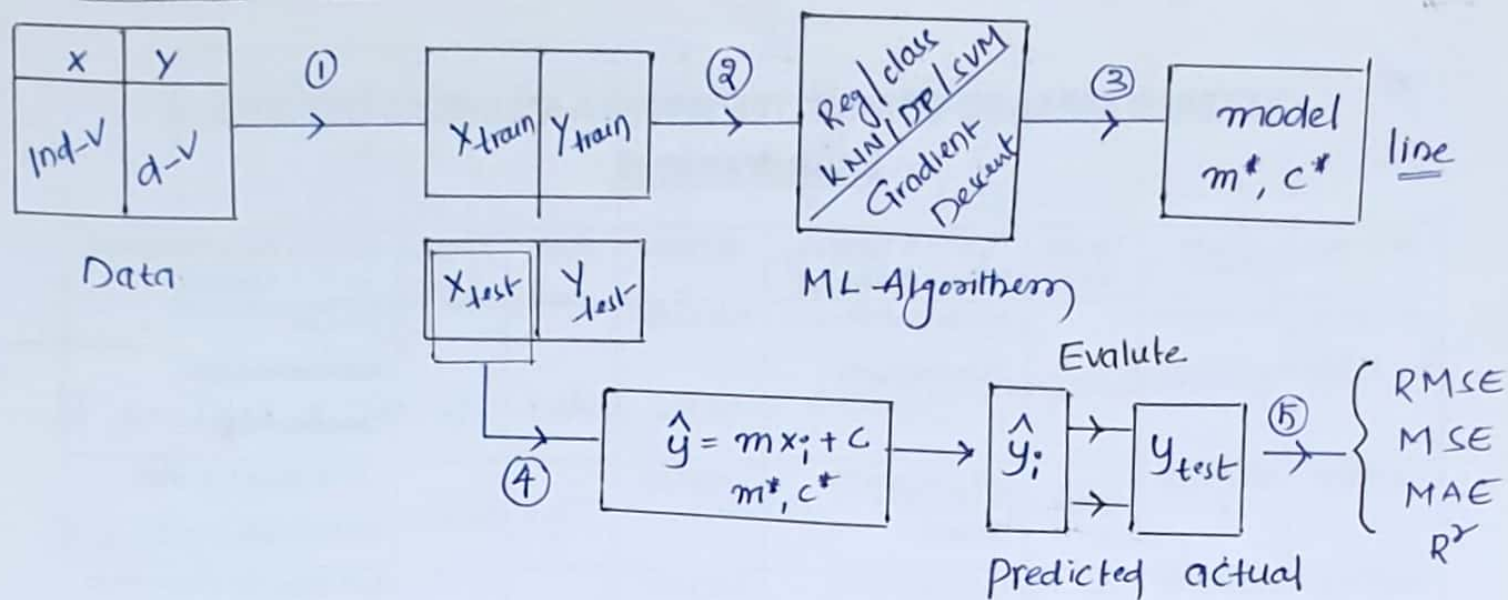


Graphical Representation :-



Simple linear Regression (S.L.R) :-

It is a Statistical method that allows us to Summarize and Study Relationships b/w two Continuous variables. Simply if we have one Independent variable available in data we called as Simple linear Regression.

Eqn - $y = mx + c$
 dep-var \rightarrow Ind-var

Multiple linear Regression :- It is a Statistical technique that uses several explanatory variables to predict the ~~model~~ outcome of response variable. Simply we can say, if we have more than one Independent variables available in data we called as multiple linear regression.

The goal of MLR is to model the linear relationship b/w Independent

-l variables and dependent variables.

$$y = m_1x_1 + m_2x_2 + c \rightarrow \text{slope} \quad \text{Intercept}$$

$$y = m_1x_1 + m_2x_2 + m_3x_3 + m_4x_4 + c$$

\rightarrow features (x_1, x_2, x_3, x_4)

In MLR we have some problems that is multicollinearity

x_1	x_2	x_3	x_4	x_5	x_6	y
		2	4			
		4	16			
		5	25			
		7	49			

$$X = (x_1, x_2, x_3, x_4, x_5, x_6)$$

'X' is Independent variable

'y' is dependent variable

if we have a dataset which contain multiple Independent Variables & One dependent Variable (y), if Independent variables any two variable Correlate to each other i.e. depend on other variable let's take (x_3, x_4) Variable, here x_4 is depend on x_3 variable (i.e. Square of it) then there is facing problem called as "multicollinearity"

In this we have to do ~~Bar~~ check which variables are dependent and resolve it, make all the variables to Independent, otherwise we are face problem called as "Interpretable model". i.e. we lost the Interpretability of a model.

