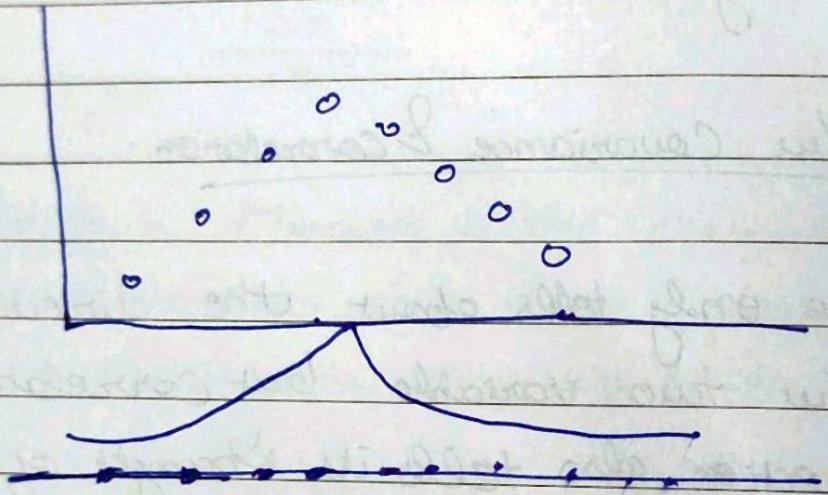
$IQR = \text{inter quartile range}$



Maximum likelihood

Goal of maximum likelihood is to find the optimal way to fit a distribution to the data.

Likelihood
of observing
the data



When someone says that they have the maximum likelihood estimates for the mean are the standard deviation or something else then they found the value for the mean or the standard deviation that maximizes the likelihood that you observed the things you observed.

Bayesian Inference

Bayesian Inference technique specify how one should update one's beliefs upon observing data.

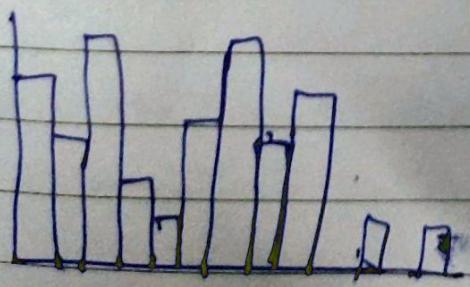
- Machine learning is mainly concerned with prediction and prediction is very much concerned with probability.
- The frequentist definition of probability is based on frequencies of events whereas the Bayesian definition of probability is based on our knowledge of events.

Ex we have 50 data

Sample x_1, x_2, \dots, x_{50}

i.e

Exponential distribution
 $(\lambda e^{-\lambda x})$



What is estimating?

Point estimation

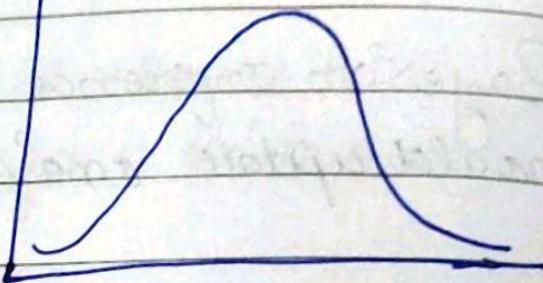
(Eg. - maximum likelihood estimation)

Single best estimation

$$\lambda = \frac{1}{(x_1 + x_2 + \dots + x_{50})/50}$$

Internal estimation

e.g. - Bayesian Inference

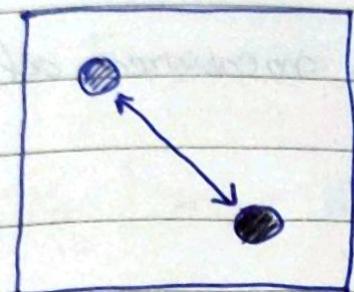


So,

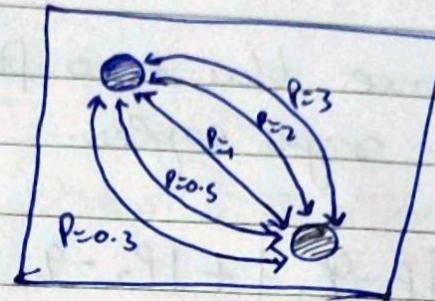
- Bayesian inference is an internal estimation
- It is a process of improving distribution function -
 $p_r(\lambda) \Rightarrow p_r(\lambda | \text{observed data})$

Distance Matrices

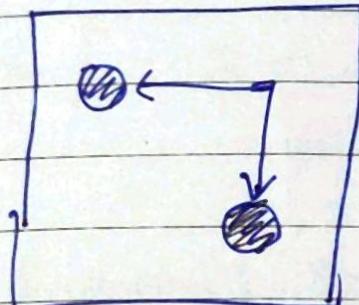
1. Euclidean Distance
2. Manhattan Distance
3. Minkowski Distance
4. Hamming Distance



