

# Classification Tree

## Classification:

Classification is the action of classifying something according to shared qualities or characteristics.

In classification world....

- To explain is called as profiling.
- To predict is called as classifying.

| Exp | state | Qualification | Hire |
|-----|-------|---------------|------|
| 5   | TN    | B.Tech        | Yes  |
| 5   | AP    | B.Tech        | Yes  |
| 10  | UP    | Metric        | No   |
| 5   | UP    | Ph.D          | No   |
| 7   | TN    | B.Teh         | Yes  |

→ profiling  
and  
Training

## New Incoming data

| Exp | state | Qualification | Prediction |
|-----|-------|---------------|------------|
| 5   | TN    | B.Teh         | Yes        |
| 8   | UP    | Metric        | No         |

→ classifying  
and  
prediction

{Independent variable} → {Predictor, dimensions, columns, features, Input,  $x_1 - x_2 - x_3 - x_4 - \dots$ }

{Dependent variable} → {Predicted, output, Target, Y}

# Popular Classification Techniques:

\* CHAID → Chi-squared Automatic Interaction Detector: CHAID is not a binary classification tree. It is used for categorical variables.

(\*) CART → Classification And Regression Tree.  
CART is a binary decision tree.

\* C4.5 → It uses the concept of information entropy.

→ Both CART and C4.5 can handle categorical and numerical data. CART is most widely used algorithm.

CART: - (classification and regression tree):

CART is a classification algorithm that uses GINI as the measure of impurity.

GINI Index is the measure of impurity.

GINI Impurity / Gini Index: It tells us the probability of misclassification of an observation.

Gini at a node:

$$\text{Gini}(t) = 1 - \sum_j [P(j|t)]^2 \quad \begin{cases} P(j|t) \text{ is the relative frequency of class } j \text{ at node } t. \end{cases}$$

Gini of a split:  $[2.0 \times (0.1)] + [2.0 \times (0.1)]$

$$\text{Gini(split)} = \sum_{i=1}^k \frac{n_i}{n} \cdot \text{Gini}(t_i) \quad \begin{cases} n_i = \text{no. of records at child } i, \\ n = \text{Total no. of records in parent node} \end{cases}$$

(\*) Gini gain = Gini(t) - Gini(split)

- Lower the Gini (split) better means lower the chances of misclassification.
- We can also say higher the gini gain greater the split priority will be.

calculation:

| Gender | Occupation    | Age | target |
|--------|---------------|-----|--------|
| M      | salary        | 22  | 1      |
| M      | salary        | 22  | 0      |
| M      | self-employed | 23  | 1      |
| M      | self-employed | 23  | 0      |
| M      | self-employed | 24  | 1      |
| M      | self-employed | 24  | 0      |
| F      | salary        | 25  | 1      |
| F      | salary        | 25  | 0      |
| F      | salary        | 26  | 0      |
| F      | self-employed | 26  | 0      |

| Node Gender    | Gini computation Formula   | Gini Index |
|----------------|--|------------|
| Overall        | $1 - \left[ \left( \frac{4}{10} \right)^2 + \left( \frac{6}{10} \right)^2 \right]$                               | 0.48       |
| Gender (M)     | $1 - \left[ \left( \frac{3}{6} \right)^2 + \left( \frac{3}{6} \right)^2 \right]$                                 | 0.5        |
| Gender (F)     | $1 - \left[ \left( \frac{1}{4} \right)^2 + \left( \frac{3}{4} \right)^2 \right]$                                 | 0.25       |
| Gender (split) | $\left[ \left( \frac{6}{10} \right) \times 0.5 \right] + \left[ \left( \frac{4}{10} \right) \times 0.25 \right]$ | 0.45       |
| Gini Gain      | $0.48 - 0.45$  | 0.03       |

| Node occupation            | Gini computation formula  | Gini Index |
|----------------------------|---|------------|
| overall                    | $1 - \left[ \left(\frac{4}{10}\right)^2 + \left(\frac{6}{10}\right)^2 \right]$                            | 0.48       |
| occupation (salary)        | $1 - \left[ \left(\frac{2}{5}\right)^2 + \left(\frac{3}{5}\right)^2 \right]$                              | 0.48       |
| occupation (self-employed) | $1 - \left[ \left(\frac{2}{5}\right)^2 + \left(\frac{3}{5}\right)^2 \right]$                              | 0.48       |
| occupation (split)         | $\left[\left(\frac{5}{10}\right) \times 0.48\right] + \left[\left(\frac{5}{10}\right) \times 0.48\right]$ | 0.48       |
| Gini Gain                  | $0.48 - 0.48 = 0.00$  | 0.00       |

Since the gender has high gini gain and

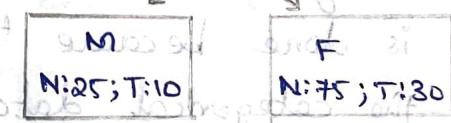
low gini split when compared to occupation, Gender has the high priority and low mis-classification and it takes the first split (root node).

