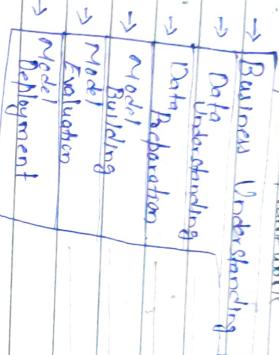


## Machine Learning - I

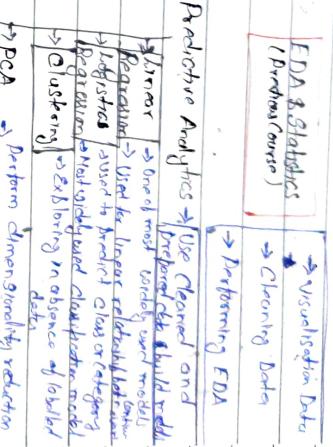
M	T	W	T	F	S	S
Page No:						YOUVA
Date: 23/12/21						

### Career-Dev Framework



### Pre-requisites

1. Exponentials
2. Logarithms
3. Equation of a straight line
4. Basis Transformation
5. Eigen decomposition
6. Inferential Statistics
7. Hypothesis Testing

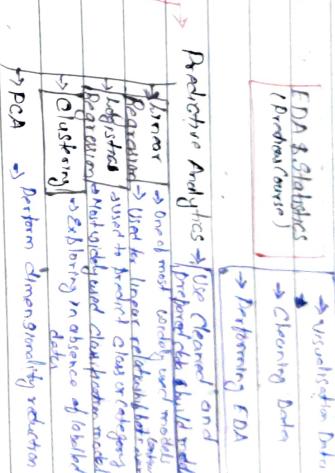
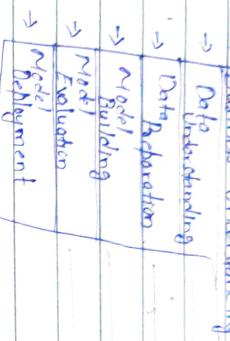


### Types of Machine Learning Algorithms

1. Regression → O/P variable to be predicted is a Continuous Variable
2. Classification → Predict a Categorical Variable.  
Ex: Spam or Ham
3. Clustering → No notion of a label is allocated to group/  
Cluster formed. Ex: Customer Segmentation.
4. Classification  
Clustering
5. Labels (Semantically correct) → N/C labels Ex: Type of cluster

## Machine Learning

### Crash-DL framework



### Pre-requisites

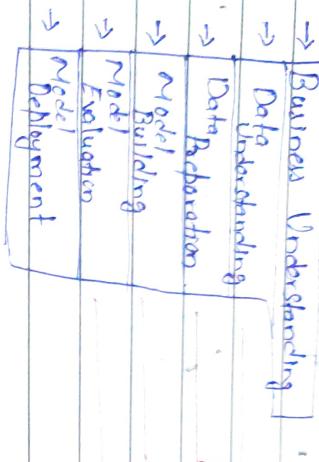
1. Exponentials
2. Logarithms
3. Equation of a straight line
4. Basis Transformation
5. Eigen decomposition
6. Inferential Statistics
7. Hypothesis Testing

### Types of Machine Learning Algorithms

1. **Regression** → Output variable to be predicted is a continuous variable
2. **Classification** → Predict a Categorical Variable
  - Ex: Spain or France
3. **Clustering** → No notion of a label is allocated to groups/cluster formed. Ex: Customer Segmentation.
- Classification | Clustering
- Labels (Sun & Moon example) → NO labels Ex: Type of Cluster

# Machine Learning - I

## Straight-DM Framework



PDA & Statistics (Prerequisites)	→ Predictive Analytics → Use Descript. and Inferential PDA
EDA	→ visualisation Data
Statistics	→ Cleaning Data
Probability	
Regression	→ Use Descript. and Inferential PDA
Classification	→ Used for linear classification Logistics → Used to predict class or category Decision Tree → Multiclass classification model
Clustering	→ Clustering → Existing no. of labelled data → Perform dimensionality reduction PCA → Perform dimensionality reduction

## Pre-requisites

1. Exponentials
2. Logarithms
3. Equation of a straight line
4. Basis Transformation
5. Eigen decomposition
6. Inferential Statistics
7. Hypothesis Testing

Types of Machine Learning Algorithms

1. Regression → Predict a Continuous Variable
2. Classification → Predict a Categorical Variable.
  - Ex: Spam or Ham
  - Ex: Label is allocated to groups / cluster formed.
3. Clustering → No notion of a label. Ex: Customer Segmentation.

Classification

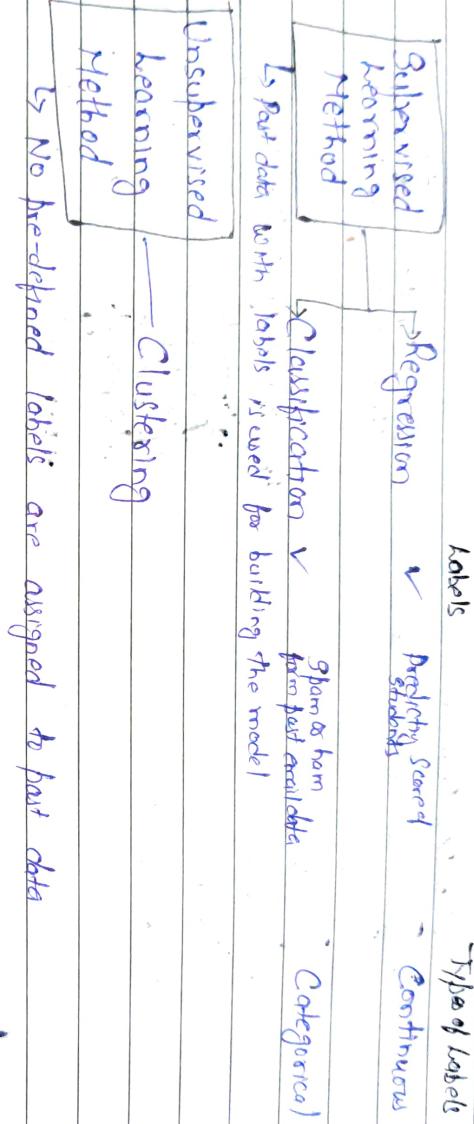
C | cuttin

Classification

C | cuttin

Labels (seen in the result) → NO labels. Ex: Type of Cluster

## Supervised and Unsupervised learning



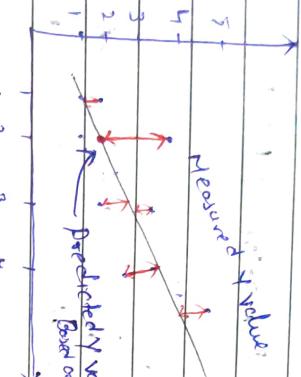
→ It has labeled data

### Types of Regression

1. Simple linear Regression
  - ↳ Only 1 independent variable
2. Multiple linear Regression
  - ↳ Model with more than 1 independent variable.