Euclidean



Minkowski



	1	0	1	1	0	0
A	1	0	1	1	0	0
B	1	1	1	0	0	0

Hamming

⇒ Euclidean Distance

Euclidean Distance b/w these represents ^{Shortest} distance
b/w two points.

$$d = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

⇒ Manhattan Distance

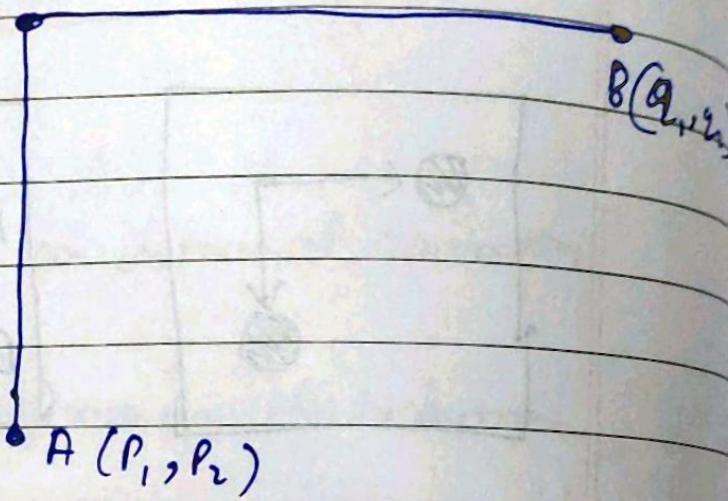
- It is the sum of absolute differences b/w points across all the dimensions.

The distance b/w two points measured along axes at right angles.

> $d = |P_1 - Q_1| + |P_2 - Q_2|$

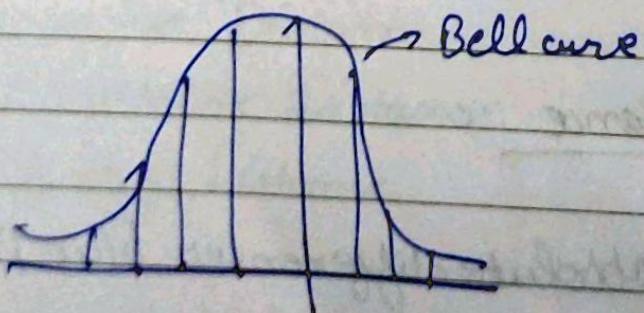
> And the generalized formula for an n -dimensional space is given as

$$D_n = \sum_{i=1}^n |P_i - q_i|$$



Normal / Gaussian Distribution

$$n \approx G.D(\mu, \sigma^2) \rightarrow \text{standard deviation}$$



Empirical Formula

$$\textcircled{1} \quad \Pr[\mu - \sigma \leq n \leq \mu + \sigma] \approx 68\%$$

$$\textcircled{2} \quad \Pr[\mu - 2\sigma \leq n \leq \mu + 2\sigma] = 95\%$$

$$\textcircled{3} \quad \Pr[\mu - 3\sigma \leq n \leq \mu + 3\sigma] \approx 99.7\%$$

Hypothesis Testing -

Evaluates 2 or more mutual exclusive statement on population using sample data

- Null hypothesis - It is regarding the assumption that there is no anomaly pattern or behavior according to the assumption made
- Alternate hypothesis - Contrary to the null hypothesis it shows that observation is the result of real effect.

Steps:-

- i) Make initial assumption (H_0 = null hypothesis)
- ii) Collect data (evidence)
- iii) Gather evidence to reject or not reject null hypothesis.

	H_0	H_1
Do not Reject	OK	Type 2 error
Reject	Type 1 error	OK

Type 1 error: False Positive (Reject H_0 when it is true)

Ex - The test says you have Covid but you actually don't

Type 2 error: False negative (Accept H_0 when it is false).

Ex - The test result says you don't have Covid but you actually do.

- * P - value - Significance value.
It tells us how likely our data could have occurred under the null hypothesis.
- Correctly it is taken as 5%
- If p value is 0.05 that means that 5% of the time you would see a test statistic at least as extreme as the one you found if the null hypothesis were true.

Dealing with missing values

- Ignore the tuple
- Filling Manually
- Using Global Constants "NULL" / "N/A"
- Use Attribute Mean / Median
- Mean Normal distribution
Median Skewed distribution (Asymmetric)
- Use of data mining Algorithm
- Inference tool, decision tree, clustering algorithm.