Problem:

{Find the gini gain for }  $\Rightarrow$    
 {the given tree split}

$$\text{Overall: } 1 - \left[ \left(\frac{60}{100}\right)^2 + \left(\frac{40}{100}\right)^2 \right] \Rightarrow 0.48$$

Root Node  
 $T: 100; T: 40$



{Gender} {Male} :  $1 - \left[ \left(\frac{10}{25}\right)^2 + \left(\frac{18}{25}\right)^2 \right] \Rightarrow 0.48$

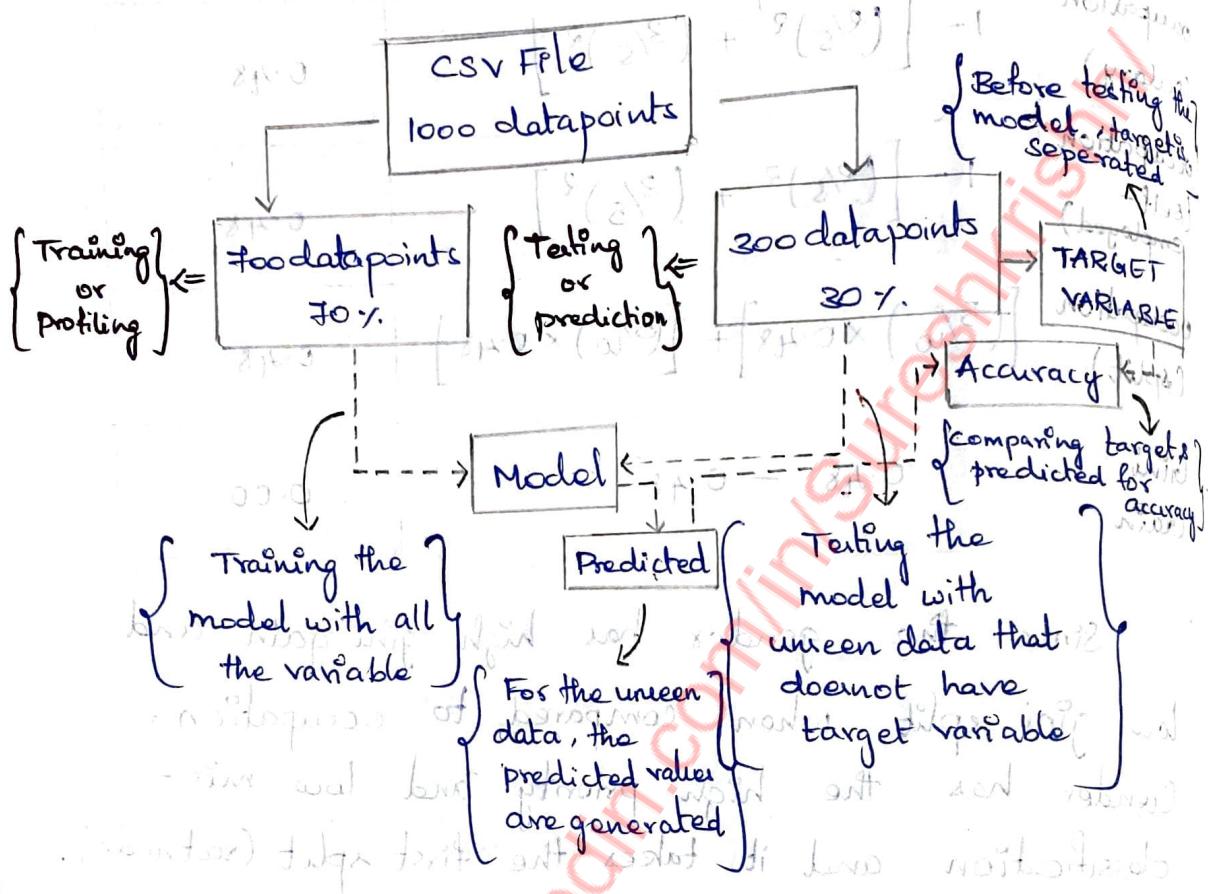
{Gender} {Female} :  $1 - \left[ \left(\frac{30}{75}\right)^2 + \left(\frac{45}{75}\right)^2 \right] \Rightarrow 0.48$

{Gender} {Split} :  $\left[\left(\frac{25}{100}\right) \times 0.48\right] + \left[\left(\frac{75}{100}\right) \times 0.48\right] \Rightarrow 0.48$

{Gini Gain} :  $0.48 - 0.48 \Rightarrow 0.00$

# Training and Testing

Training and testing in CART is also called as profiling and prediction.