### 1. Simple

### Best Fit Line



→ At  $x=2$ , measured  $y$  value = 2

Predicted  $y$  value = 2

⇒ Difference b/w measured and predicted value is called **Residuals**

$$\text{Ex: } 4 - 2 = 2 \quad (\text{residual value})$$

→ Syntax  $\Rightarrow [e_i = y_i - y_{\text{pred}}]$  It is also called as Error

### Ordinary Least Squares Method

O-L-S is a method of minimising the total error square.

$$e_1^2 + e_2^2 + e_3^2 + \dots + e_n^2 \quad (\text{RSS Residual Sum of Squares})$$

## Coefficient of Least Squares Regression Line are determined by Ordinary Least Squares method.

Page No.:	M	T	W	T	F	S
YUVVA						

- We need to find the best value of RSS which will best fit to the eqn of straight line:

$$y = \beta_0 + \beta_1 x$$

Intercept      ↗ Slope

error  
foret. Point →  $y_i - \hat{y}_i$

$$e_i = y_i - (\beta_0 + \beta_1 x_i)$$

$$e_i = y_i - \hat{y}_i$$

$$\text{RSS} = (y_1 - \beta_0 - \beta_1 x_1)^2 + (y_2 - \beta_0 - \beta_1 x_2)^2 + \dots + (y_n - \beta_0 - \beta_1 x_n)^2$$

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

- We need to find the best possible of RSS and we know  $y_i$  and  $x_i$ .
- but we need to find the <sup>optimal</sup> value of  $\beta_0$  and  $\beta_1$ .
- $y_i$  is an absolute quantity, so changing the unit of dependent variable, RSS changes.

### \* Ways to minimize Cost Function (RSS)

#### 1. Differentiation

#### 2. Gradient Descent Function

→  $y_i$  is an optimization algorithm which optimizes the objective function for linear regression (Also Cost Function)

→ reach to the optimal soln.

#### # TSS

- TSS is an absolute quantity, so changing the unit won't affect the TSS.
- Dependent variable RSS changes.
- Hence we need a relative approach instead of RSS (absolute approach).

→ Calculate the ~~Y<sub>avg</sub>~~ of all the given points.

•  $R^2 = \frac{\text{RSS}}{\text{TSS}} = \frac{\sum (y_i - \bar{y})^2}{\sum (y_i - \bar{y})^2}$  :  $R^2$  tells how good our model is.

$$\text{Intrinsic value} = \bar{y}$$

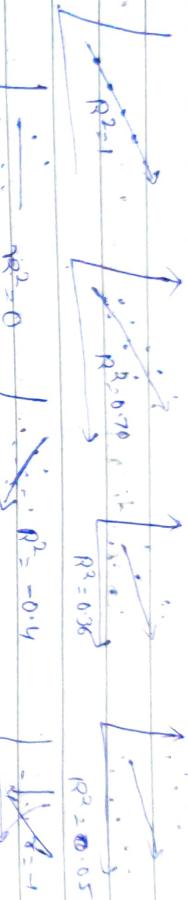
$$\text{TSS}$$

$$R^2 = \frac{1 - \text{RSS}}{\text{TSS}}$$

High value for  $R^2$  is good.

→ To correlation coefficient and R squared value is same.

Page No.	YOUNA
Date	



### RSE (Residual Sum of Square)

$$\rightarrow \text{RSE} = \sqrt{\frac{\text{RSS}}{df}}, \quad df = n - 2 \quad (n = \text{no. of data points})$$

df = degree of freedom

→ It also has some disadvantage of RSS.

→ R value lies between 0 - 1.

- 1- Variance in the data being explained by the model
- 0 - None of the variance value is being explained by the model

→  $R^2$  (Correlation)  $\Rightarrow$

$\Rightarrow$  Slope determining the correlation between variable  $\Rightarrow$  a straight line

determine the correlation between variable

lies between (-1, 1).

→ Correlation Coefficient lies between (-1, 1).

→  $R^2$  = How much variability in  $y$  (Dependent variable) we are able to explain

with the model:

With the model Sum of Squared

$$R^2 = \frac{SSE}{SST} = 1 - \frac{RSS}{TSS} \quad (\text{Total Sum of Squared})$$

(Variability in  $y$  explained by model / (total variability in  $y$ ))  $\times 100$  unrounded

$$MSS + RSS = TSS$$

## Gradient Descent (Steepest Descent)

M	I	N	R	S	S
Page No.	Date	24/10/21	Yousaf		

- It is an iterative 1st-order optimisation algorithm used to find local minimum/maximum of a given function.
- It is a method to minimise a cost/loss function.
- Cost function quantifies the error between predicted value and expected value and present it in the form of a single real no.
- It was originally proposed by CAUCHY in 1847.

### # Function Requirements

- 91) does not work for all func.
- There are two specific requirement

1. Differentiable

#### 2. Convex

- If a function is differentiable it has a derivatives for each point in its domain. Not all func. meet these criteria.

$$\frac{df(x)}{dx} = 2x$$

$$f(x) = 3x^2$$

$$\frac{df(x)}{dx} = 6x$$

$$f(x) = x^3 - 5x$$

$$\frac{df(x)}{dx} = 3x^2 - 5$$

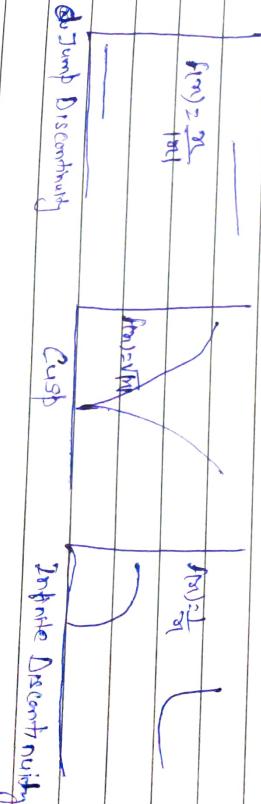
$$\frac{df(x)}{dx} = \text{undefined}$$

$$f(x) = \frac{x}{|x|}$$

$$f(x) = \sqrt[3]{x}$$

$$f(x) = \frac{1}{x}$$

↑ Differentiable functions ↑



## Gradient Descent (Steepest Descent)

- It is an iterative 1st-order optimisation algorithm used to find local minimum of a given function.
- It is a method to minimise a cost/loss function.
- Cost function quantifies the error between predicted values and expected values and present it in the form of a single real no.
- It was originally proposed by CAUCHY in 1847.

### # Function Requirements

- It does not work for all functn.
- There are two specific requirement

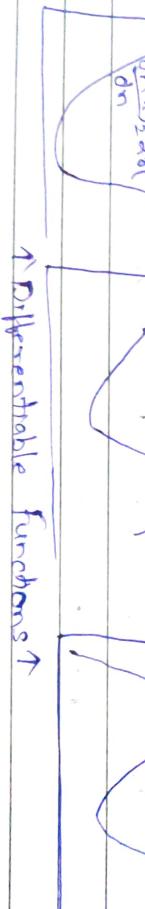
1. Differentiable
2. Convex

### 1. Differentiable

- If a function is differentiable it has a derivatives for each point in its domain. Not all functions meet these criteria.

$$\begin{cases} f(x) = \sin(x) \\ f'(x) = \cos(x) \\ \frac{d^2f}{dx^2} = -\sin(x) \end{cases}$$

$$\begin{cases} f(x) = x^2 \\ f'(x) = 2x \\ \frac{d^2f}{dx^2} = 2 \end{cases}$$



↑ Differentiable functions ↑

$$f(x) = \frac{x}{1+x}$$



↓ Jump Discontinuity

Cusp

Infinite Discontinuity

↑ Non-Differentiable Functions ↑

## 2. Convex

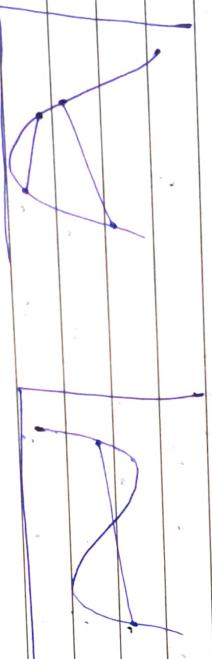
①  $\rightarrow$  For a univariate function, it means that the line segment connecting 2 function's points lies on or above its curve (it doesn't cross it).