Regression Equation:-

- (g) From the given ~~equation~~ data, calculate regression equations taking deviation of items from the mean of X and Y series

X	1	2	3	4	5
Y	2	5	11	8	17

Theory Introduction to Regression Equation:-

- i) Regression eqⁿ of Y on X
- ii) Regression eqⁿ of X on Y

Y on X → Variation in Y when the change in X is given.

X on Y → Variation in X when the change in Y is given.

ii) Regression equation of X on Y

$$x - \bar{x} = \frac{\sigma_x}{\sigma_y} (y - \bar{y})$$

→ regression coeff

$$b_{xy} = \frac{\sigma_x}{\sigma_y} = \frac{\sum ny}{\sum y^2}$$

iii) Regression eqn of y on x

$$y - \bar{y} = \frac{\sigma_x}{\sigma_y} (x - \bar{x})$$

$$b_{yx} = \frac{\sigma_y}{\sigma_x} = \frac{\sum ny}{\sum x^2}$$

γ = correlation coeff

$$\gamma = \sqrt{b_{xy} \cdot b_{yx}}$$

<u>x</u>	$x = (x - \bar{x})$	x^2	y	$y = (y - \bar{y})$	y^2	xy
1	-2	4	2	-6	36	12
2	-1	1	5	-3	9	3
3	0	0	4	3	9	0
4	1	1	8	0	6	0
5	2	4	14	6	36	12
$\sum x = 15$.	$\sum x^2 = 10$	$\sum y = 40$		$\sum y^2 = 27$	$\sum xy = 27$

$$\bar{x} = \frac{\sum x}{N} = \frac{15}{5} = 3$$

$$\bar{y} = \frac{\sum y}{N} = \frac{40}{5} = 8$$

1) Regression equation of X on Y

$$(x - \bar{x}) = \frac{\sigma_x}{\sigma_y} (y - \bar{y})$$

$$(x - \bar{x}) = \frac{\sum xy}{\sum y^2} (y - \bar{y})$$

$$(x - 3) = \frac{27}{90} (y - 8)$$

$$(x - 3) = 0.3 (y - 8)$$

$$x - 3 = 0.3 y - 2.4$$

$$x = 0.3 y + 0.6$$

2) Regression equation of Y on X

$$(y - \bar{y}) = \frac{\sigma_y}{\sigma_x} (x - \bar{x})$$

$$(y - \bar{y}) = \frac{\sum xy}{\sum x^2} (x - \bar{x})$$

$$(y - 8) = \frac{27}{10} (x - 3)$$

$$y - 8 = 2.7 x - 8.1$$

$$y = 2.7 x - 0.1$$



Tadip

$$r = \frac{\sum b_{xy} \cdot b_{yx}}{\sqrt{\sum y^2 \sum x^2}} = \frac{\sum xy \times \sum xy}{\sqrt{\sum y^2 \sum x^2}} = \frac{\sum xy^2}{\sqrt{\sum y^2 \sum x^2}}$$

T-test

A t-test is a type of inferential statistic used to determine the significant diff b/w the means of two groups which may be related to certain features.

→ A t-test is only used when comparing the means of two groups also known as a pairwise comparison

Assumption

1. A t-test assumes your data are independent.
2. Assumes your data are normally distributed.
3. Assumes data have a similar amount of variance within each group.

Type of t-test

i) One sample t-test

we compare the average of one group against the set average.

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

t = t-statistic

\bar{x} = mean

μ = population mean

s = standard deviation

n = group size

ii) Unpaired or independent t-test -

Compare the means of two different groups of sample

$$t = \frac{\bar{x}_A - \bar{x}_B}{\sqrt{\frac{s_A^2}{n_A} + \frac{s_B^2}{n_B}}}$$

\bar{x}_A & \bar{x}_B = means of two diff groups

n_A & n_B = sample sizes

$$S^2 = \frac{\sum (n - m_A)^2 + \sum (n - m_B)^2}{n_A + n_B - 2}$$

Degree of freedom

Refers to the values in a study that can vary and are essential for assessing the null hypothesis importance & validity.

Paired t-test

We measure one group at two different times.
 We compare different means for a group at two different times or under two different conditions.

$$t = \frac{m}{\frac{s}{\sqrt{n}}}$$

Chi-square Test:-

This helps to determine whether there is a notable difference b/w the normal frequency and observed freq in one or more classes or categories.

Properties

- Two times the number of degree of freedom is equal to the Variance.
- The number of degree of freedom is equal to mean distribution.
- Chi-square distribution curve approaches the normal distribution when the degree of freedom increases.

Formula

$$\chi^2 = \sum \frac{(O \text{ (Observed Value)} - E \text{ (Expected value)})^2}{E}$$

ANOVA → Analysis of variance test to generate the T-tests for multiple groups.

