

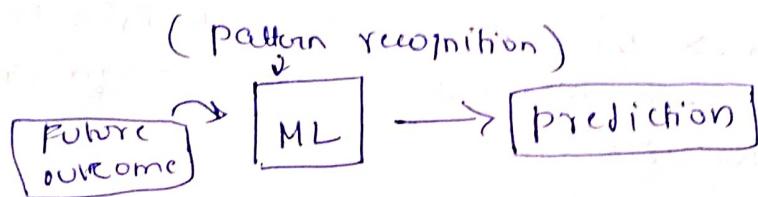
# Machine Learning

## Notes

By -> Rishabh Dwivedi

## Machine-Learning

- ↳ ML is the ability of machines to learn from past experience & predict the future outcome.



### How it works?

- ↳ past experience is not but **[DATA]**

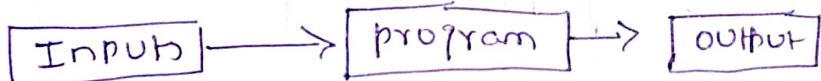
Any information about anything.



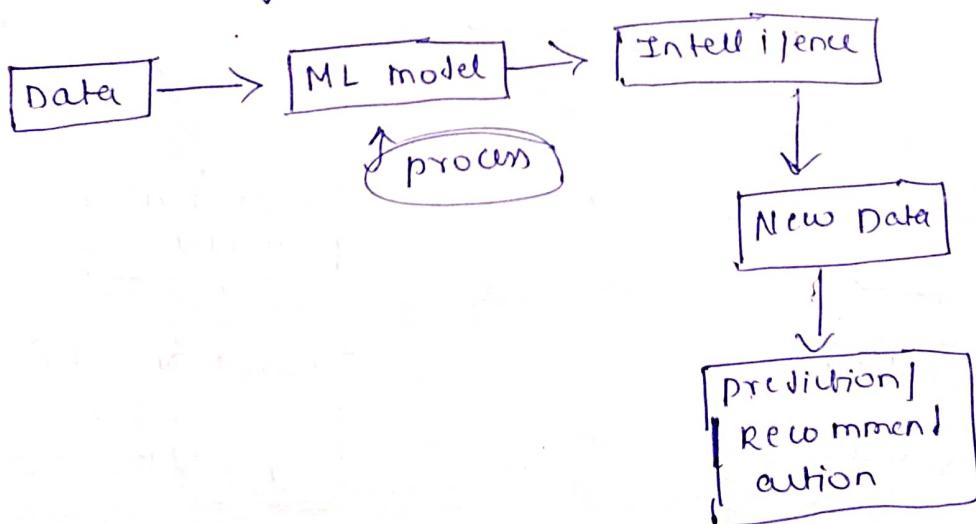
SIZE of data can be huge

∴ for normal human

being it is impossible  
to derive pattern from it.



↓↓↓ similarly



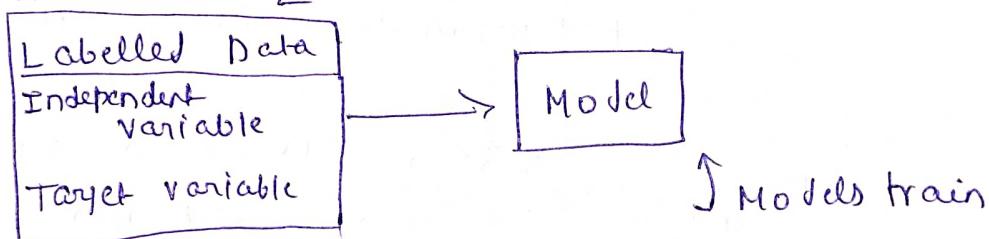
## TYPE OF ML

- Unsupervised
- Supervised

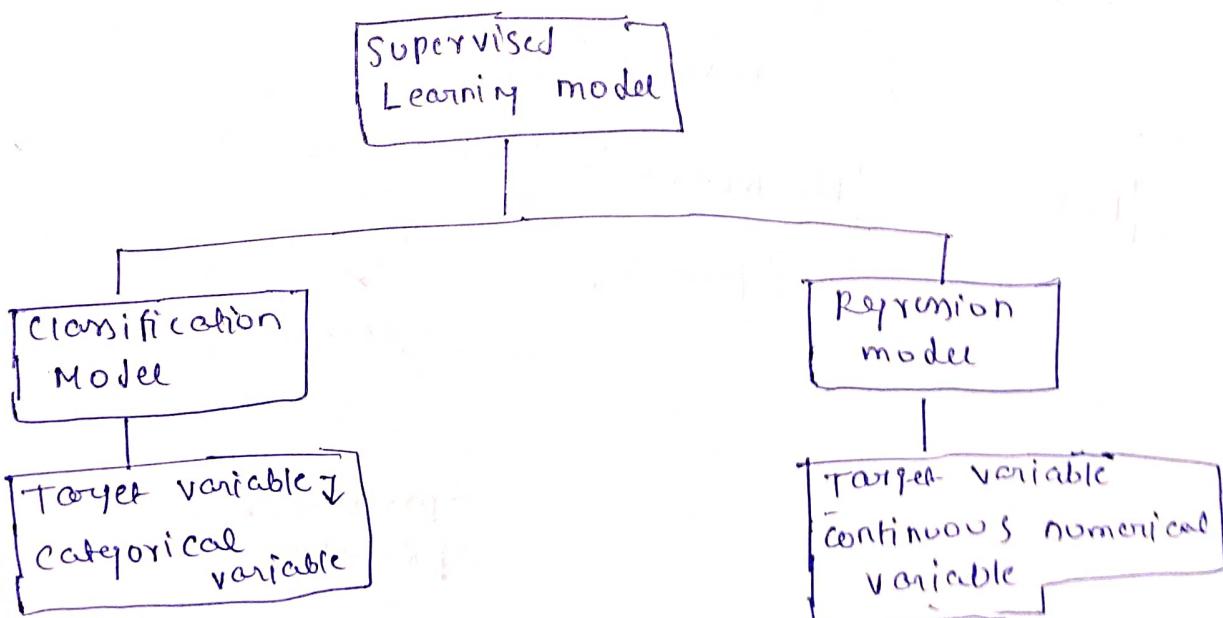
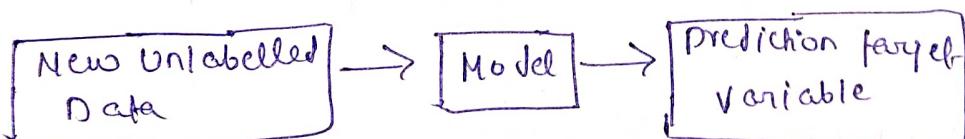
**Labelled Data** → The data which contains a target variable or an output variable that answers a question of interest is called Labelled Data.

**Unlabelled Data** → UL Data is one which contains information about something but does not have a predefined target variable.