## Encoding: (one hot encoding)

One hot encoding is used to convert the categorical variable into quantitative variable. This is done because the ML model can understand the categorical data.

| Gender | Gender | Target | Target |
|--------|--------|--------|--------|
| Male   | 1      | No     | 0      |
| Female | 0      | Yes    | 1      |
| Female | 0      | Yes    | 1      |
| Female | 0      | No     | 0      |
| Male   | 1      | Yes    | 1      |
| Female | 0      | No     | 0      |
| Male   | 1      | No     | 0      |
| Male   | 1      | Yes    | 1      |

| Cust. | Occupation |
|-------|------------|
| C1    | self       |
| C2    | Non-self   |
| C3    | Govt       |
| C4    | self       |
| C5    | private    |
| C6    | self       |

| Cust | self | Non self | Govt | Private |
|------|------|----------|------|---------|
| C1   | 1    | 0        | 0    | 0       |
| C2   | 0    | 1        | 0    | 0       |
| C3   | 0    | 0        | 1    | 0       |
| C4   | 1    | 0        | 0    | 0       |
| C5   | 0    | 0        | 0    | 1       |
| C6   | 1    | 0        | 0    | 0       |

{ one-hot } converts categorical variable into 0's and 1's.

{ Label } : Assigns numerical values (0, 1, 2, ...) to the categorical variable. By this the converted values are considered as weights. (ie.  $0 < w_1 < w_2 < w_3 < \dots$ ). So it is not an effective method.

When to use one-hot encoding and when to use label encoding?

One hot encoding when:

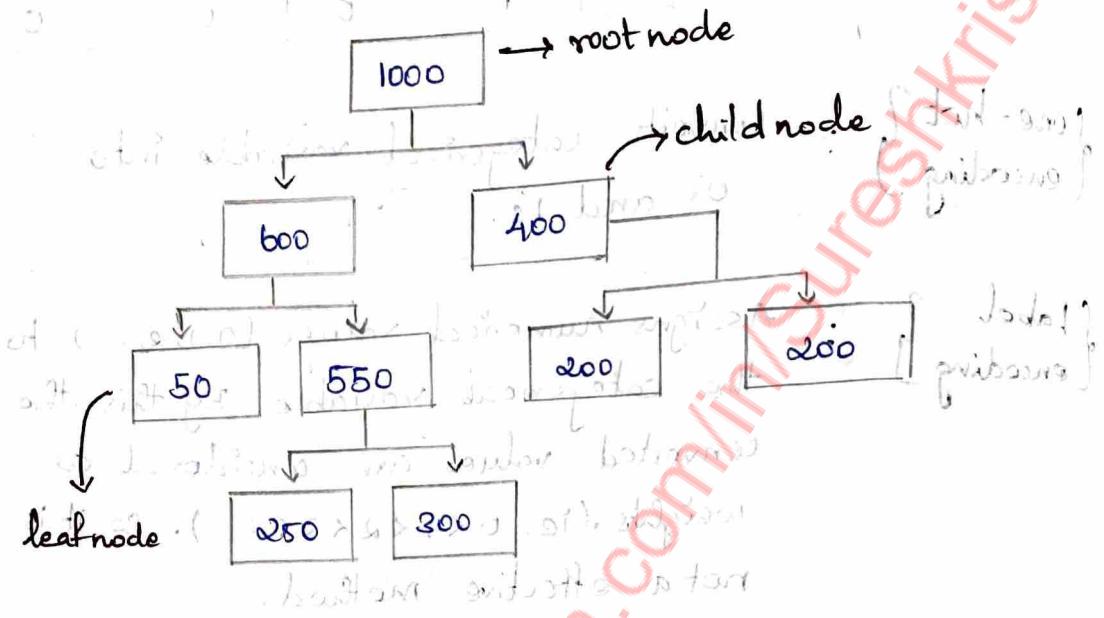
- \* The categorical feature is not ordinal. (eg: countries, color, job occupation etc.)
- \* When the no. of categorical feature is less.

