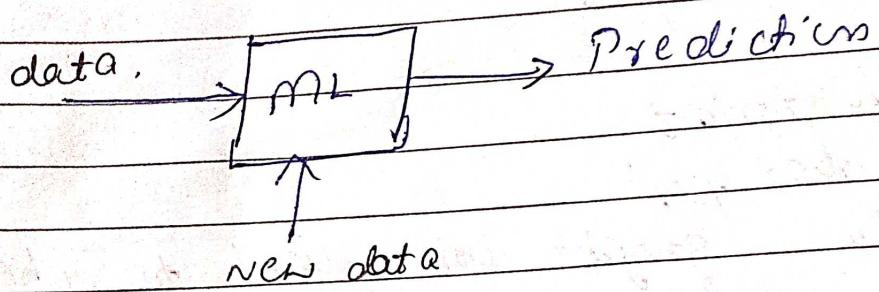


Machine learning

Machine learning is a branch of computer science which has several machine algorithms which learns from past data and makes predictions.



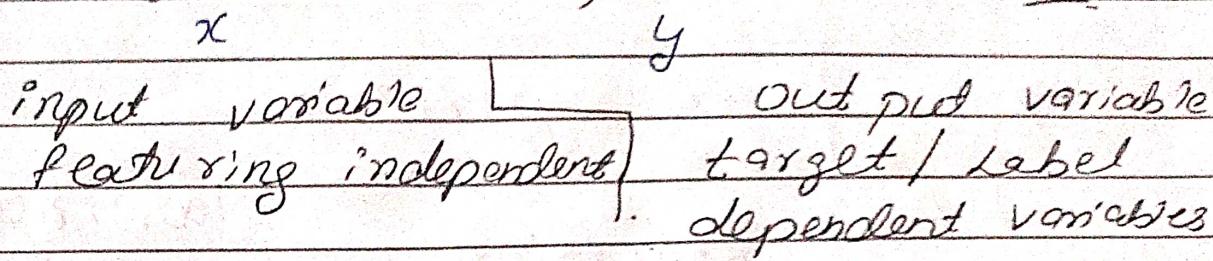
Types of Machine Learning

- ① Supervised Learning
- ② Unsupervised Learning
- ③ Reinforcement [Self driven or, learns from punishment and rewards]

① Supervised Learning:

Train machine learning algorithms with both input & output variables such algorithms are called as supervised learning.

Ex: predicting house prices, credit score etc



② Unsupervised Learning:

If we train machine learning with only input variable, such algorithms are unsupervised learning.

Ex: Market basket analysis, clustering of customers feedback, market segmentation.

③ Reinforcement learning:

It is a machine learning training method is based on rewarding and punishing

Ex: self driven car.



Supervised

Regression

Classification

- a) Regression :- If target is continuous
ex:- house price prediction
flight price prediction
- b) Classification :- If target is discrete
or categorical.
ex:- Diabetic, Ratings, covid19

Linear Regression :-

- ⇒ It comes under supervised Learning
- ⇒ It comes under regression and it is used to predict the continuous target variable
- ⇒ ^{Simpl} Linear Regression is used to find linear relationship b/w independent (x) and dependent (y) variables by fitting straight line.

Types of Regression:

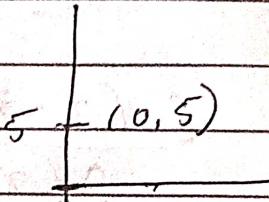
① Simple linear regression: If an algorithm uses single input variable and single output variable such an algorithm is called as Simple linear Regression.

$$y = mx + c$$

m = slope

$$m = \frac{dy}{dx}$$

c = y -intercept



② Multi linear Regression: If an algorithm uses multiple input variables and single output variables, such an algorithm is called as multi Linear Regression.

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

Assumption of linear Regression:

- ① Linearity: There should be a linear relationship b/w dependent & independent variable.
- ② Data should be normal.
- ③ Little or no multicollinearity
[There should be very less correlation among input variables.]