### ① Supervised



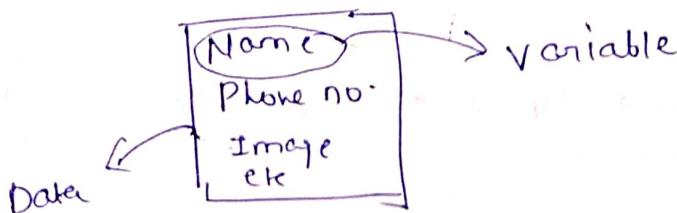
∴ Now its a time for testing



## Unsupervised Learning

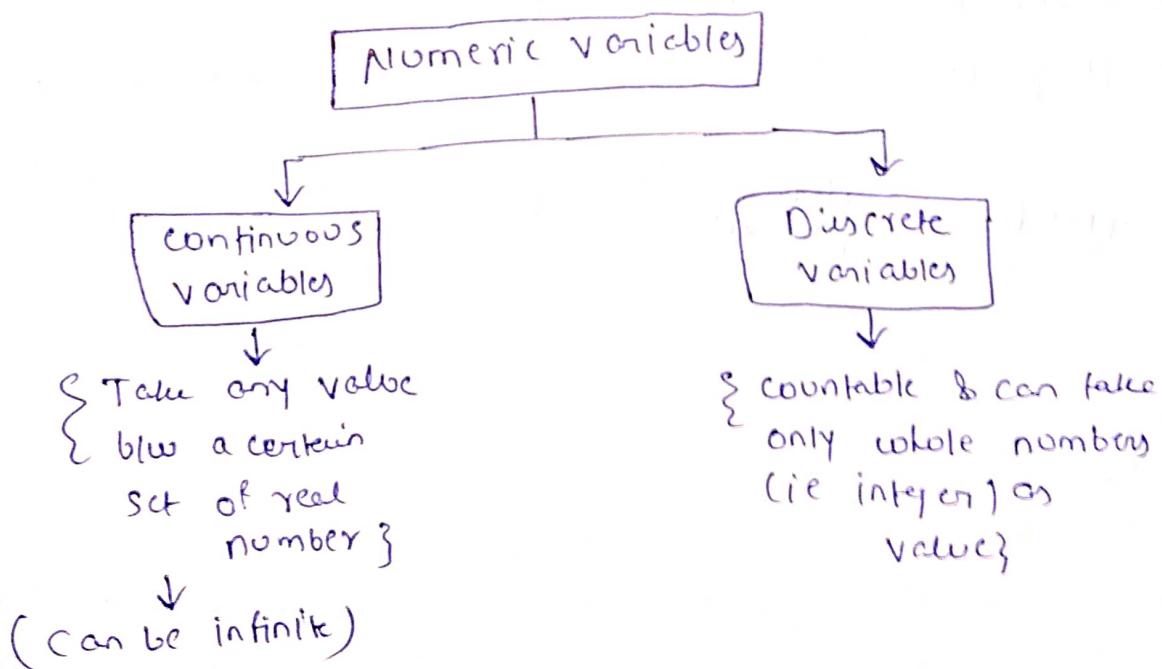


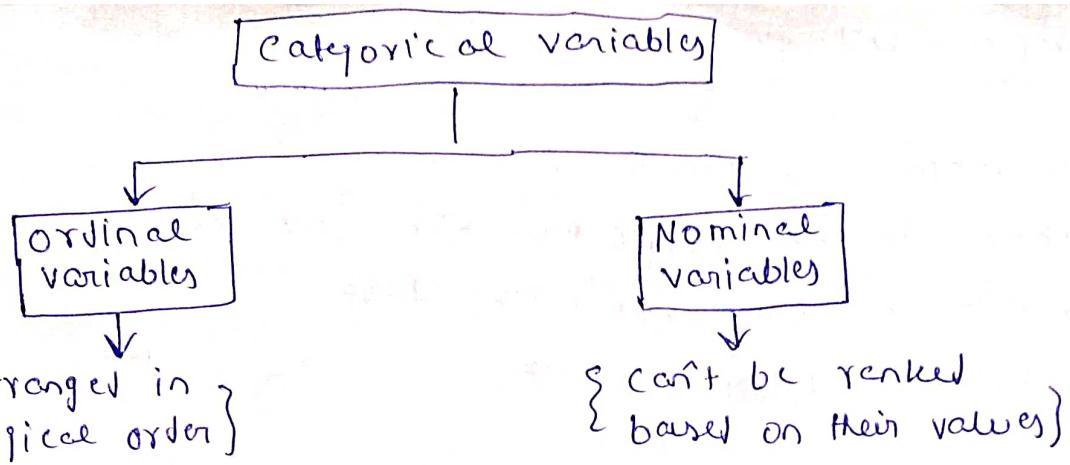
- complement of supervised learning
- There is no target variable involved.
- works on only unlabelled data
- identifies if some pattern exist in the data.
- Variable [can take any value ~ No. or String]



\* A variable represents one specific characteristic of the data or tell one specific information about the data under consideration.

- Numeric or Quantitative variable
  - ↳ { Measured or counted }
  - ↳ Numbers → exg 1, 2, ...  
(Also known as quantifiers)
- Categorical or Qualitative variable
  - ↳ { can't be measured }
  - ↳ Expresses a feeling or a characteristic





Eg: poor & Average & good

Eg: colors

### Independent variable

↳ are ones whose value don't depend on any other variables.

### Dependent variable

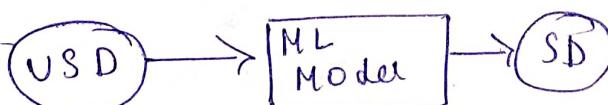
↳ whose value depend on other independent variable & cannot be changed easily.

### Structured data → Arranged in structured format.

Roll no	Name
300	Rishabh
500	Rohit

### Unstructured data → Not arranged in tabular form.

↑ Eg: msg received in chats.



⇒ Graphical representation of data.



Data Analysis

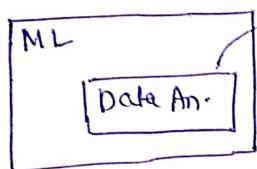
↳ Process of studying the available data and drawing valuable insights or information from it with the help of a software.

Sub	Marks
A	20
B	30
C	60

→ Analyse it

↳ various Software like

↳ python, R, Excel, Matlab  
etc.

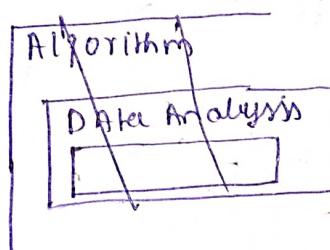


→ Ifs a part of NL which analyses your data & bring out insights before you apply your algo'

↳ Exploratory Data Analysis (EDA)

↳ To find out the pattern & relationship among the data variables.

\* EDA gives summary about dataset



**Descriptive Analysis**  
{what happened}

Step 1

**Diagnostic Analysis**  
{why did it happen?}

2

**Predictive Analysis**  
{what will happen?}

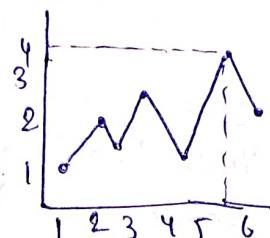
3

**Prescriptive Analysis**  
{what should we do?}

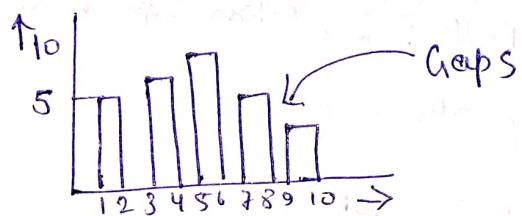
4

### Graph

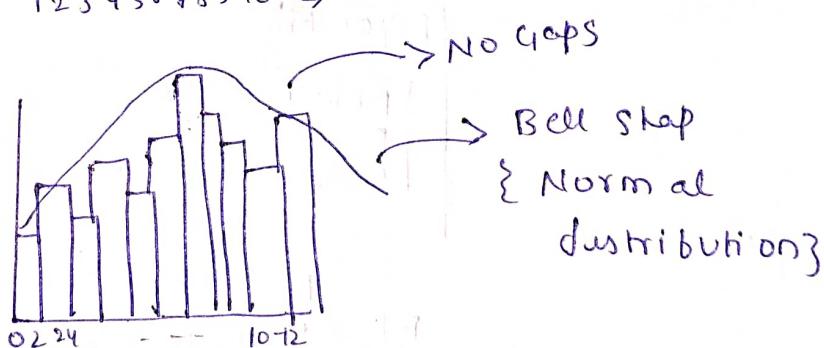
#### ① Line Graph



#### ② Bar Graph



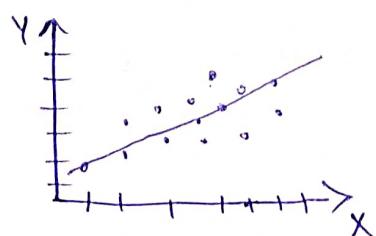
#### ③ Histogram



#### ④ Pie chart



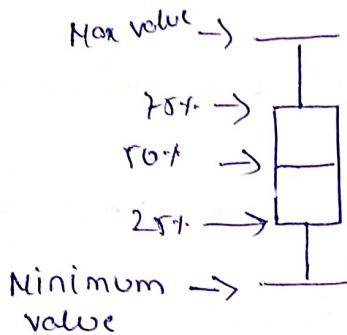
#### ⑤ Scatter plot



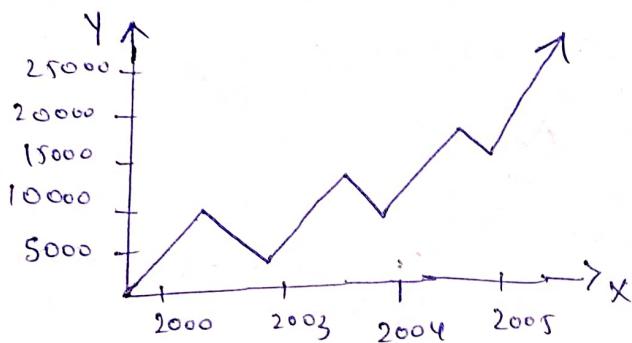
⑥

## ⑥ Box & whisker plot

↳ Graphical representation of the Statistical Summary of numeric variables.