Label encoding when:

- \* The categorical feature is ordinal (eg: rank, class, job type etc.,)
- \* When the number of categories is quite large.

Code to build CART tree:

```
from sklearn.tree import DecisionTreeClassifier
model_dt = DecisionTreeClassifier(max_depth=8, criterion='gini', min_samples_split=100, min_samples_leaf=10)
```



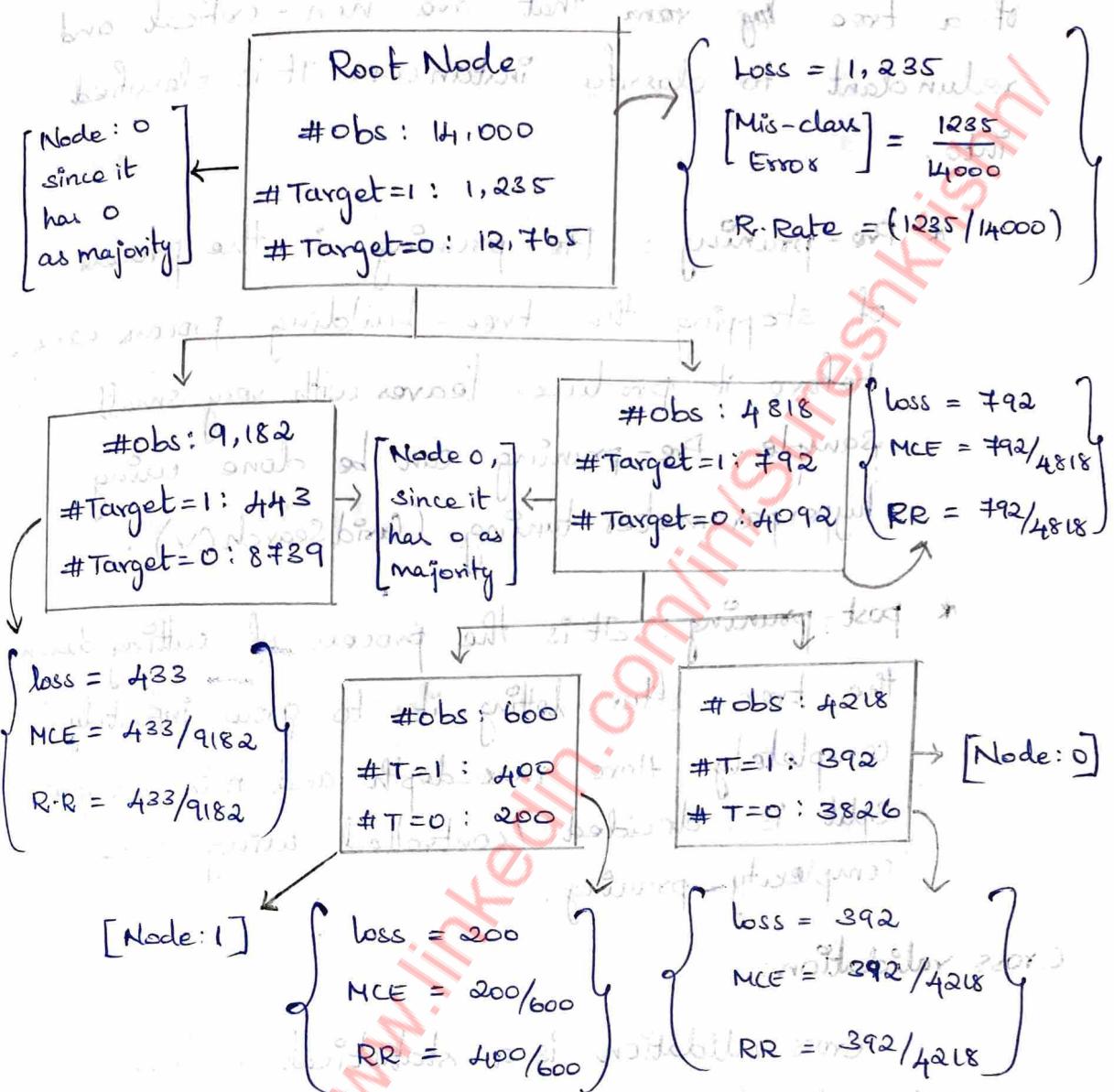
\* **min\_samples\_split**: The minimum number of observations that should present in a node in order to split any further.