Implementation of Simple linear Regression :-

- ① Import libraries
- ② load the dataset
- ③ Exploratory Data Analysis [find the linear Relationship should exist]
- ④ Split the data into train & test
- ⑤ Do Normalization (or) Standardization
- ⑥ train the model Using X_{train} , Y_{train} (i.e. Apply / fit Algorithm)

while train data we have to take some assumptions of linear Reg.

- (a) linear relationship should exist between input & output
 - (b) Errors / Residuals on training data should follow Normal distribution
ie mean is zero and variance is anything E of train data $\sim N(0, \sigma^2)$
 - (c) Errors / Residuals on training data (x_{train}, y_{train}) expected to be Independent of each other
 - (d) Errors / Residuals on training data expected to be "HOMOSCEDASTICITY"
- * Residuals also called as Errors

NOTE :- linear Regression task can be solved by two approaches -

- (1) Geometric way approach
- (2) probability way approach

In modern days we are using Geometric way approach to solve linear reg. task

then we use only assumption (a) if follows or not

Statisticians generally using probability way of approach to solve linear R. task
in this case we have to check if follows 4 assumption or not

- (7) predict on x_{test} - Using the model $(m^*, c^*) \rightarrow y_{test-pred}$

- (8) Evaluate the model - y_{test} (vs) $y_{test-pred}$ by using matrices we

calculate RMSE, MSE, MAE etc.

NOTE :- Accuracy in only Classification task
NO Accuracy for Regression problems.

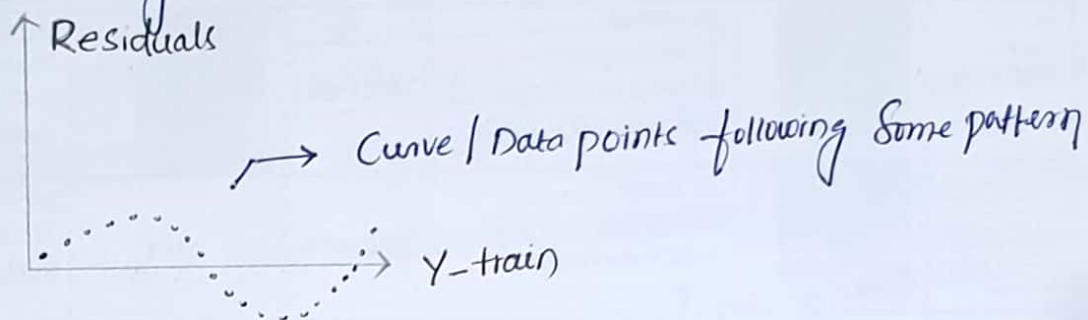
Assumption (b, c, d) can also known as Residual Analysis.

Residual Analysis :- it is nothing but assumptions of linear

Regression (b, c, d) brief explanation.

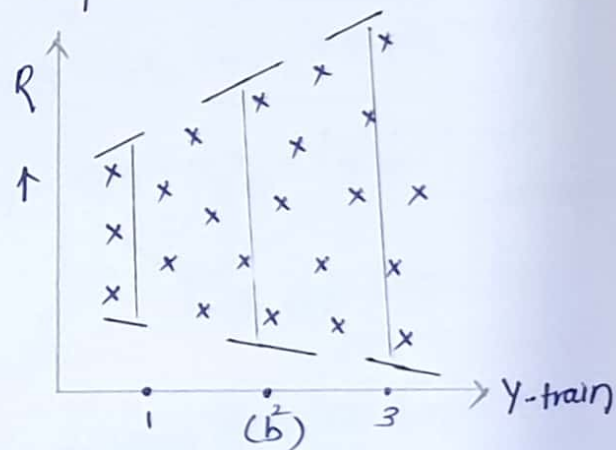
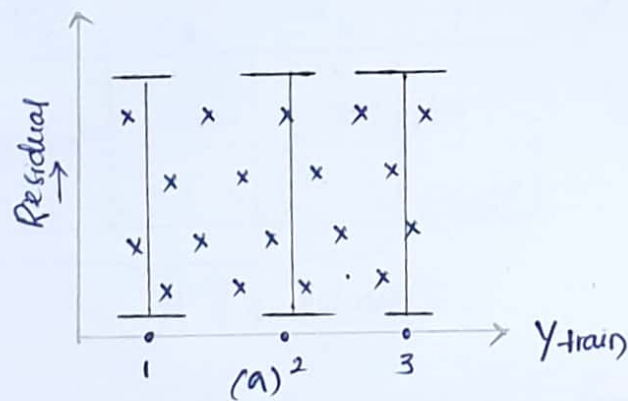
Ass (b) tells us residual of training data, when plot follow normal dist
ie bell shape Curve with '0' mean.