$\rightarrow$  It means that it has a local minimum which is not a global one.

$\rightarrow$  A point  $x_1, x_2$  laying on the function's curve this condition is expressed as:

$$f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda)f(x_2)$$

λ = Denote a point's location on a section line and its value has to be between left point and right point



Convex

Non-Convex

2. If a univariate function is convex is to calculate the second derivatives and check if its value is always bigger than 0.

$\frac{\partial^2 f}{\partial x^2} > 0$  (strictly convex)

$\rightarrow$  Quasi-convex function can be used with Gradient descent algorithm and they are so called saddle point called minimax points, where algorithms can get stuck.

### 1st Derivatives

$$f(x) = x^4 - 2x^3 + 2$$

$$\frac{df(x)}{dx} = 4x^3 - 6x^2 = x^2(4x - 6) \quad \left\{ \begin{array}{l} x=0 \\ x=1.5 \end{array} \right.$$

Note: At  $x=0, x=1.5$  the first derivatives is 0 so these places are candidates for function's extrema (min or max)

→ Slope is 0 there

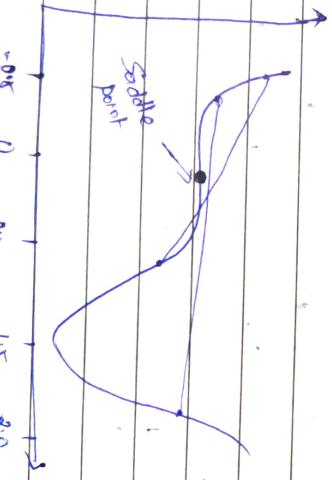
### 2nd Derivatives

$$\frac{d^2 f(x)}{dx^2} = 12x^2 - 12x = 12x(x-1)$$

$$\frac{d^2 f(x)}{dx^2} = 12x(x-1) \quad \left\{ \begin{array}{l} 0, x=0 \text{ and } x=1 \end{array} \right.$$

→ The local where second derivative is 0 are called an inflection point, a place where the curvature changes sign (mean it changes from convex to concave and vice-versa).

- \* For  $x < 0$ : funkt is convex
- \* for  $x < 1$ : funkt is concave
- \* for  $x > 1$ : funkt is convex again



At  $x=0$ , 1st 2nd derivative  
is 0 so it's called a  
saddle point  
 $x=1.5$  global minimum.

→ It is a slope of a curve at a given point in a specified direction.

→ In case of a univariate func., it is simply the 1st derivatives at a selected point.

→ In case of a multivariate func., it is a vector of derivatives in each main direction (along variable axis).

→ A gradient for an n-dimensional func.  $f(x)$ ,

$$\nabla f(p) = \begin{bmatrix} \frac{\partial f(p)}{\partial x_1} \\ \vdots \\ \frac{\partial f(p)}{\partial x_n} \end{bmatrix}$$

Example:-  $f(x) = 0.5x^2 + y^2$

$$\frac{\partial f(x,y)}{\partial x} = x; \quad \frac{\partial f(x,y)}{\partial y} = 2y$$

$$\nabla f(x,y) = \begin{bmatrix} x \\ 2y \end{bmatrix} \quad \begin{bmatrix} [x_0, y_0] \Rightarrow [x_0] \\ [x_0] \end{bmatrix}$$

→ Slope is twice steeper along the Y axis.

### Gradient Descent Algorithm

- It iteratively calculates the next point using gradient at the current position, then scale it (by a learning rate) and subtracts obtained value from current position.
- It subtracts the value because we want to minimise the function (to maximise it would be added).

$$\boxed{p_{n+1} = p_n - \eta \nabla f(p_n)}$$

$\eta$  = scales the gradient and control step size  
 (Learning rate and have strong influence on performance)

Note:-

- If the learning rate is too high, we might OVERSHOT the minima and keep bouncing, without reaching the minima.
- If learning rate is too small, the training might turn out to be too long or may reach maximum iteration before reaching the optimum point.

In Summary, Gradient Descent Method's steps are:

1. Choose a starting point (initialisation).
2. Calculate gradient at this point.
3. Make a scaled step in the opposite direction to the gradient (objective minimise)
4. Repeat points 2 and 3 until one of the criteria is met
  - \* max no. of iterations reached
  - \* step size is smaller than the tolerance.

## Simple Linear Regression, i.e., $y^{\text{pred}}$

### Making Predictions

$$\hat{y} = \beta_0 + \beta_1 x$$

C.  $\wedge$  mean estimate i.e. estimated value of  $y$ ,  $\hat{y}$

$$\begin{aligned}\epsilon &= y_i - \hat{y}_{\text{pred}} \\ \text{RSS} &= \sum_{i=1}^n (\hat{y}_i - \hat{y}_{\text{pred}})^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2.\end{aligned}$$

\* Model's Formulation  $\rightarrow Y = \beta_0 + \beta_1 X + \epsilon$

\* Fitting a line: minimize the sum of squared residuals

\* Assessing goodness of fit:

$$\text{R-squared} = \frac{\text{RSS/TSS}}{\text{TSS}} = 1 - \frac{\text{RSS/TSS}}{\text{TSS}}$$

### Making a prediction

For a new instance, predict the output as  $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$   
 $\Rightarrow t = \hat{\beta}_1 / \text{SE}(\hat{\beta}_1) \Rightarrow t\text{-value for a predictor coefficient is given by its standard error}$

### # Assumptions of Simple linear Regression

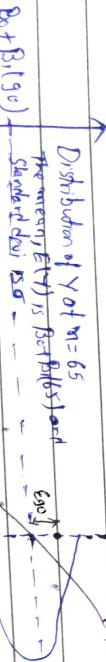
$\rightarrow$  We try to find the state of population by using sample data to make a model.

$\rightarrow$  Do so which we add some uncertainty in the output from the model.

\* In linear regression, at each  $x$ , finds the best estimate for  $y$ .

- At each  $x$ , there is a distribution on the value of  $y$ .  
 - Model predicts a single value, therefore there is a distribution of error terms.

$\rightarrow$  True Regression line,  $E(y) = \beta_0 + \beta_1 x$



$\rightarrow$  Distribution of  $y$  at  $n=90$   
 The mean,  $E(Y)$  is  $\beta_0 + \beta_1(90)$  &  
 Standard dev. is  $\sigma$  (standard variance)

$$\beta_0 + \beta_1(65) \rightarrow \text{True Reg. Line}$$

$$\beta_0 + \beta_1(90) \rightarrow \text{Predicted Value}$$

Note:- There is no assumption on the distribution of  $\epsilon$  only.