## ⑦ Time series



\* Some python coding \*

↳ Numpy, pandas, Matplotlib & seaborn.

⇒ Descriptive Statistics.

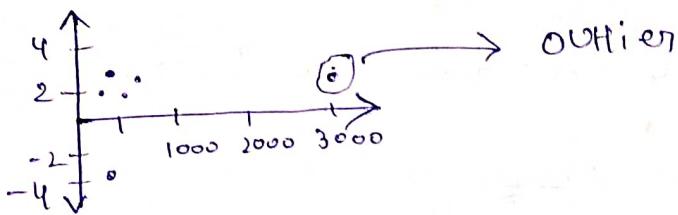
- ↳ count of values for a variable
- ↳ Mean
- ↳ Standard deviation
- ↳ Minimum value
- ↳ Maximum value
- ↳ Three percentiles - 1<sup>st</sup>, 2<sup>nd</sup> & 3<sup>rd</sup>

→ DS is calculated only for Numerical variable only.

↳ Uses of DS

- (i) Build a better understanding about data
- (ii) Identify & treat missing values
- (iii) Identify any anomalies

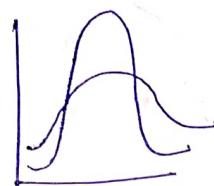
(7)



\* standard deviation tells the spread from the mean or average value or expected value.

↳ If SD is low  $\rightarrow$  most of the no. are close to the average value.

↳ If SD is high  $\rightarrow$  most of the no. are spread out.



model describe()

↳ count, mean, std, min, 25%, 50%, 75%, max

\* calculating SD using pandas & Numpy will show different result because  $\downarrow$

SD using pandas

$\hookrightarrow$  Formula used

$$SD = \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}}$$

Each value  
 $n =$  mean  
in the dataset  
 $\bar{x} =$  Mean of all  
value  
 $N =$  no. of values  
in dataset

SD using numpy

$\hookrightarrow$  Formula used

$$\sigma = \sqrt{\frac{\sum (x - \bar{x})^2}{n}}$$

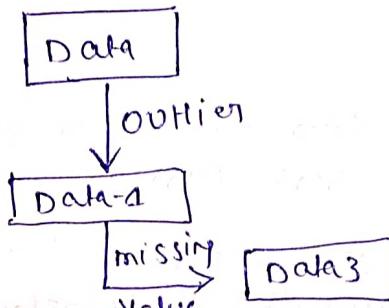
Fix it with the help of degree of freedom (df) = 1

$\downarrow$   
 $np.std(model['Sale Price'])$ , df = 1)

$$\hookrightarrow \sigma = \sqrt{\frac{\sum (x - \bar{x})^2}{n - \text{degree of freedom}}}$$

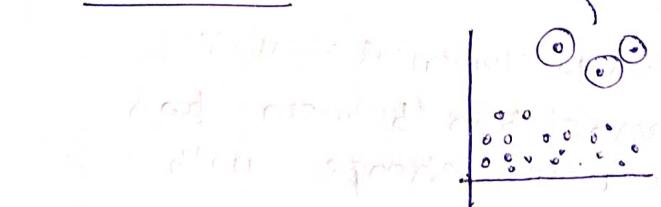
## ⇒ Data Exploration

- ↳ why?
- ↳ There may be anomalies / outlier or missing value present in the data, which can create wrong prediction. So to avoid it we perform data exploration.



As follows

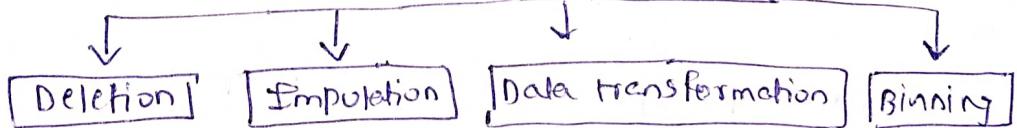
### Outliers ?



An outlier is a data point that is distinct from other data point.

∴ Any data point which is higher than upper limit or lower than lower limit is an outlier.

### Treating Outlier



↑ For target variable, independent variable

### Missing value

↳ NaN form

### Treating Outlier Missing value

Deletion

Imputation

## Variable transformation (Vt)

Vt refers to the process of making changes to a variable in a way that it become more useful and meaningful for analysis & modelling process.

why Vt?

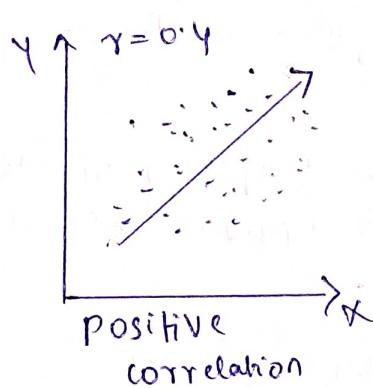
### (i) Outlier treatment

(ii) when a variable does not represent the information it contains in the best possible ways.

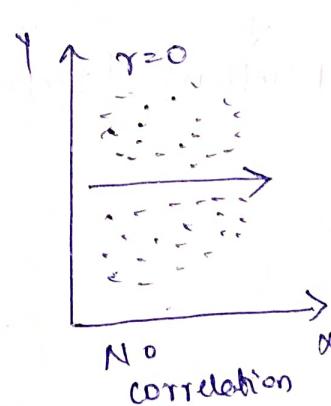
(iii) To derive a new variable from an existing variable or combining two variable to transform a variable so that it become more useful.

⇒ correlation ~ {Data exploration for Numerical Variable}

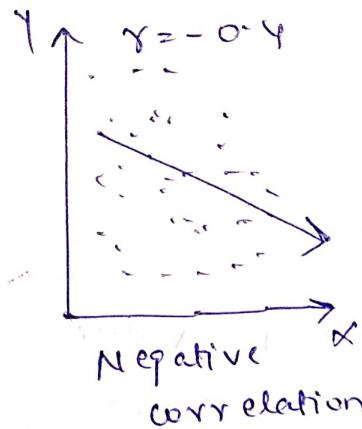
- ↳ It is measure of dependence or association between two variables i.e. how does one variable change with change in other.
- ↳ Correlation refers to how close two variables are to have a linear relationship with each other.



↓  
If 1 variable inc.  
other inc. to.



↓  
Random changes  
& no pattern.



↓  
With increase in  
1st v variable  
other v variable  
decreases.

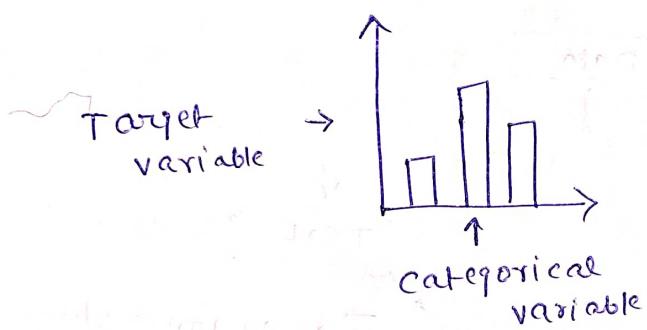
\* Any random value b/w -1 to +1.

Eg: model['Sale Price'].corr(model['Flat area'])

### using numpy

↳ np.corrcoef(data[-1], data[-1])

\* only those variables should be selected for building the model which have high correlation (above a certain cut-off value), positive or -ve, with the target variable.