\* **min\_samples\_leaf**: The minimum number of observations in any terminal leaf node (node that will be terminated if it contains such observations). If **min\_samples\_split** is only specified then **min\_samples\_leaf** will be **min\_samples\_split**/3. If **min\_samples\_leaf** is specified then **min\_samples\_split** will be **min\_samples\_leaf** × 3.

\* **max\_depth**: The maximum depth of a tree. If none is specified, the tree grows till all leaves are pure or until all leaves contain less than **min\_samples\_split**.

\* Criterion - 'gini' for gini impurity and 'entropy' for information gain.

Loss, Mis-classification error and response rate:



Loss : Number of cases mis-classified in the given node.

Mis-classification error: Ratio of total number of cases misclassified to the total number of observations at the node. (Have to find for full tree)

Response Rate: Ratio of total number of responders (Target=1) to the total number of observation at node.

Nodes with high response rate are always best and that what we are interested in about these

**Pruning:** Pruning is the process of reducing the size of decision tree by removing the sections of a tree that are non-critical and redundant to clarify instances. It is classified into,

\* **Pre-pruning:** Pre-pruning is the process of stopping the tree-building process early, before it produces leaves with very small sample. Pre-pruning can be done using hyperparameter tuning (Grid Search CV).

\* **post-pruning:** It is the process of cutting down the tree after letting it to grow infinitely/completely. Here max-depth and min-sample-split is decided/controlled using cost-complexity-pruning.

## Cross validation:

Cross validation is a statistical method used to estimate the skill of machine learning model on new (unseen) data. Common methods of cross validation are,

- \* K-Fold cross validation.
- \* stratified cross validation.

K-fold CV uses random sampling so the chance of bias in sampling is possible. stratified CV divides based on strata and does random sampling on each strata.

Typical value of k-fold cross validation is set to 10.  
 This number represents the number of folds.

| K Fold CV | P <sub>1</sub> | P <sub>2</sub> | P <sub>3</sub> | P <sub>4</sub> | P <sub>5</sub> | P <sub>6</sub> | P <sub>7</sub> | P <sub>8</sub> | P <sub>9</sub> | P <sub>10</sub> |
|-----------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|-----------------|
| Fold-1    | Train          | Test           | Test            |
| Fold-2    | Train          | Train          | Train          | Train          | Train          | Train          | Test           | Test           | Train          | Train           |
| Fold-3    | Train          | Train          | Train          | Train          | Test           | Test           | Train          | Train          | Train          | Train           |
| Fold-4    | Train          | Train          | Test           | Test           | Train          | Train          | Train          | Train          | Train          | Train           |
| Fold-5    | Test           | Test           | Train           |

$$\begin{aligned} \text{obs} &= 10 \\ k &= 5 \\ \text{Test per fold} &\text{ is equal} \\ &\text{to } \text{obs}/k \\ \text{obs}/k &= 10/5 = 2 \end{aligned}$$

Accuracy metrics:

\* Confusion matrix - Precision, Recall & F1 score

|             |                   | Actual      |                   | Predicted +ive rate (PPR) | $\frac{A}{(A+B)}$ |
|-------------|-------------------|-------------|-------------------|---------------------------|-------------------|
|             |                   | Positive    | Negative          |                           |                   |
| Predicted   | Positive          | A (TP)      | B (FP)            |                           |                   |
|             | Negative          | C (FN)      | D (TN)            | Predicted -ive rate (PNR) | $\frac{D}{(C+D)}$ |
| sensitivity |                   | specificity |                   |                           |                   |
|             | $\frac{A}{(A+C)}$ |             | $\frac{D}{(B+D)}$ |                           |                   |

(Precision)

→ (Recall) ⇒ True positive Rate

$$\text{Accuracy} = \frac{A + D}{(A+B+C+D)} \quad (\text{or}) \quad \frac{(TP + TN)}{(TP + FP + FN + TN)}$$