M	T	W	T	F	S	S
Page No.:	YOUVA					
Date:						

→ We know for every value of  $X$  we can have different value of  $y_i$  but the model gives one value of  $\hat{Y}$  when there are many and this introduces some error.

### \* Assumption

1. → There is a linear relationship between  $X$  and  $Y$ .

2. → Error terms are normally distributed <sup>(with mean zero)</sup> ( $\epsilon \sim N(0, \sigma^2)$ )

• Mean for one value of  $X$  we can get different value of  $Y$ . we chose one value of  $Y$ , so there is an error for each chosen value of  $Y$  that follows a normal distribution.

• For every  $x_i$  there is a distribution of error.

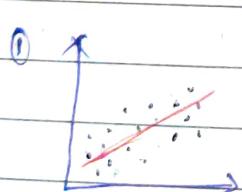
3. → Normal distribution of error <sup>the std. deviation</sup> is same across different values of  $X$ . i.e. Std. dev.  $\sigma$  is same at  $x=65$  or  $90$ .

• Error has constant deviation or variance across different values of  $X$ .

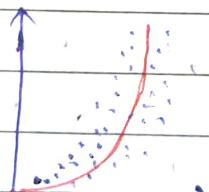
4. → Error term is independent of each other

• Mean previous doesn't have any effect on present error.

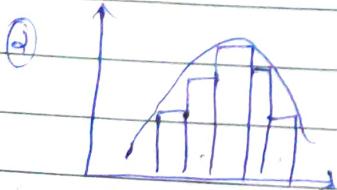
5. → Error terms have constant variance (homoscedasticity)



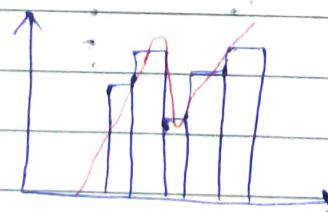
Linear pattern



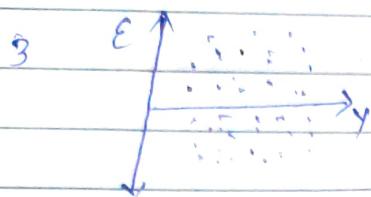
Non-linear pattern



Error term  
normally distributed

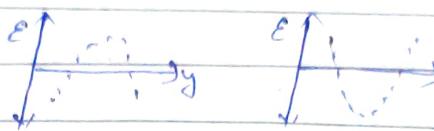


Error term not normally distributed.



No visible pattern

Error terms independent

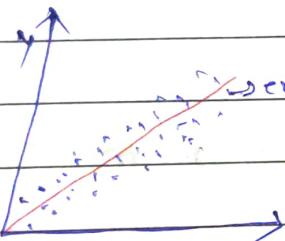


visible Pattern - Error terms dependent

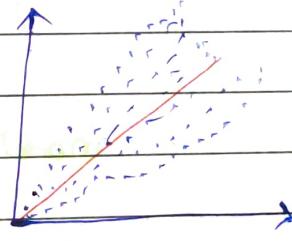
→ Not understand

5. → The variance should not increase (or decrease) as the error value changes

→ The variance should not follow any pattern as the error terms change.



Constant variance (Homoscedastic)



Changing variance (Heteroscedastic)

↳ Inferences made on the model

would be unreliable

↳ Even if we fit a line through the model, we can't make inference about model.

↳ Parameters used to make inferences will become highly unreliable.

## Syllabus

1. Reading and understanding Models  
 2. Training the model  
 3. Residual Analysis  
 4. Prediction and evaluation on the test set.

M	T	H	R	F	S	S
YOUVA						Date 26/12/21

## Building a linear Model

TV      Radio      Newspaper      Sale  
 ↗ Predictable Variable      ↗ Target Variable

$$y = c_0 + c_1 x_1 + c_2 x_2 + \dots + c_n x_n$$

$y$  = is the response

$c$  = Intercept

$c_i$  = Coefficient of  $i^{th}$  feature

$c_{n+1}$  = Coefficient of  $n^{th}$  feature

" $c$ " values are called "model coefficient" or "model parameters"

Note:- linear regression model can be built with

"sklearn" and "sklearn" packages/module/library

### Steps to Create a Model

1. Create predictable variable( $x$ ) and target variable ( $y$ )
2. Create train and test (unseen or hidden data) test sets  
(70-30 or 80-20)
3. Train model on training set.(i.e **Learn-the-Coefficients**)
4. Evaluate the model (Training set, test set).
  1.  $X = df[["x1"]], y = df["y1"]$
  2. import Sklearn,

from sklearn.model\_selection import train-test-split

$X\_train, X\_test, y\_train, y\_test = train-test-split$

train-test-split(X,y, trainSize=0.70, randomState=100)

Size of training set - 70%.

Size of test set

Only this will

$$P = .001 (\log \text{loss})$$

M	T	W	T	F	S	S
Page No.:						YOUVA
Date:						

### 3. Train model

Note: 1. "statsmodel" by default exclude "c" (Intercept) while training the model.

3-2. We need to explicitly add "c" value to the model.

# from statsmodel.api as sm

x-train-sm = sm.add\_constant("x-train")

it will add a column "const" to the dataframe of "x-train" with Const value = 1.

Const	x-train	if Const is not added then the line will pass through origin (by default)
1.0	2	
1.0	5	
:	:	
1.0	10	

$$y = c \cdot \text{const} + m \cdot x_1$$

### 3.3. Fitting the model

(OLS = Ordinary Least Square)

- lr = sm.OLS(y-train, x-train-sm)

↳ lr is a linear regression object (or instance of RLR)

↳ No learning is done over here only object of class OLS is created.

- lr-model.parameters

→ lr-model = lr.fit() // learning of coefficient

(lr-model.coef\_) // returns C and m (intercept & slope)

### 3.4. lr-model.summary()

- statsmodel gives detailed summary which is not present with sklearn / scikit learn

→ R-squared → 0.816 i.e. 81% of variance in the target variable is

1. due to computable variables.

→ Coef and P-value

→ R-squared is 81.6% very high

→ F-statistics) Prob ↓ means the fit is not purely

by chance!

→ If Prob(F-stat) < 0.05 then overall model fit is significant

- R-squared varies from 0 to 1
- R-squared 0 means none of the variance in the data is explained.
- R-squared implies that all of the variance in the data is explained.

M	T	F	S
Page No.:			
Date:			YOUVA

## 4. Residual Analysis and Prediction & Evaluation

### Residuals

- Now we need to verify our assumption that the (residual) error is normally distributed with mean equal to 0.
- Error should be independent of each other.
  - Look for patterns in residuals (we should not be able to identify)

$$y_{\text{train-pred}} = \text{lr-model.predict}(X_{\text{train-sm}})$$

$$\text{res} = y_{\text{train}} - y_{\text{train-pred}}$$

sns.distplot(res) // normal distribution of "res"

plt.scatter(X\_train, res)

### # Predictions

- Predict the y value of test set

$$X_{\text{test-sm}} = \text{sm.addconst}(X_{\text{test}}) \rightarrow$$

add const value to  
test set to make  
it compatible for  
model predict

$$y_{\text{test-pred}} = \text{lr-model.predict}(X_{\text{test-sm}})$$

Test Set  $\text{r2-score}(y_{\text{true}}=y_{\text{test}}, y_{\text{pred}}=y_{\text{test-pred}})$

Train Set  $\text{r2-score}(y_{\text{true}}=y_{\text{train}}, y_{\text{pred}}=y_{\text{train-pred}})$

- r2-score value tells us how well the model is performing on test set.