\* The method of finding out whether a categorical variable is having an impact on the numerical variable

→ ANOVA

↓  
Analysis of variance

ANOVA

↳ It checks if the mean i.e. Average value of the target variable across different levels or unique values of a categorical variables are equal or not.

### 2 TYPES OF ANOVA

1-way ANOVA

↳ It deals with one independent variable only.

2-way ANOVA

↳ It deals with two or more independent variables.

## Splitting of data into train - test

↳ Before splitting, we need to separate the dataset into on the basis of dependent & independent variable

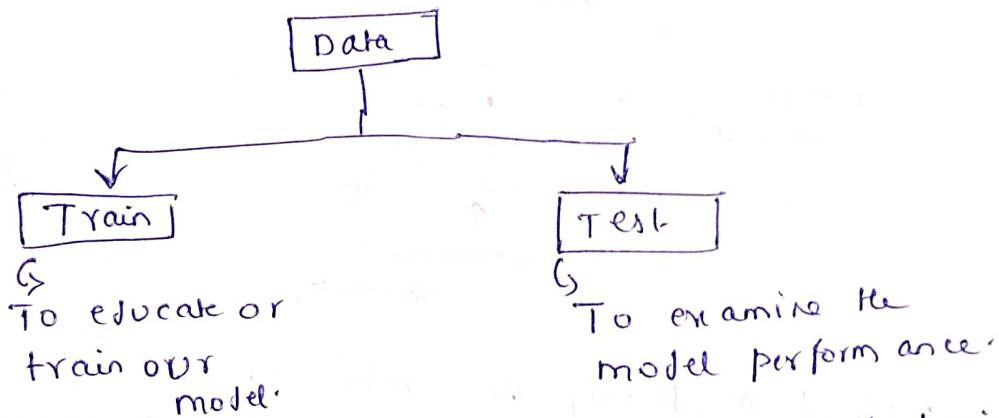


For doing this we can use "iloc" function

e.g.  $X = \text{data.iloc[:, 0]}$  ← dependent

$y = \text{data.iloc[:, 1:3]}$  ← independent

## Now train & test



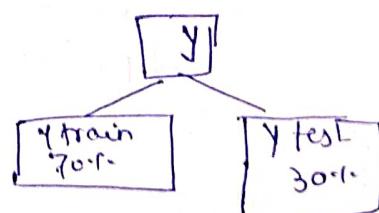
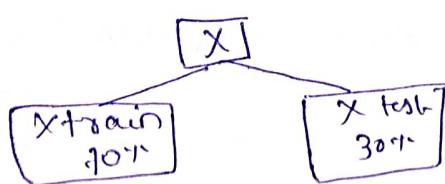
\* The data on which a model is built is called training data.



\* Once the model is trained, it is examined as to how well it has learnt using another subset of the original data which is called Test Data.

ideal split data ratio = 70 : 30

↑      ↓  
Train    Test



## ⇒ Feature Scaling

↳ Scaling the feature variables (ie all the independent variable) into the same range.

\* Basically sometimes the variation in magnitude or range may be present in dataset.

- ↳  $\downarrow$  problem it can create
  - Variable with higher magnitude will have more impact compare to one which is smaller.
  - Model might not predict properly because of unequal weight.

↳ Gradient descent algorithm that is used to find the coefficient of linear algo regression may take a long time to converge.

- ∴ with Feature scaling this problems will overcome
- ↳ The variables are scaled to have similar magnitude or range so that model is not biased towards a particular variable.
- ↳ Required for algo like → Linear regression, Logistic Regression, k Nearest Neighbour etc
- ↳ Not required for tree based Algo like → Random Forests, Decision tree etc.

How to perform?

- ↳ Standardization       $\left\{ \begin{array}{l} \text{Standard Normal Distribution} \\ \rightarrow \text{mean} \rightarrow 0 \\ \text{SD} \rightarrow 1 \end{array} \right\}$
- ↳ Min Max scaling       $\left\{ \begin{array}{l} \text{lies b/w [0,1]} \\ \text{Range -1 to 1} \end{array} \right\}$
- ↳ Normalisation       $\left\{ \begin{array}{l} \text{Range -1 to 1} \end{array} \right\}$

$$\text{Stand.} \rightarrow \frac{x - \mu(\text{Mean})}{\sigma(\text{SD})} \quad \left| \begin{array}{l} \text{Min max} \rightarrow \frac{x - \min(x)}{\max(x) - \min(x)} \\ \text{Nor} \rightarrow \frac{x - \text{mean}}{\max(x) - \min(x)} \end{array} \right.$$

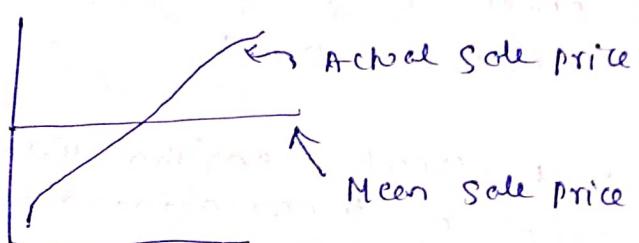
## Linear - Regression (AI90 - 1)

↳ Before moving, first understand some common terms  
↳ error.

(i) Build your first predictive model with mean prediction.

↳ Mean will come out to be same for all rows.

Hence mean sale price is same for even higher price house or lower price



To improve this model, take grade wise mean.

Now plot residual plot for both.

(Mean-scale, Mean-grade-scale)

↓

More scatter

Less scatter

→ Model evaluation Metrics

(calculating errors)

### (1) Mean error

↳ Mean of whole model -- but but but it will come out to be zero -- why? because some value are +ve & some are -ve, hence overall zero.

↓ To overcome we have absolute error

## ② Mean Absolute Error.

Take the absolute value of each individual error term, sum it up ~~all~~ for all the data points & then take a mean.

$\hat{y}$  = predicted value

$y$  = Actual value

$n$  = Total number of data points.

$$\rightarrow \text{MAE} = \frac{1}{n} \times \sum_{i=0}^n |\hat{y}_i - y_i| \quad \{ \text{or } |\hat{y} - y| \}$$

import mean\_absolute\_error.

## ③ Mean square error (MSE)

↳ same concept, but greater penalties for higher values.

$$\text{MSE} = \frac{1}{n} \times \sum_{i=0}^n (\hat{y}_i - y_i)^2$$

## ④ Root mean square error (RMSE)

↳ gives the expected average distance b/w the actual points & the prediction.

$$\text{RMSE} = \sqrt{\text{MSE}}$$

Now

$$\boxed{\text{MAE} = 0} \rightarrow \text{100% accurate}$$

$$\boxed{\text{MAE} = 250,000} \rightarrow \text{How much % accurate}$$

∴ To describe it we have concept of  $R^2$ .

$\Rightarrow R^2$

It give us the relative error of a regression model w.r.t  
the simple mean regression model.

$\bar{y}$  = Is the mean of the "Sale Price"

$\hat{y}_i$  = Is the predicted values from a regression model.

$y_i$  = represent actual "Sale-Price"

$$R^2 = 1 - \frac{\sum_{i=0}^n (\hat{y}_i - y_i)^2}{\sum_{i=0}^n (\bar{y}_i - y_i)^2}$$

Represent  $n \times (\text{MSE of regression model we want to evaluate})$   
 }  
 Represent  $n \times \text{MSE of simple mean model.}$

$$R^2 = 1 - \frac{\text{MSE(model)}}{\text{MSE(mean)}}$$

$\text{MSE(model)}$
$\text{MSE(mean)}$

$\rightarrow$  The lower this ratio is, the better is the model is.

When  $\hat{y} = y$ , then  $R^2 = 1 \rightarrow 100\%$  or perfect model

If  $\hat{y} = \bar{y}$ , then  $R^2 = 0 \rightarrow$  For mean prediction value  $R^2 > 0$

\* For any model which is better than mean prediction model but is not a perfect model.

$$0 < R^2 < 1$$

$$0 < R^2 < 100\%$$

(16)

## In reality

$$-\infty < R^2 < 1$$

↓

model we are trying to evaluate is even worse than the mean regression model.