Ass (c) tells us training data expected to be Independent of each other.



By the above fig. shown as the train data is not Independent when we plot graph on (Residuals - Y-train) it follow some pattern. (dependent)

Ass (d) tell us residuals on training data expected to be "homoscedasticity"



If we observe figure (a) - let's take point '1' on Y-train the variance or spread at '1' point same as '2' point as well '3' point also if we

observe the fig (a) maintain constant variance / spread at the end.

this is called as "HOMOSCEDASTICITY", if the plots are not constant

variance [fig (b) shown] is called as "HETROSCEASTICITY"

Simple linear Regression :- it is a regression model that estimates

the relationship between One Independent Variable and One dependent Variable. Using straight line, both variables should be Quantitative.

X	Height	Weight	Y
	Ind-var	Dep-var	

$\Rightarrow Y = mx + c$
 \swarrow Intercept
 \searrow Slope
 $\rightarrow \in \mathbb{R}$
 $w_1x_1 + w_2x_2 + w_0 = 0$

Multiple linear Regression :- it is a statistical technique that estimates the relationship b/w two or more explanatory variables (Independent) and a response variable (Dep-var) both variables should be Quantitative.

Multiple regression is the extension of Ordinary least-Squares (OLS) Regression because it involves more than one explanatory variable.

X_1	X_2	X_3	Y
Ind-var	Ind-var	Ind-var	Dep-var

$\rightarrow \in \mathbb{R}$
 $w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 = 0$
 $y = m_1x_1 + m_2x_2 + m_3x_3 + c = 0$
 \swarrow Intercept
 \searrow 3-diff Slopes
 X (3-Ind-var)

Note :- ① All the four Residual analysis (Assumptions) hold on Simple linear.

Regression as well as Multiple linear Regression.

② The Independent variable (X) must not be exist- "multicollinearity" if multicollinearity exist then we have to drop one of the column.

We solve this regression problems by ① Geometric way
 ② by using prob & Statistic way.

Recursive feature Elimination (RFE) :-

if we have a data set which consist of multiple Independent variables and One dependent variable. if there are multiple variables then there high chances to get multicollinearity problem. (One column maintain some relation with other column or two are more columns derive other column).

if there is multicollinearity problem then we loose the better Interpretability. So In Order to get best interpretability we have to select features that are most important to our model and remove unwanted features.

In this process there are multiple different ways to select the features, in this the basic method is to "Recursive feature elimination".

there are different kind of feature elimination methods available.

they are \rightarrow (1) Automatic way \rightarrow R.F.E

(2) Manual way \rightarrow (1) forward feature Selection.

(2) Backward feature Selection.

(3) Mixed way approach \rightarrow In this we do both Automatic and manual way approach.

VIF (Variance Inflation factor) :- VIF is a tool to help identify the

degree of multicollinearity in a set of multiple regression variables.

$$VIF = \frac{\text{Overall model Variance}}{\text{Variance of a Single Independent variable.}}$$

- Assumptions of linear Regression :-

if we have a data that contain some i/p variable (X) and o/p variable (Y) and if the o/p variable is continuous ($\mathbb{R} \rightarrow \text{Real Number}$)

Then we have to perform regression task

Procedure to Complete task :-

Step-1 We have to split the data into trained data, test data.

Step-2 take trained data and passed to machine learning algorithm (KNN, linear Reg., logistic Reg., DT, SVM, etc) and it gives us to a model

Step-3 if we use linear Reg we get model as (m^*, c^*) which is a line
in step-2 we have to use GD (Gradient descent) to optimize model

Step-4 we have to do predictions on test data (x_{test}) that gives

$$y_{\text{pred values}} (\hat{y}) = mx_i + c$$

Step-5 finally we have to evaluate the y -pred values with y_{act} values
here ($y_{\text{act}} = y_{\text{test}}$) in this we use to matrices to calculate performance of model, here some matrices are used for regression task are -

RMSE, MSE, MAE, R^2 .