→ If r2-score of test set value lies with 10% difference of r2-score of train set then we can say that the model is very nicely able to predict the output on the test set.

Note:- ① p-value need to be very small (less than 0.05/0.01), then we can reject the Null Hypothesis and conclude that the coefficient is significant.

② Sum of residual should be equal to zero, as residuals are normally distributed around zero (mean=0)

(P.T.G.)

Note:- ① How value of RMSE (Root Mean Square Error) is better as it tells the deviation of the predicted value of a model from the actual value.

### Linear Regression using SKlearn

→ Sklearn is a goto package for predictive analysis and machine learning.

# Steps

1. Create an object of linear regression
2. Fit the model
3. See the params, make predictions (train, test)
4. Evaluate ( $R^2$ , etc.)

1.

$X\text{-train-lm} = X\text{-train-lm.values.reshape}(-1, 1)$

y-train need to be single series

ignore all 0s  
train 1 col 0

Note:- ② R-Squared Value gives the extent of the fit. i.e. how much variance in the data is being explained by the model.  
It doesn't tell anything about the significance of fit.

### ④ F-statistics

→ It determines the overall model fit better or not.

→ The idea behind F-test is that it is relative comparison between the model that we have built and the model without any of the coefficient except  $B_0$ .

↑ F-stats then Prob(F-stats) ↓ → Model is significant

↓ F-stats then Prob(F-stats) ↑ (0.05) → Model fit is insignificant and the intercept is only model can provide a better fit.

5. t-value of the coefficient:

6. Popular Package to built linear Regression Model:

1. Statsmodel-api
2. SKLearn
3. Scipy

7. Correlation Coefficient value shows how strongly or loosely the <sup>relationship between 2</sup> variables are related to

8. Null hypothesis states that the coefficient  $B_i$  is equal to zero. (no statistically significant relationship)

Alternate hypothesis is  $B_i \neq 0$

9. Correlation coefficient should be in range [-1, +1]

A value beyond this range indicates an error in measurement

10. R-Squared why it is called that. (In Simple Linear regression)

(Correlation)  $\rightarrow$  Corr = np.corrcoef(X-train, y-train)  $\rightarrow$  Correlation coefficient

$$\Rightarrow R\text{-Squared} = \text{Corr}[0,1] ** 2$$

Correlation (Pearson) is also called "r" or "Pearson's R"

11. There is no notion of standard or good RMSE value.

$\rightarrow$  RMSE depends on the units of the Y variables

$\Rightarrow$  RMSE is not normalised measure

For 2 different model With same dataset the model With lower RMSE is better

12. Scaling does not affect the model statistics and the model goodness

Will remain the same.

$\rightarrow$  Scipy helps with interpretation  $\rightarrow$  Scipy was gradient descent at background interpretation becomes easy when we have different feature all on a Similar Scale.  $\rightarrow$  gradient descent will be faster <sup>variable</sup> not if we have, have Similar Scales.

Mean  $\approx$  0 or StdDev  $\approx 1$

Note: When we Scale X & Y then the slope or Beta coefficient will be same to Correlation beta X & Y.

## Multiple Linear Regression

- It gives the relationship between two or more independent input variables and a response variable.
- It is used when one variable is not sufficient to create a model and make accurate predictions.
- Adding the R-squared value of 2 or more variable doesn't decreases, either it will increase or remain same.

### \* Formulation of MLR

Extensions of Simple Linear Regression to 'adds' more factor effect

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \epsilon$$

### \* Interpretation of the Coefficient

→ Change in  $Y$  ( $E(Y)$ ) after per unit increase in the variable ( $X_i$ )  
When other predictors are held constant.

- \* A lot of ideas are same as or are simple extensions of SLR.
- Model now fits a "hyperplane" instead of a line
- Coefficient still obtained by minimizing sum of Squared error (Least Squares criterion)
- For inference, the assumption from SLR still hold
  - Zero mean, independent, <sup>Errors are</sup> Normally distributed that have Constant variance.

## Moving from SLR to MLR: New Considerations

### \* New Considerations:

1. Adding more isn't always helpful
  - a. Model may "overfit" be becoming complex
  - (i) Model fits the train set "too well" (memorizing the detail) and doesn't perform well on test set or other set.
  - (ii) Symptoms: high train accuracy, low test accuracy
- b. Multicollinearity
  - i. Associations b/w predictor variables

2. Feature Selection becomes an important aspect. So that 'overfit' should no happen and multiple variable increases the accuracy of the model in a generic way.

## Multicollinearity

- It refers to the phenomenon of having related predictors variables in the input dataset.
- In simple terms, in a model which have been built using several independent variables, some of those variables might be interrelated, due to which the presence of that variable in the model is redundant.
- We can drop some of these correlated independent variables as a way of dealing with multicollinearity.

## Effect

### \* Interpretation:-

- Does "Change in Y, when all the other are held constant" apply? No as the other variable are dependent.

### \* Inference

- Coefficients swing wildly, signs can invert.
- p-value are, therefore, not reliable.

## Does not Effect

- The predictions, precision of the predictions
- Goodness-of-fit Statistics

Multicollinearity is a big issue when we are trying to understand the model (but detailed explanation)

Two Basic way of dealing Multicollinearity

M	T	W	T	F	S
Page No.					YOUNA

Date: 20/12/21

1. Pairwise Correlation
    - looking at the correlation of independent variables between different pairs
    - Ex: → Scatter plot to visually inspect.
    - Correlations do quantify the linear association.
  2. Variance Inflation Factor (VIF)
    - Sometimes pairwise correlation is not enough
    - Instead of one variables, the independent variable might depend upon a combinations of other variables
    - VIF calculates how well one independent variable is explained by all the other independent variables combined
- \* Common heuristic for VIF values
- >10: Definitely high VIF value and the variable should be eliminated
  - >5: Can be okay, but is worth inspecting
  - <5: Good VIF value. No need to eliminate the variable.
- Note: VIF doesn't have to do anything with target variable.
- Ex:- Suppose variable =  $\alpha_1, \alpha_2, \dots, \alpha_{10}$
- VIF of  $\alpha_1$  calculated after dropping all independent variables except  $\alpha_1$  (predictor)
- will change (decrease) in general.

- We should have minimum no. of variables to make to build a good model
- > In case of interpretation we must deal with multicollinearity

NAME:	M	T	W	T	F	S
YOUVA						

## Dealing with Multicollinearity

### 1. Dropping Variables

- > Drop the variable that is highly correlated with others or other variables gives the same information.
- > Pick the business interpretable variable (if interpretation and explicability important)

Ex - Suppose 2 variable are strongly related with one and another, then we should drop a variable which is of less important from business perspective.

2. Create a new variable using the interactions of the older variables. (Some sample given above)

- > Add interaction features. drop original features  
Ex. From a variable, make a new feature. Which represents the information of both the variables. (It has some implications)
- > Variable Transformations. (Get all variable and make a new variable from it)
  - \* PCA (Principle Component Analysis)
  - \* PLS (Partial Least Square).

- Notes:-
1. Multicollinearity makes some variable redundant hence the p-values changes.
  2. Predictive power given by R-squared value is not affected because even though we might have redundant variable in the model.