- \*  $R^2$ : The degree to which your input variable explain the variation of your target or output variable.

$R^2 = 0.8 \rightarrow 80\%$  of the variable in the output

variable is explained by the input variable

and so on for the other input variables

- \* The higher the  $R^2$  is, the more is the variation explained by the input variable & hence better is the model.

- \* As more input variable are added to the regression model, the value of simple  $R^2$  can no longer be relied upon.

- \* The best linear regression model is the one that has the least error.

## Parameter of Linear regression

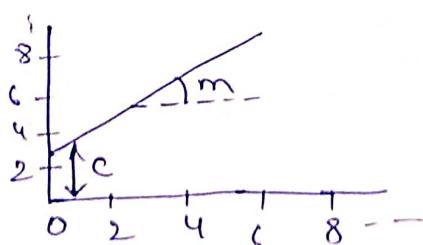
$$y = mx + c \quad \text{where } m, c \rightarrow \text{Parameters of lines.}$$

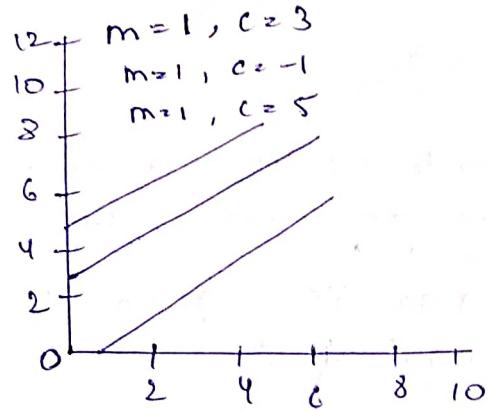
$y$  = Dependent variable

$x$  = Independent variable

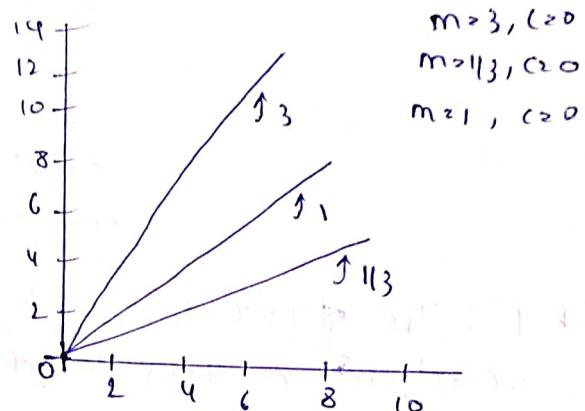
$m$  = Slope

$c$  = Intercept (what point the line cuts the Y axis)





Value of 'c' increases, line shifts upwards but angle of line does not change.



As the value of 'm' increases, the rate at which the value of  $y$  increases with respect to  $x$  also increases.

## Gradient Descent

It is an optimization algorithm that works iteratively & aims to find the minimum value of a convex function with respect to a set of parameters.

↳ Just like hit & trial.

→ Following steps occurs.

**Step-1:** Random initialization

↳ Take random value of 'm' slope & 'c' intercept.

To avoid everything from scratch, we take value of slope  $m = 0.1$

$c = \text{mean of "target value"}$

**Step-2:** Generate Prediction

$$\hat{y} = mx + c$$

(using initialize value)

Step-3 : Calculating the cost

$$T = \sum_{i=1}^n \frac{(\hat{y}_i - y_i)^2}{n} \quad \begin{matrix} \{ \text{wrt actual} \\ \text{point} \} \end{matrix} \rightarrow \text{MSE}$$

$$T = \sum_{i=1}^n \frac{(m x_i + c - y_i)^2}{n}$$

Step 4 : Updation of parameters

$$\{ m_{\text{new}} = m_{\text{old}} - z_m \}$$

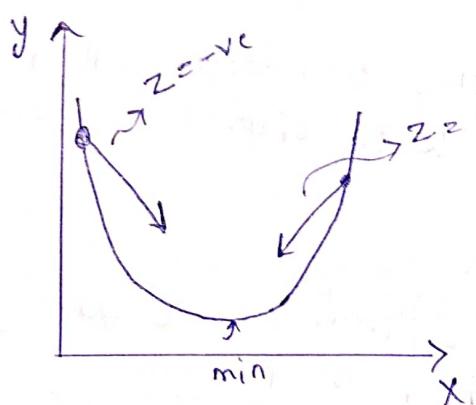
$[z \rightarrow \text{still unknown}]$   
but can be positive  
or -ve]

$$\{ C_{\text{new}} = C_{\text{old}} - z_c \}$$



If  $z$  is +ve, parameter will decrease

If  $z$  is -ve, parameter will increase



BUT,  
how to know on which side of the minimum cost we have  
on & what value should be we take for  $z$ ?

↓ so in this

Partial differentiation

$$\frac{\partial J(T)}{\partial m}$$

∴ partial derivatives of the cost function wrt  $m$  &  $c$

$$G_m = \frac{\partial J(T)}{\partial m} = 2 \sum_{i=1}^n \frac{(x_i m + c - y_i) x_i}{n}$$

$$G_c = \frac{\partial J(T)}{\partial c} = 2 \sum_{i=1}^n \frac{(x_i m + c - y_i)}{n}$$

$$\therefore Z_m = \alpha G_m \text{ & } Z_c = \alpha G_c$$

$$m_{\text{new}} = m_{\text{old}} - \alpha G_m$$

[ $\alpha \rightarrow$  Gradient]

$$c_{\text{new}} = c_{\text{old}} - \alpha G_c$$

$\alpha$  = learning rate = a constant

$\alpha \rightarrow$  High  $\Rightarrow$  cost function may explode & bounce out far from the minima.

$\alpha \rightarrow$  Low  $\Rightarrow$  parameters will take lot of iteration to converge to the optimal values

If  $\alpha = 0.01$

previous value of the cost function.

New value of the cost function

The threshold value (usually  $10^{-5}$ )

The algorithm is converge

## Logistic Regression (Algo-2)

Classification Algo

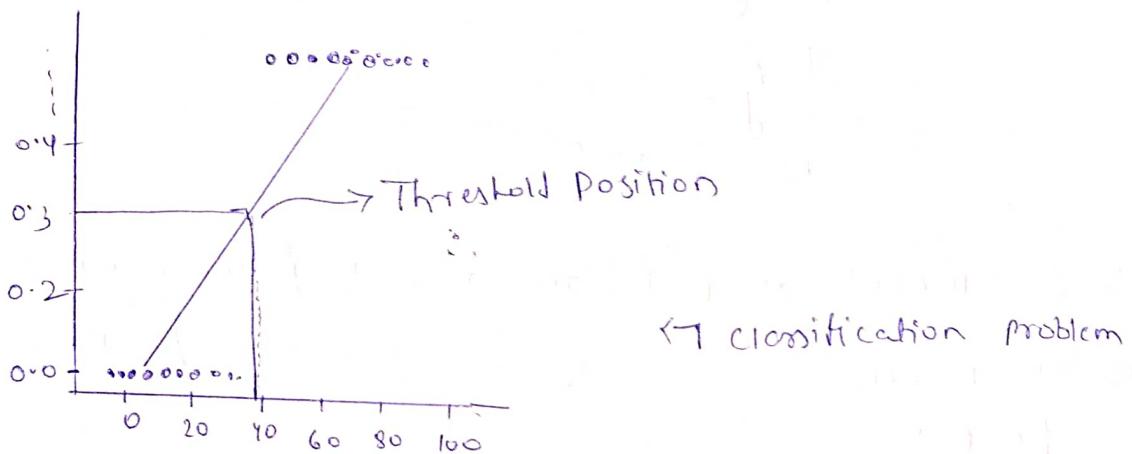
A categorical target variable

is used.

A categorical variable is one which represent a characteristic which cannot be measured or counted.

Q) If logistic regression is a classification algo., then why is there the term "regression" in its name.

↳ { Think, Think, Think -- will answer afterwards }



? why not linear regression here

- Susceptible to the outlier.
- Interpretation of the model.

→ If outliers present  
then threshold point can be shifted  
ISSUES

Here we will be working with "class" categories hence its probability lies in a range of 0 & 1.