How Linear Regression Works?

<u>x</u>	<u>y</u>
no. of hours studied	marks obtained
2	20
3	30
4	45
7	75
8	85
10	95
5	?

Linear Regression

straight line Eqn

$$y = mx + c$$

Find; $m = ?$

$$c = 7$$

$$m = \frac{\varepsilon}{\varepsilon} \frac{(x - \bar{x})(y - \bar{y})}{(x - \bar{x})^2}$$

$$c = \bar{y} - m\bar{x}$$

$\bar{x} \in \bar{y} = \text{mean}$

after finding $m=8$, $c=5$

$$\boxed{y = 8x + 5}$$

→ Regression model

Best fitting Regre

Find for 5 hours,

$$\hat{y} = 8 \times 5 + 5$$

$$\hat{y} = 45$$

-sign line.

① Error / loss = $\frac{\text{actual value} - \text{Predicted value}}{1}$

x	y	\hat{y}	predicted	Error
2	20	21		1
3	30	28		1
4	45	37		8
7	75	61		14
8	85	69		16
10	95	85		10
5	?	10		

Calculation Total Error:

(1) MAE \rightarrow mean absolute error

$$MAE = \frac{\sum_{i=1}^n |y_i - \hat{y}_i|}{n}, \quad y_i = \text{actual} \\ \hat{y}_i = \text{predicted.}$$

(2) MSE \rightarrow mean squared error

$$MSE = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}, \quad MSE \text{ best because of it square}$$

(3) RMS E: Root mean squared error no true.

$$RMSE = \sqrt{MSE}$$

IMP

Note:-

- \rightarrow If Error is less, we say its good model
- \rightarrow If Error is high, we say its bad model.

Data preprocessing :-

- checks for missing values and handle them
- checks for duplicates
- checks for outliers and handle them

How to handle outliers :-

IQR

3 - Sigma Rule

- | | |
|---|--|
| <ul style="list-style-type: none"> • This method is used to find outliers when data is skewed / not normal | <ul style="list-style-type: none"> • This method is used to identify outliers when data is normal |
|---|--|

μ μ

μ

- Lower limit = $Q_1 - 1.5 \times IQR$
- Upper limit = $Q_3 + 1.5 \times IQR$

$$\text{Lower limit} = \mu - 3\sigma$$

$$\text{Upper limit} = \mu + 3\sigma$$

$\mu \Rightarrow \text{mean}$

$\sigma \Rightarrow \text{standard deviation}$

- $x < \text{Lower}$ and $x > \text{Upper}$ then x is outliers

- $x < \text{Lower}$ and $x > \text{Upper}$ then x is outliers

$$x \notin [Q_1 - 1.5 \times IQR, Q_3 + 1.5 \times IQR]$$

$$x \notin [\mu - 3\sigma, \mu + 3\sigma]$$

- WTF
- If we are using IQR replace with median

- WTF
- If we are using 3-sigma rule replace with mean

IQR:

① Find Q_1 and Q_3

② $IQR = Q_3 - Q_1$

③ Find lower limit and upper limit

lower limit = $Q_1 - 1.5 \times IQR$, upper = $Q_3 + 1.5 \times IQR$

VVI

→ handle only outliers i.e less than 5%, other wise don't touch it.

Scaling:-

Scaling is one of the preprocessing method which is applied on continuous data to bring all the values into a certain scale.

- Scaling is applied only on input variables

- There are two types

(1) min max scaler

(2) standard scalar

- It is called as normalization

- It is called as standardization

- $\text{Minmax} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$

- standardization

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

- range of minmax = $[0, 1]$
It will reduce all the values to lie bw $[0, 1]$

- range = $[-3, 3]$

Sklearn :-

Sklearn is one of the scientific library which includes all the fn/packages related to machine learning, data science, preprocessing etc.