Assumptions-

- 1) Data for every level of the factor is distributed generally.

2. Case independent - The sample cases must be independent from each other.
3. Variance b/w the group needs to be around equal.

Types :-

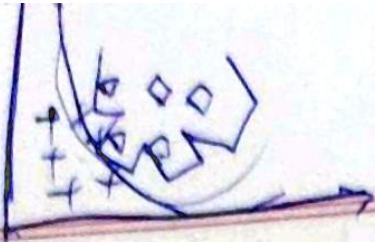
- i) One-way → Analysis of Variance Test which has only one independent variable
- ii) Two-way → ANOVA that has two independent variables. Also known as factorial anova test.
- iii) N-way - Use more than two independent variables.

Missing Value Treatment

Mainly missing values can be treated in two ways.

- i) Deletion -

To either delete the entire row or column in which missing value is present.
(Not recommended).



ii) Imputing

- Replacing with an arbitrary value
- Replacing with mean
- Replacing with mode
- Replacing with median
- Replacing with previous values

For Categorical data imput most frequent value.

Outlier detection

- Whether user - labeled examples of outliers can be obtained
 - Supervised, Semi-Supervised and unsupervised methods
- Assumptions about normal data and outliers
 - Statistical, proximity-based and clustering-based methods
- Among the data objects, on which does not obey general behaviour.

* Outlier Detection

Process of detecting outliers and Subsequently removing them

1. Statistical Approach :

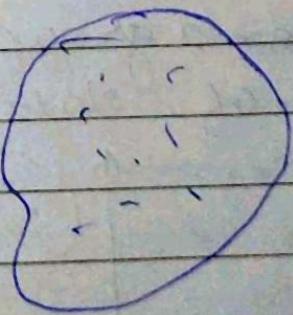
based on probability of the data points low probability = outliers

- Parametric methods
- Non Parametric methods

2 Proximity Approach:

based on location of data points

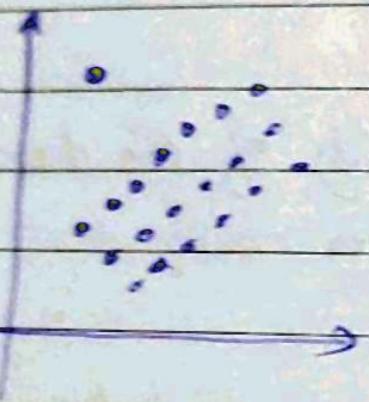
- Density based approach
- Distance " "
- Grid " "
- Deviation " "



* Types of outliers.

① types of outliers

1. Global / Pair outliers

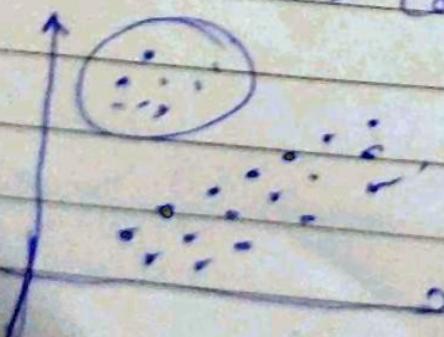


When a single data object deviates from the rest of data.

Point → Global / Pair outliers

2. Collective outliers:

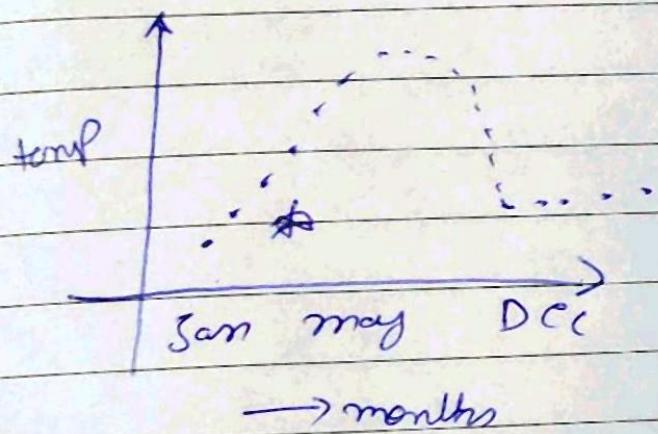
When a group of data points deviates from the rest of data called collective outliers.



3.

Contourual / conditional outliers:

Data objects deviates from others because of any specific condition - Called contour outliers.



Feature engineering

It is the process of selecting, manipulating, and transforming raw data features that can be used in supervised learning.

It is a machine learning technique that leverages data to create new variables that aren't in the training set. It can produce new features for both supervised and unsupervised learning → with the aim of simplifying and

Speeding up data transformations while also enhancing model accuracy.

Feature creation : creating features in values
Creating new variables which will be most helpful