Eqn for linear regression is

$$[Z = mx + c]$$



In order to restrict this value b/w 0 & 1 we use Logit (or Sigmoid) function.

The Logit function

↳ A func. which have a range b/w 0 & 1, irrespective of any input, it will always give output b/w 0 & 1.

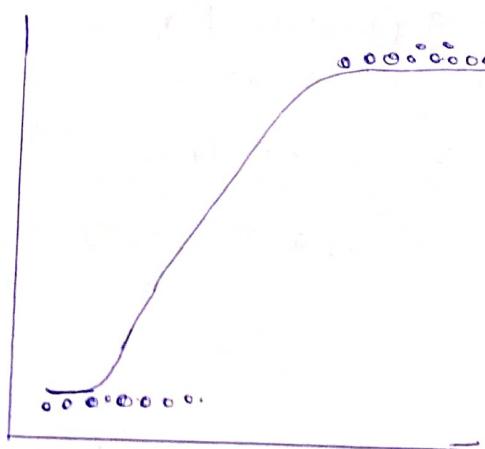
Given as:

$$g(x) = \frac{1}{1 + e^{-x}}$$

$$\therefore \hat{y} = g(z)$$

$$\hat{y} = \frac{1}{1 + e^{-(mx+c)}}$$

- If  $(mx+c)$  is high then overall  $\hat{y}$  comes out to be 1
- If  $(mx+c)$  is less (or -ve) then overall  $\hat{y}$  comes out to be 0.



S' shaped graph

Hence the main reason for using term "regression" is:  
→ Prediction continuous probability values b/w 0 & 1  
↳ Regression.

Eg:

probability of 0.82 → means he belongs to class-1.

→ Loss function

↳ In this tech. we use log loss function.

$$J = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 \quad \text{For linear regression}$$

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

where

$\hat{y}$  = predicted probability of class 1

y = Actual Class

n = Number of observation

\* → probability b/w 0 & 1

→ Middle point - 0.5 (threshold)

→ If probability > 0.50, the observation is regarded

as class 1

→ If probability < 0.50, the observation is regarded

as class 0.

## → Evaluation Metrics

### ① Confusion Matrix

It is used to interpret the model predictions systematically.

→ It is  $n \times n$  matrix,  $n \rightarrow$  no. of distinct classes in the target variable.

→ For  $2 \times 2$  matrix

↳ Class 1: Positive class

Class 0: Negative class

		Prediction outcome	
		positive	negative
Actual Value	positive	TP	FN
	Negative	FP	TN

\* Confusion matrix serves as the basic platform of representation for most of the classification metrics.

\* The metrics which derives from the confusion matrix are:

① Accuracy

② precision

③ Recall

\* Precision & recall handle the imbalanced dataset efficiently.

$$\text{precision} = \frac{\text{Predictions Actually Positive}}{\text{Total Predicted Positive}}$$

$$\Rightarrow \frac{TP}{TP + FP}$$

Precision :

It is used when avoiding false positives is more essential than encountering false negative.

Eg:

4	2
8	486

$$\Rightarrow \text{precision} = \frac{4}{4+8} \approx 33\%$$

∴ Model having 98% is only 33% precise

$$\text{Recall} = \frac{\text{Predictions Actually Positive}}{\text{Total Actual Positive}}$$

$$\Rightarrow \frac{TP}{TP + FN}$$

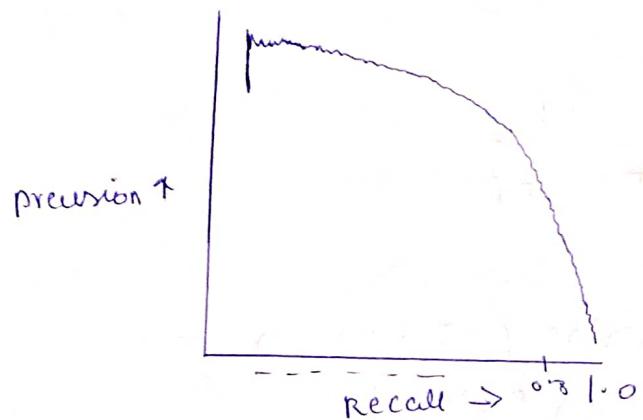
From Above Eg

$$\Rightarrow \frac{4}{4+2} \approx 66\%$$

only 66% of actual +ve value were predicted as +ve.

## ∴ Recall

It is used when avoiding false -ve is prioritized over encountering false +ve.



If precision ↑ then  
Recall ↓.

But if we are not sure from graph (or method) then  
we use F1 metrics which is nothing but harmonic  
mean of precision & Recall.

$$F1 = \frac{2}{\frac{1}{\text{precision}} + \frac{1}{\text{recall}}}$$

Max. when precision = Recall.

? But if precision, Recall or F1 score yield the same result, then how do we distinguish which one is better.

↓  
We use log loss

↓  
Cost fun. of Logistic Regression.

\* The smaller the value of log loss, the better the model.

⇒ The farther the points from actual probability (0 & 1) the more is the loss function

⇒ AUC-ROC Curve

It is a performance measurement for classification problems at various threshold settings.

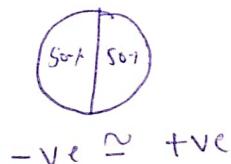
AUC - Area Under Curve

ROC - Receiver Operating Characteristics

\* Higher the AUC, the better the model is distinguishing between 0 and 1.

(AUC - ROC often used interchangeably)

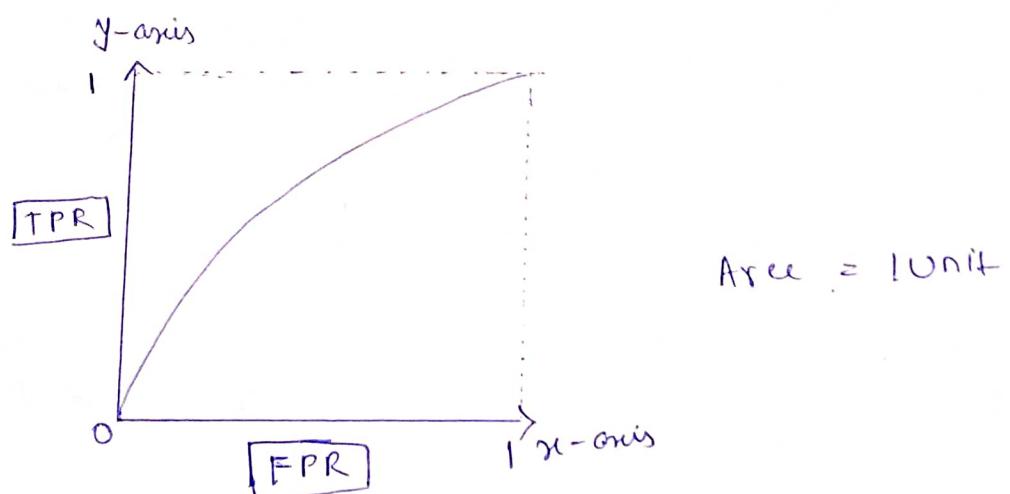
\* AUC-ROC works well for balanced data set.



\* AUC-ROC is not suitable for imbalanced data.