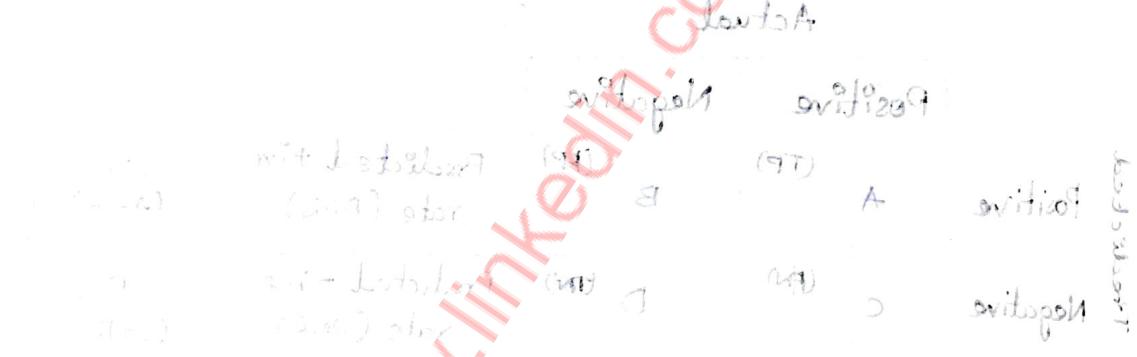
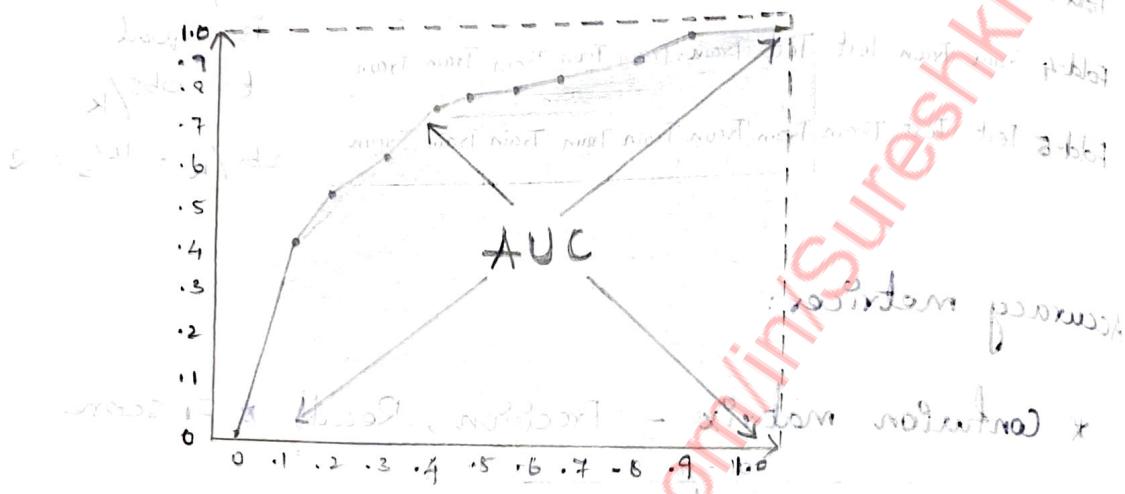
$$\text{F1-score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$\text{False positive Rate} = 1 - \text{specificity}.$$

\* Area under the curve: (AUC)

Area under the curve shows how many mistakes (wrong prediction) are made by the model before it makes correct prediction.

For AUC, False positive rate is plotted against the true positive rate.



$$\text{AUC} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$\Rightarrow \text{AUC} = \frac{\text{TP}}{\text{TP} + \text{FN}}$

$$\text{AUC} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{\text{TP}}{\text{TP} + \text{FN} + \text{TN} + \text{FP}}$$

$$\text{AUC} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{\text{TP}}{\text{TP} + \text{FN} + \text{TN} + \text{FP}}$$

$\Rightarrow \text{AUC} = \frac{\text{TP}}{\text{TP} + \text{FN}}$

## Random Forest

Random forest is a supervised learning classification algorithm. The forest that it builds is an ensemble of decision tree, usually trained with a bagging method. Feature importance can be easily measured using random forest.

**Ensemble:** Ensemble technique uses multiple learning algorithms to obtain better predictive performance.

Eg: Random forest is built using multiple decision trees.

### Bootstrap Aggregation :- (Bagging):

Bagging is the process of creating multiple decision trees from the given training set. For each of the decision tree the observations are picked by doing random sampling with replacement.

- \* In decision tree algorithm we have fixed proportion of target variable. But in random forest due to bagging each decision tree has a varying proportion.
- \* The models are fitted using the bootstrap samples (decision tree) and combined by averaging the voting. Eg: For a test datapoint if the output proportion is predicted as 1 in 700 trees out of 1000 trees, then the final predicted output will be 1.

\* Individual trees in the random forest may have high over-fitting but averaging the output removes the bias. Hence randomforest helps reduce over-fitting.