Feature Engineering:-

- *) Drop irrelevant column
- *) Create New columns
- *) Select best features with target
 - Find corr
 - plot heat map
- *) Include the columns which has high correlation with target
- *) drop columns which has less correlation (<30%) with target
- *) ~~Drop~~
- No (②) little multicollinearity
- *) make sure the correlation b/w input variables should be less
- *) If correlation is high among input variables drop one of them.

random_state = 4, (shuffle data)

classmate

Date _____

Page _____

→ Split data for training and testing

70% , 30%
train (on) > test
80% , 20%

$x = \begin{cases} x - \text{train} \\ x - \text{test} \end{cases}$ $y = \begin{cases} y - \text{train} \\ y - \text{test} \end{cases}$

→ Training (x -train, y -train)

→ testing / prediction $\rightarrow x$ -test.

① R-squared / R² Score

• It is used to measure the strength of the model.

• It tells us about how good is the model.

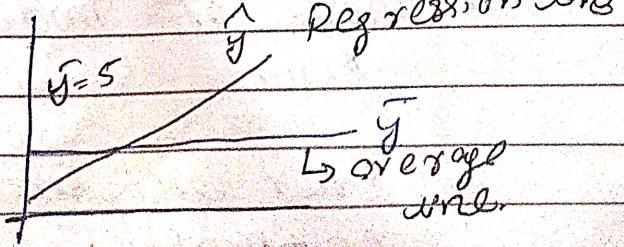
$$R^2 = 1 - \frac{RSS}{TSS} = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

• RSS \rightarrow Residual sum of Squares.

• TSS \rightarrow Total sum of squares.

• R² score compares regression line with average

$$R^2 \text{ score} = [-1, 1]$$



⑤ Adjusted R^2 -score:

Disadvantage of R^2 :

R^2 score increases as the number of independent variables increases which has very less relationship with target variable.

To overcome above issue we use adjusted R^2

$$\text{adj} = 1 - \frac{(1-R^2)(N-1)}{(N-P-1)}$$

$N \Rightarrow$ no of observation in test data

$R^2 \Rightarrow$ R^2 -score

$P \Rightarrow$ no of independent variables

if Adjusted R^2 -score < R^2 -score.

Then the model is good.

→ Cost function:

$$J = \frac{1}{2n} \sum (y_i - g_i)^2$$

- Goal of model is to reduce cost function error.

- Error can be minimised by choosing optimal values from m & c.

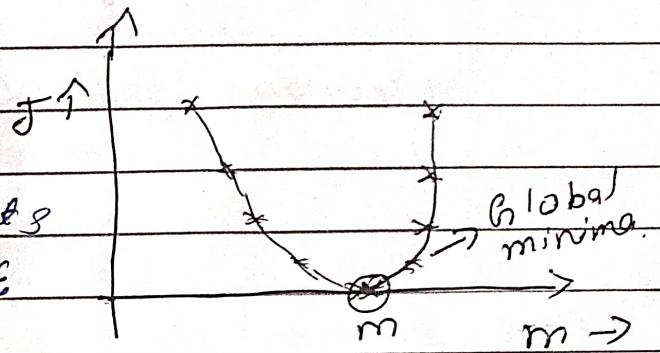
- How to choose optimal values for m & c

Gradient descent:-

Gradient descent is one of the ~~option~~ optimization method which helps in selecting optimal values for m & c such that error is minimum.

How gradient descent works?

- plot graph of error vs slope



- Gradient descent starts with random slope & works iteratively to reach global minima where error is minimum.

- m value at the global minima will be the optimal value.

$$m_{\text{new}} = m_{\text{old}} - \eta \frac{\partial J}{\partial m}$$

$\eta \rightarrow$ Learning rate

$$(c_{\text{new}}) = (c_{\text{old}}) - \eta \frac{\partial J}{\partial c}$$

(@)

η - Learning rate: It tells us about how many steps it has to take to reach global minima.

when η is large :- It overshoots and never reaches global minima.

when η is small :- It takes more time to reach global minima.

How to choose η ?

Initially take larger steps and then take smaller steps when near to global minima.