-ve > +ve

(90%) (10%)



$$FPR \text{ (False pos rate)} = \frac{FP}{TN+FP}$$

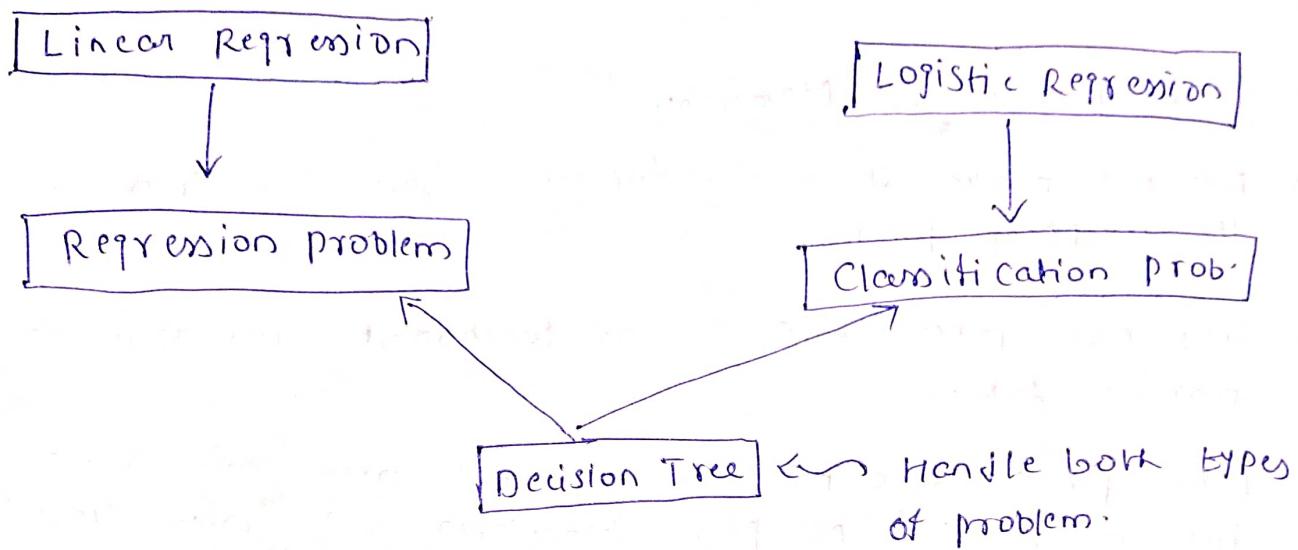
$$TPR \text{ (True pos rate)} = \frac{TP}{TP+FN}$$

\* on this method, create a plot at different threshold & calculate the area.

The higher the AUC-ROC, better is the performance of model.

\* Remember ↓  
If  $AUC-ROC > 0.95$ , there could be something wrong with the model or the dataset.

## Decision-Tree (AI10-3)



Linear & Logistic Regression

: Make restrictions or assumptions.

Decision tree : does not make any assumptions.

↳ Predictive model

↓ 2 types:

① Parametric

② Non parametric

① Parametric

→ Make strong assumptions about the form of the mapping function.

→ very simple & interpretable.

→ set of parameters does not depend upon the amount of data.

\* A learning model that summarizes data with a set of parameters of fixed size (independent of the number of training examples) is called parametric model.

Linear req. → parametric Algo.  
 Logistic req. → non-parametric Algo.

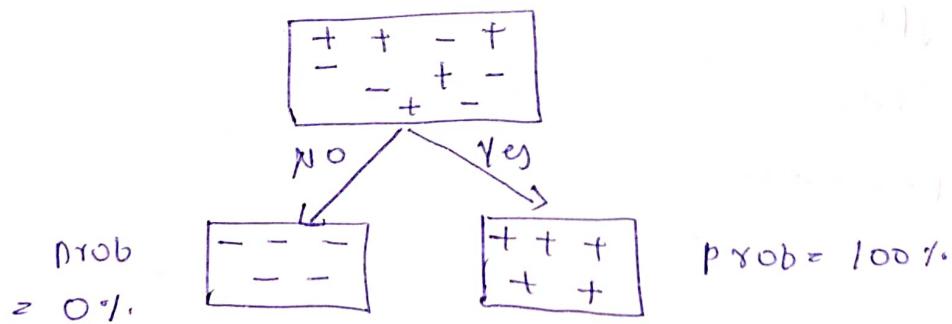
## ⑥ Non parametric Algorithm

- Do not make strong assumptions about the form of the mapping function.
- They are free to learn any functional form from the training data.
- Non parametric methods are good when you have a lot of data & no prior knowledge and when you don't want to worry too much about choosing just the right features.

## → Decision Tree

Decision to separate the classes

All the +ve to be one side & all the -ve on the other side.

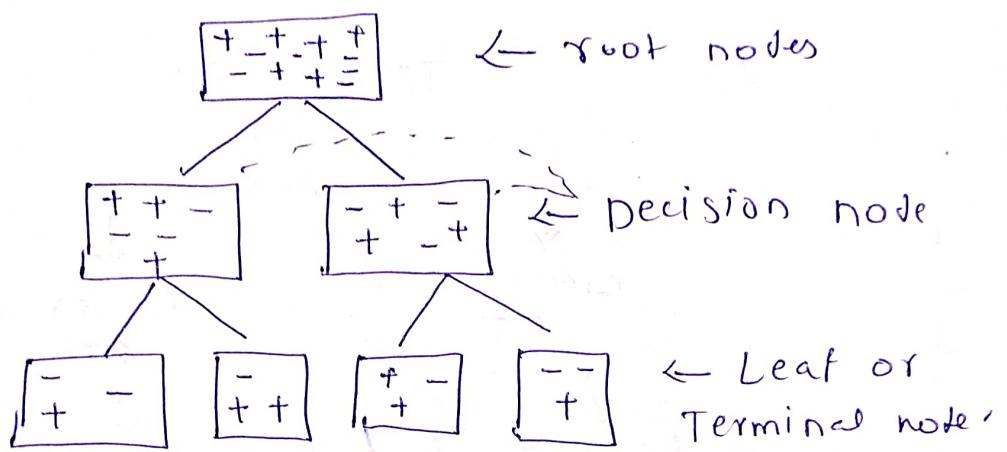


\* An ideal case - [Pure node] after split

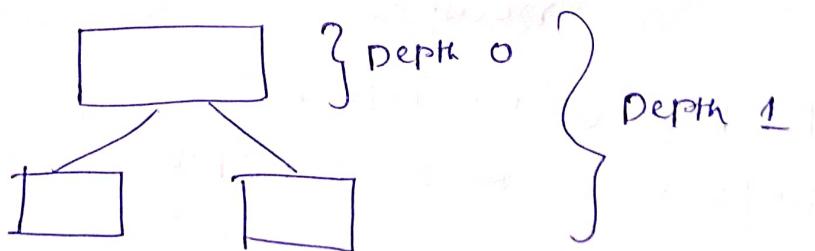
A pure node is in which all the data points exhibits the desired behaviour or have the same class.

The objective of decision tree: Is to have pure nodes (or close to pure nodes)

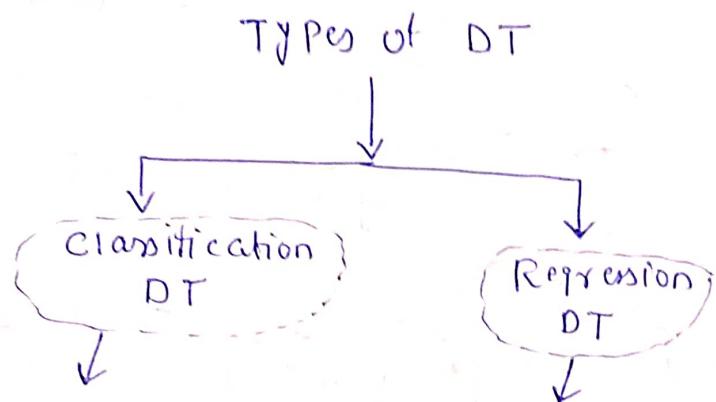
## → Various terminology



- (i) **Splitting** : A process of dividing a node into 2 or more nodes.
- (ii) **Decision Node** : A sub node that is further divided into sub nodes. (Initial Sub Node)
- (iii) **Leaf Node** : A node which does not split further.
- (iv) **Parent Node** : A node which is divided into sub nodes.
- (v) **Child Node** : A sub node of a parent node.
- (vi) **Depth of tree** : The length of the longest path from the root node till the leaf node.



## → Logic behind DT



Target Variable  
is categorical  
in nature.

Target Variable is  
continuous in nature.

↳ Criteria to be used to split a decision node

- (i) Gini impurity
- (ii) Information gain using entropy.

### (i) Gini impurity

↳ It measures the impurity of node.

$$GI = 1 - \text{gini}$$

(gini) → It is the probability, whose ranges are from 0 to 1.

↳ Higher the value, the more will be the purity of the node. Lesser the value means lesser purity nodes.

### Property of Gini impurity

- (a) Helps to decide the best split.

$$\text{GI} = 1 - \frac{\text{impurity}}{\text{purity}}$$

- ⑥ works only for categorical variable & not for continuous variable.
- ⑦ Always performs a binary split.