**RandomForest is a Blackbox Technique**

Inside the Random Forest black box:

**Step - 1: \* Random sampling with Replacement.**  
Here we create "n" number of trees by using sampling with replacement.

| ID    | $x_1$ | $x_2$ | $x_3$ | $x_4$ | ... | y |
|-------|-------|-------|-------|-------|-----|---|
| $R_1$ |       |       |       |       |     |   |
| $R_2$ |       |       |       |       |     |   |
| $R_3$ |       |       |       |       |     |   |
| .     |       |       |       |       |     |   |
| $R_4$ |       |       |       |       |     |   |
| .     |       |       |       |       |     |   |
| $R_n$ |       |       |       |       |     |   |

$\Rightarrow$  Training set

| ID       | $x_1$ | $x_2$ | $x_3$ | $x_4$ | ... | y |
|----------|-------|-------|-------|-------|-----|---|
| $R_{10}$ |       |       |       |       |     |   |
| $R_3$    |       |       |       |       |     |   |
| $R_{15}$ |       |       |       |       |     |   |
| .        |       |       |       |       |     |   |
| $D_1$    |       |       |       |       |     |   |

| ID       | $x_1$ | $x_2$ | $x_3$ | $x_4$ | ... | y |
|----------|-------|-------|-------|-------|-----|---|
| $R_1$    |       |       |       |       |     |   |
| $R_1$    |       |       |       |       |     |   |
| $R_{12}$ |       |       |       |       |     |   |
| .        |       |       |       |       |     |   |

$D_2$

$D_{...}$

| ID       | $x_1$ | $x_2$ | $x_3$ | $x_4$ | ... | y |
|----------|-------|-------|-------|-------|-----|---|
| $R_1$    |       |       |       |       |     |   |
| $R_1$    |       |       |       |       |     |   |
| $R_{12}$ |       |       |       |       |     |   |
| .        |       |       |       |       |     |   |

| ID       | $x_1$ | $x_2$ | $x_3$ | $x_4$ | ... | y |
|----------|-------|-------|-------|-------|-----|---|
| $R_1$    |       |       |       |       |     |   |
| $R_1$    |       |       |       |       |     |   |
| $R_{12}$ |       |       |       |       |     |   |
| .        |       |       |       |       |     |   |

Decision tree (first iteration)

$D_1, D_2, \dots, D_n \Rightarrow$  Decision tree created by bagging.

step-2: \* Building the tree for each sample with only partial set of 'm' variables (predictors) being considered at each node.

\*  $m < M$ , where  $M =$  total # of predictor variable

| ID              | $x_1$ | $x_5$ | $x_9$ | Y |
|-----------------|-------|-------|-------|---|
| R <sub>10</sub> |       |       |       |   |
| R <sub>3</sub>  |       |       |       |   |
| R <sub>15</sub> |       |       |       |   |
| .               |       |       |       |   |
| .               |       |       |       |   |

→ {partial list of variable considered for splitting based on best variable from partial list.}

| ID              | $x_2$ | $x_6$ | $x_9$ | Y |
|-----------------|-------|-------|-------|---|
| R <sub>10</sub> |       |       |       |   |
| R <sub>3</sub>  |       |       |       |   |
| R <sub>15</sub> |       |       |       |   |
| .               |       |       |       |   |
| .               |       |       |       |   |

| ID              | $x_1$ | $x_4$ | $x_8$ | Y |
|-----------------|-------|-------|-------|---|
| R <sub>10</sub> |       |       |       |   |
| R <sub>3</sub>  |       |       |       |   |
| R <sub>15</sub> |       |       |       |   |
| .               |       |       |       |   |
| .               |       |       |       |   |

{A different set of partial list of variable considered}

step-3: \* classifying.

→ Based on 'n' samples, 'n' trees are built.

→ Each record is classified based on 'n' tree. Final class record is decided based on voting.

kle don't have pruning in RandomForest

Out of Bag error rate:

out of bag error rate (oob) is computed

from averaging the loss on each decision tree.

{Number of mis-classification at each node of tree}

Based on oob, the number of trees to build can be decided.

**Variable Importance:** Decisions on the output of random forest is purely driven by the predictive power of the variable. The variable with high predictive power contributes a lot on final decision. In some cases variables with low predictive power are omitted from the model, making it simpler and faster to fit & predict.

There are two measures of importance given for each variable in random forest,

\* Accuracy based importance : This method finds how much accuracy of the model decreases when a variable is excluded.

\* Gini based importance : This method finds out the average of highly used variable based on gini gain. The variable with high gini gain (or) low gini impurity has high importance.

Code to build Random Forest :

```
from sklearn.ensemble import RandomForestClassifier  
RandomForestClassifier(bootstrap = True, criterion = 'gini'  
n_estimators = 100,  
max_features = 'auto',  
max_leaf_node = 10,  
min_sample_split = True)
```

Which one is best? RF or decision Tree?