## Dealing with Categorical Variables

### Handling Categorical Variables with few levels

Value	(dummy) Indicator Variable	Value	Indicator Variable	
Gender	Female	Furnishing Status	Furnished	Un-furnished
Male	(dummy value) 0	Furnished	0	0
Female	(dummy value) 1	Semi-furnished	1	0
		Unfurnished	0	0

$$x_{li} = \begin{cases} 1 & \text{if } i\text{th person is female} \\ 0 & \text{if } i\text{th person is male} \end{cases}$$

$$y_i = \beta_0 + \beta_1 x_{li} + \epsilon_i \quad \begin{cases} \beta_0 + \beta_1 + \epsilon_i & \text{if } i\text{th person is female} \\ \beta_0 + \epsilon_i & \text{if } i\text{th person is male} \end{cases}$$

(Base State)  $\rightarrow$  When  $x$  is 0 with any no. of variable.

If the person is a female it has an extra effect ( $\beta_1$ ) compared to base state.

Note:- Creation of dummy variables to convert a categorical variable into a numeric variable is an important step in data presentation.

Formula:- N no. of Categorical Variable

$(N-1)$  no. of dummy variable will be created

- In case of multiple variable with different levels then we need  $(n-1)$  dummy variable [  $n = \text{level of Categorical Variable}$  ]

### Feature Scaling

→ Why we scale feature?

1. Ease of interpretation

When we have multiple variable with different scale level then it becomes difficult to for co-efficient to signify importance of them.

But if we get all the variable on one scale then it becomes easy to compare the coefficient of one with another hence interpreted become easy.

M	T	W	T	F	S	S
Page No.:	YOUVA					
Date:						

## 2. Faster convergence for gradient descent methods

When we bring (scale) all the variable in a particular (nearly) range i.e. (-1, 1), (0, 1) (-3, 3) which is nearby then it, convergence for gradient descent method becomes faster.

Note:- Neither  $\beta$ -values nor Model Accuracy changes With scaling, it just changes the co-efficient which is upto our interpretation.

### Popular Feature Scaling Methods

1. Standardization ( $x' = (\bar{x} - \text{mean}(x)) / \text{std.dev}(x)$ )
2. MinMax Scaling ( $x' = (x - \min(x)) / (\max(x) - \min(x))$ )  
(Normalisation) ( $0-1$ ) range [Mean=0, std.dev=1]

2. Scaling should be done after train-test split because if you don't want test dataset to learn anything from the train data. Else before then test data will have information regarding like minimum & maximum values.

Read the extra note given at the end of this section.

## Model Assessment and Comparison

- In multiple linear regression, we may build more than one model and compare them to check which one yields optimal result.
- Shall we compare the F statistics, R-square of different model
- Selecting the best Model
- Trade off between explaining highest variance and Keeping it simple. (Bias vs Variance Trade-off)
  - Key idea: Penalize models for using higher no. of predictors
- We can have a perfect model with many variable but it may overfit or keep simple model with less required no. of variable and predict the model by keeping it simple
- Now we need to make the trade off by penalizing the model with high no. of variable.

Suppose:- we have 2 model with same R-Squared one with high no. of variable other with less no. of variable, so in that case we reduce the R-squared with less no. of variable (adjust the R-squared)

$$\text{Adjusted } R^2 = \frac{1 - (1 - R^2)(N-1)}{N-p-1}$$

(↓ Redundant variable is Not good  
(More redundant variable with same R-squared is a better option)

### Akaike Information Criteria (AIC)

$$AIC = n * \log(\text{RSS}) + 2p$$

(lower value of AIC is better)  
hence "true" penalty  
penalty increase linearly with ↑ no. of iff variable in the model)

→ It is used in automatic model Selection

→ BIC, Mallows' CP and many more.

Note:- BIC is similar to AIC, it just put harder penalty on the predictor.

N, n = Sample Size

p = no. of predictor variables.

Read about AIC, BIC, Mallows' CP from given link:

M	T	W	T	F	S	S
Page No.:	YOUVA					
Date:						

Note: → R-squared doesn't penalise the model for having more no. of variables.

- Adding variable to the model, the R-squared will always increase (remain constant if correlation b/w that variable and dependent variable is zero).
- R-squared assumes that that any variable added to the model increases the predictive power.
- Adjusted R-squared on the other hand, penalise the models based on the no. of variables present in it.
- So if we add a variable and the Adjusted R-squared drops, we can be sure that the variable is insignificant to the model and shouldn't be used.
- In case of multiple linear regression, we should always look at the adjusted R-squared value in order to keep redundant variables out from your regression model.

## Feature Selection

- In case of multiple linear regression we may have a few potential predictor variables, selection those is an important task.
- Try all possible combinations?  $2^P$  models for  $P$  features

### 1. Manual Feature Elimination:

- Build Model
- Drop features that are least helpful in prediction (high  $p$ -value)
- Drop features that are redundant (using Correlation, VIF)
- Rebuild model and repeat.

Note:- Manual feature Elimination is a good choice when we ~~have~~ have less feature suppose 10 or 20 but it is not practical approach once we have a large no. of features like, 50, 100 etc and so on.

### 2. Automated Approach

- Top 'n' feature: RFE (Recursive Feature Elimination)
- Forward/Backward/Stepwise selection: Base on AIC
- Regularization (Lasso)  
(It makes the Coefficient 0 for the variable which is not needed for the model)

Note:- 1. Backward and Stepwise Selection gives same output but latter method is more popular.

- A balanced approach: use a combination of automated (coarsegrained) elimination + manual (fine tuning) selection.

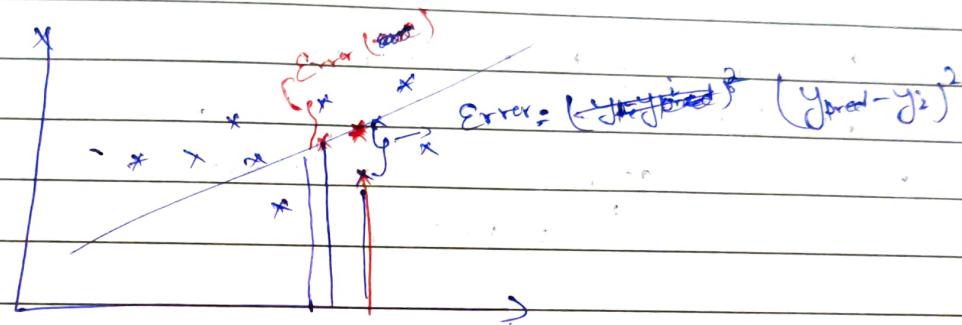
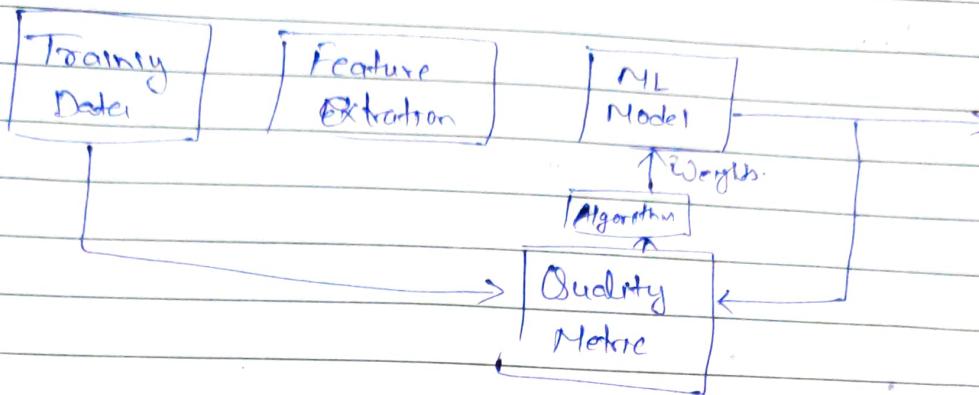
Note:- After automatic elimination of variable we need to use our expertise and Subjectivity to eliminate a few other features (fine tuning).