↳  $\boxed{\text{Gini impurity} = 1 - \text{gini}}$

$$\text{gini} = (P_1^2 + P_2^2 + P_3^2 + \dots + P_n^2)$$

↓

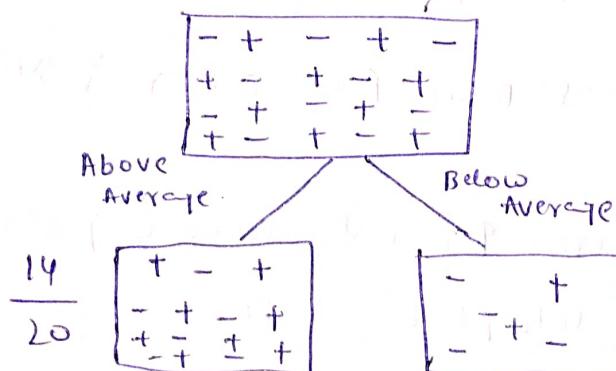
The probability of any 2 random data points belonging to class 1.

→ {For complete tree}

→ Use the weighted impurity of both sub nodes.

Weight of a node =  $\frac{\text{number of samples in that node}}{\text{Total number of samples in its parent node.}}$

Eg:



$$\text{Prob Play} = 0.57$$

$$\text{Prob not Play} = 0.43$$

$$\text{Prob Play} = 0.33$$

$$\text{Prob not Play} = 0.67$$

$$gini(\text{sub node}) (\text{Above Ave.}) = 1 - [(0.57)^2 + (0.43)^2] = 0.49$$

$$gini(\text{sub node}) (\text{Below Ave.}) = 1 - [(0.33)^2 + (0.67)^2] = 0.44$$

weighted gini impurity : performance in class

$$\left(\frac{14}{20}\right) \times 0.49 + \left(\frac{6}{20}\right) \times 0.44 = 0.475$$

Now check for other split like this.  
And choose the tree with less G.I.

### (ii) Information gain using entropy

**Information gain**: The difference b/w information needed to describe the parent node & information needed to describe the children node.

\* The more is the homogeneity of the nodes, the more is the information gain.

∴  $\uparrow I_G \rightarrow \text{homogeneity} \uparrow$

$$I_H = 1 - \text{Entropy}(\text{subnodes})$$

$$\text{Entropy} = -P_1 \log(P_1) - P_2 \log(P_2) - P_3 \log(P_3) - \dots - P_n \log(P_n)$$

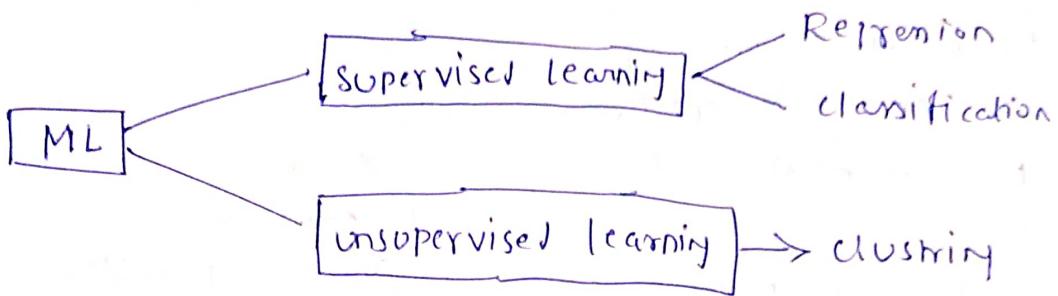
\* Lower entropy means more pure nodes & thus more information gain.

\* It also only works for categorical target variable.

↳ Steps to calc. entropy of a split:

- cal. the entropy of the parent node.
- cal. the entropy of each child node.
- Cal. the weight average entropy of the split.
- If the weighted entropy of the child node is greater than the parent node, ignore the split.
- Less entropy (or higher  $I_G$ ) leads to more homogeneity.

## ⇒ Introduction to clustering.



\* Supervised → Develop predictive model based on both input & output data.

\* unsupervised → Group & interpret data based only on input data.

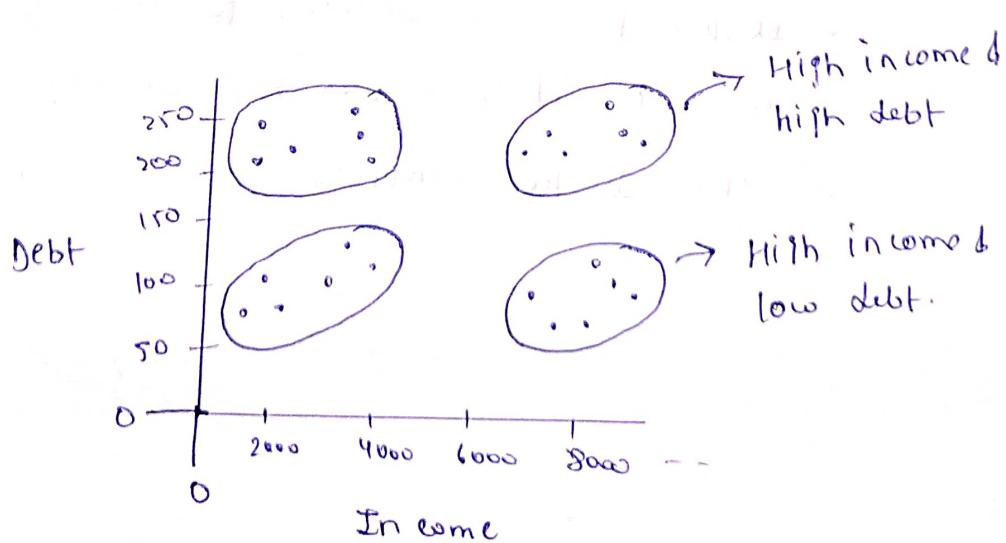
Clustering ; we try to find patterns in the observation based on their features are found.

; Target variable are not defined.

→ Properties of clustering.

→ All the data points in a cluster should be similar to each others.

→ The data points from different clusters should be as different as possible.



## → Application of clustering

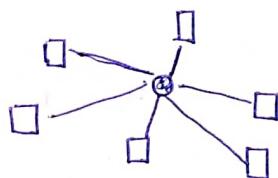
- Customer segmentation.
- Document clustering.
- Image Segmentation.
- Recommendation Engine.

{ because visualising & prediction  
is not possible}

## → Evaluation metrics

### ① Inertia

It tells us how far the points within a cluster are.



The distance of a point from the centroid within the cluster is called Intra cluster distance.

∴ Inertia gives us sum of intra cluster distance.

\* Lesser the inertia value, better the model is.

$$\therefore \text{Inertia} = \sum_{i=1}^m \sum_{k=1}^K w_{ik} (||x_i - \mu_k||)^2$$

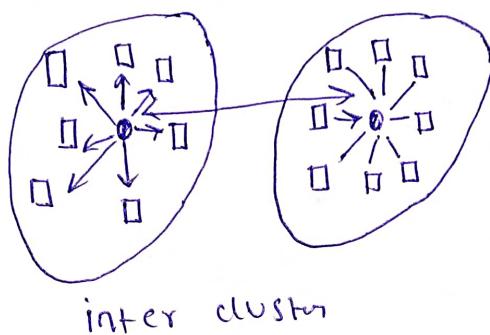
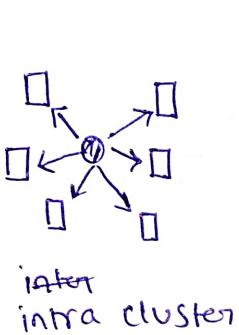
→  $w_{ik} = 1$ , if  $x_i$  belongs to cluster  $k$

→  $w_{ik} = 0$ , if  $x_i$  does not belong to cluster  $k$ .

→  $\mu_k$  is the centroid of  $x_i$ 's cluster.

## ② Dunn index

↳ use to evaluate the clusters.



$$\text{Dunn index} = \frac{\min(\text{inter cluster distance})}{\max(\text{intra cluster distance})}$$

∴ The higher the value of the Dunn index, the better the cluster are.