M	T	W	T	F	S	S
Page No.:						
Date:						YOUVA

Note:- → High p-value is indicative of the variable not having predictive power.

→ It might not necessarily mean that the variable is being described by one or more of the feature variable.

SGL-ML-C34-CC-Gr04

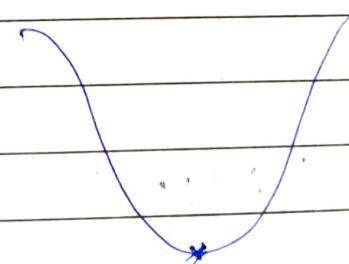


$$\text{All data} = \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

(RSS) → we amplify the error and out gradient descent can work properly on it.

### Gradient Descent

- $f(x)$  →



$$f'(x) \Rightarrow \frac{df(x)}{dx} = 0 \quad \left\{ \text{find } x \text{ will give minimum} \right.$$

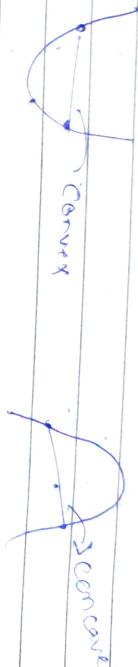
$$f''(x) \quad \frac{d^2f(x)}{dx^2}$$

\* Linear regression start point desired  
effect do local global minima

M	T	W	T	F	S	S
Page No:						YOUVA

## Convex & Concave

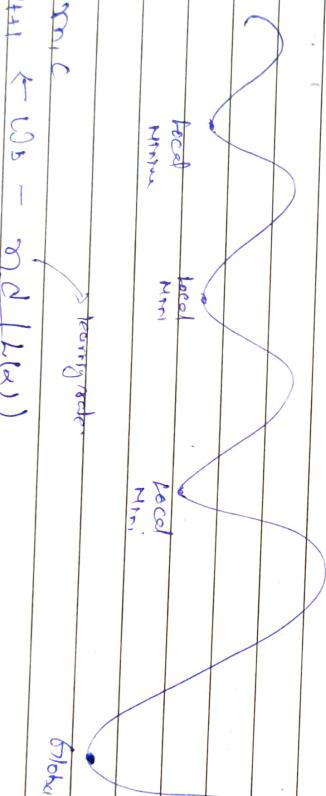
Minima concave non convex



\* Can help to find global minima.

GD doesn't know it is going to find minima or maxima.

\* Gradient descent  $\rightarrow$  global minima:  
 $\rightarrow \text{mini} - f(x)$



\* Linear regression is a convex curve

\* Based on gradient method we can find whether the curve is convex and non-convex.

\* Logistic regression used for classification.

\* Stats model is well for linear regression model but not so for logistic regression as it is bit hard to implement hence, Sklearn is used.

\* Linear regres → Close

$$\omega^* = (X^T X)^{-1} X^T Y$$

$$Inverse = O(n^3)$$

## S. Difference b/w Standardization & Min-Max

M	T	W	T	F	S	S
Page No.:	YOUVA					
Date:	02/11/22					

→ Linear regression is used in the various field such as real estate, telecom, e-commerce etc to build predictive model.

- We use mean, median, mode to <sup>impute</sup> missing value.
- The main method to graphically visualize Categorical variable and continuous variable is Boxplot.
- To convert a column from series to dataframe enclose it with in '[]'.
- For Categorical Variable "mode" is used to <sup>impute</sup>.

### Rescaling

1. To bring the coefficient of different variable in a certain range by bring all the variable in a comparable scale (Interpretability) (Interpretability).
2. If we rescale the variable within the range of 0-1 then the optimisation happening behind the scene becomes much faster i.e. minimisation routine becomes much more fast when we train a network using gradient descent funkt.
3. **Min-Max** rescaling should be used morecs if take care of outlier present in the dataset input by putting the outlier point to 1 and non-outliers in b/w 0-1.

Note:- The advantage of standardisation over the other is that it doesn't compress the data b/w a particular range as in Min-Max scaling. It is useful when there are extreme data point (outlier).

- Sklearn comes with a package called "Preprocessing" which has MinMaxScaler function in it.
- Fit on data
  - # fit: learns  $x_{min}, x_{max}$  on from data set.
  - # transform(): calculates  $(x - x_{min}) / (x_{max} - x_{min})$
  - # fit-transform(): does the work of above 2 methods in one method

### 3. Training the Model

- Plot heat map to check the correlation between all the numeric variables.
- Check the variable which has high correlation then keep on adding other called "Bottom-up" Approach

Note: Even 2 to 3% increase in the value of R-squared value after adding a variable to the model means a significant change.

#### Method to Drop Variable

- ① → P(t) below 0.05 (low) and vice-versa
- ② → Significance (P-value)
- ③ VIF

- High P-value, high VIF (Drop this variable)
- High P-value, low VIF (Remove these first)
- Low P-value, high VIF (Remove these after the one above)
- Low P-value, low VIF (Keep it)

Note: Drop a variable one by one, because after dropping the variable with highest VIF the other high VIF's variable will also get affected due to the drop of VIF variable with highest VIF.

- Never fit on Test Sets because test set are unseen data for the model so we don't want the model to know about it (i.e. its name, max, min).

## ML C34 - Linear Regression Module

Speaker → Karan

$$\text{RSS} = \sum_{i=1}^N (\text{Actual-output} - \text{predicted-output})^2$$

Residual

Total Square Sum

$$\text{TSS} = \sum_{i=1}^N (\text{Actual-output} - \text{average-of-actual-output})^2$$

It tells how much the data point move around the mean.

$$R^2 = \frac{\sum (y_p - \bar{y})^2}{\sum (y - \bar{y})^2}$$

$$RSE = \text{Square root } (\text{RSS} / \text{degree of freedom})$$

F-statistics → It tell whether the predictable model is better than a model with all  $\beta$  coefficient with 0 value.

Establishes prob → It tells the probability of a model with all  $\beta$  coefficient is equal to 0.

$$\left. \begin{array}{l} \beta_0 = \beta_1 = \beta_2 = \dots = \beta_n = 0 \end{array} \right\}$$

F-statistics: - If you assume all  $\beta$  coefficient is 0 then how good we are away from that model.

→ It tell the model which we have created and

## Linear Regression

1. Process of estimating relationship b/w variables.
2. Explains change in dependent variable with change in the values of predictors:
  - Simple linear Regression: Changing only one variable at a time.
  - Multiple linear Regression: Changing multiple variable at a time.
3. Uses :- **Forecasting and Prediction**
  - Regression guarantees "**interpolation**" of data not "**extrapolation**" necessarily.
4. Shows correlation, not causation.
  - Correlation does not imply causation.
5. A form of **Parametric Regression**

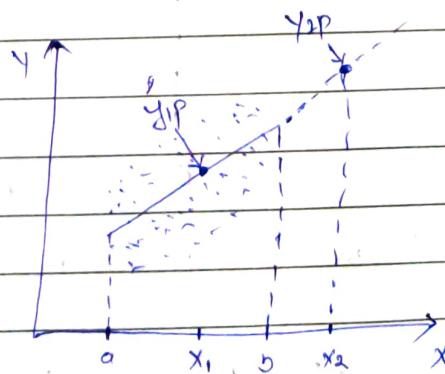
Parametric Regression	Non-Parametric Regression
1. Data follows fixed parameters	1. Data doesn't follow fixed parameter
	2. Dynamic in nature.

### Interpolation:-

It basically means using the model to predict the value of dependent variables on independent values that lie within the range of data we already have.

### Extrapolation:-

It means predicting the dependent variable on the independent values that lie outside the range of the data the model was built on.



$y_1P$  = Interpolation ;  $y_2P$  = Extrapolation

## Prediction vs Projection (Forecasting)

### Prediction

Importance of Outcome Focus: Identifying the driver variable and measuring their impact on dependent variables.

### Projection

Focus: Final projection result / forecasted value

### Assumption

No specific assumption is considered

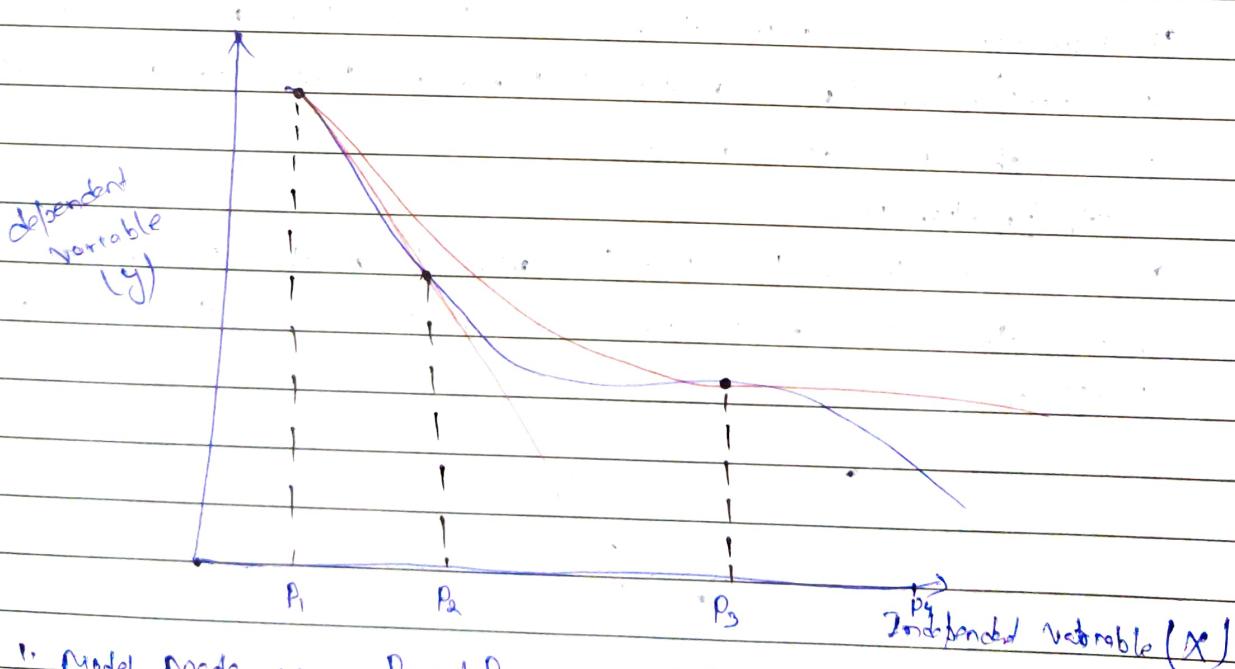
Assumes everything remains the same as today.

Forecast will change if new incident occurs.

### Complexity / Accuracy of Model

Simple models are better than complex models.

Choose accuracy over explanation.



- Model made using  $P_1$  and  $P_3$  - we can't predict accurately at value beta  $P_2$  and  $P_4$ .
- Model made using  $P_1$  and  $P_3$  - we can't predict accurately between  $P_3$  &  $P_4$ . (extrapolation)

## Assessing the Model Stability

- If R-squared and Adjusted R-Squared are very close it means ~~adjusted~~ none of the parameter are redundant in the model means no extra variable is added to the model.
- Errors bet actual and predicted OLP is normally distributed which confirms that there are no variables that could have helped to explain the model better.

1. Data set is small
  - Not advisable to separate testing and training sets.
2. Use bootstrapping
  - Choose 10-20% of sample randomly.
  - Opt for multiple iterations with replacement
  - Train model on "in-sample" group
  - Test on "out-sample" group.
3. Model is stable when R-squared is similar

### Note:- Least Square Error Method

- Used to find the best fitted line through the set of points.

### Mean Square Error

- Used to evaluate model after fitting it.
- If finds out the average of the difference b/w the actual and predicted value is a good parameter to compare various models on the same data set.

### Correlation Coefficient

- Pearson's R (used for linear relationship b/w variable)
- Spearman's R (used to determine the correlation b/w non-linear variables)

### # Bias

- Bias mean difference b/w predicted value by model and the real values.
- If is a error and need to have low bias for ML model.

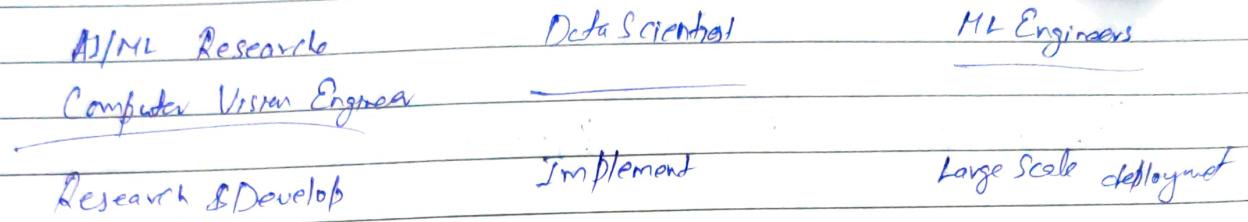
### # Variance

- Sensitivity of the model to small fluctuations in the training dataset.
- Low variance for ML algorithm.

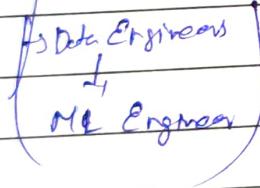
- Note:-
- Straight line mode will have low variance, high bias
  - High degree polynomial will have low bias, high variance

M	T	W	F	S	S
08				YOUVA	

## SGI-ML-C34-CC-G04



① Applied team:- Analyses the present model and check its feasibility if not



①.1 Data Scientist  $\Rightarrow$  EDA, Decision to use ML is required or not  
 ①.2 ML Engineers Build ML model

\* Data Scientist  
 $\Rightarrow$  EDA

\* ML Engineer  
 $\Rightarrow$  Build DL, ML model

Google Dataset