No algo is superior or inferior. The usage depends on domain.

RF is not used due to no clarity in structure.

RF is preferred when there is more predictor variables.

## Linear Regression

: algorithm

After b/w and glands with disease present with it  
 ⇒ Linear regression is a machine learning algorithm that finds (the best) linear fit relationship on any given data between independent and dependent variable.

⇒ The target or dependent variable in case of linear regression will be a continuous variable.

⇒ Linear regression is also called as multiple regression, multi variate regression, ordinary least square (OLS) regression.

⇒ In linear regression, the regression line predicts the pattern of data.

$y = \beta_0 + \beta_1 x$  → Increase / decrease  
 PREDICTS THE CHANGE IN "Y" WHEN  
 "X" INCREASES BY ONE UNIT

Regression line equation:  $y = mx + c$

$$y' = mx' + c$$

→ {any value of x}  
 ↓  
 {predicted} {slope} → {y intercept}  
 {value of y}

From the above equation:

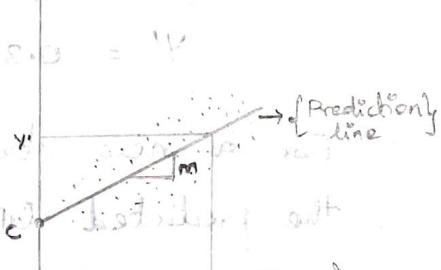
$$c = -m\bar{x} + \bar{y} \quad (\text{or})$$

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

and

$$m = r \left( \frac{s_y}{s_x} \right)$$

r - correlation coeff.  
 S<sub>x</sub>, S<sub>y</sub> - slope standard deviation of x & y.



Example:

The table represents the study time and GPA of students. Predict the student GPA from the amount of time they study each week.

| Hours (x)     | GPA (y)       |
|---------------|---------------|
| 1             | 2             |
| 2             | 1.5           |
| 3             | 2.5           |
| 5             | 3.5           |
| 6             | 3.0           |
| 8             | 4.0           |
| 10            | 4.5           |
| $\bar{x} = 5$ | $\bar{y} = 3$ |
| $s_x = 3.26$  | $s_y = 1.08$  |

We know that,  $y' = mx' + c \quad \text{--- } ①$

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

$$m = r \left( \frac{s_y}{s_x} \right) = 0.94 \left( \frac{1.08}{3.26} \right) = 0.311$$

$$c = 3 - (0.311 \times 5)$$

$$c = 1.45 \quad \text{--- } ③$$

Substitute ② and ③ in ①, we get

$$y' = 0.311(x') + 1.45$$

For a new student with study time  $x'$  hours the predicted GPA will be  $y' = 0.311(x') + 1.45$

As study time increases by one hour, we predict the GPA increases by 0.311

The example is just a numerical representation with less complex line fitting. When we have more number of datapoints with similar x value but different y value, we choose to find RSS for best line of fit.

RSS - Residual sum of squares.

$$RSS = \sum_{i=1}^n [(Actual - Predicted)^2]$$

RSS is squared to prevent the cancellation of lower order values by higher ones. [The line with least RSS value is considered as the best fit line.]

Accuracy of linear regression :

The accuracy of a linear regression model can be evaluated by a term called  $R^2$ .

①  $R^2$  - It represents the percentage variation in the data that the model can understand / capture. Higher the value of  $R^2$  better the prediction will be.

{Adjusted}  $R^2$  - Adjusted  $R^2$  gives the value of per unit change in  $R^2$  value when you are removing or adding a particular variable in the model.

② One big disadvantage of  $R^2$  is that it will keep on increases when a variable is added to a model, even if the variable is not significant. In the other hand the adjusted  $R^2$  decreases if the variable is not significant. In this way Adjusted  $R^2$  is superior than the  $R^2$ .

Code for checking statistical significance:

```
import statsmodel.api as sm
```

mod = sm.OLS(y, x)

res = mod.fit()

print(res.summary())

R<sup>2</sup>, Adjusted R<sup>2</sup>, P-values can be found using this

Code for Linear Regression model:

```
from sklearn.linear_model import LinearRegression
```

```
linreg = LinearRegression()
```

```
model = linreg.fit(x-train, y-train)
```

Code for finding model accuracy (R<sup>2</sup>):

```
print(model.score(x-train, y-train))
```

```
print(model.score(x-test, y-test))
```

# clustering

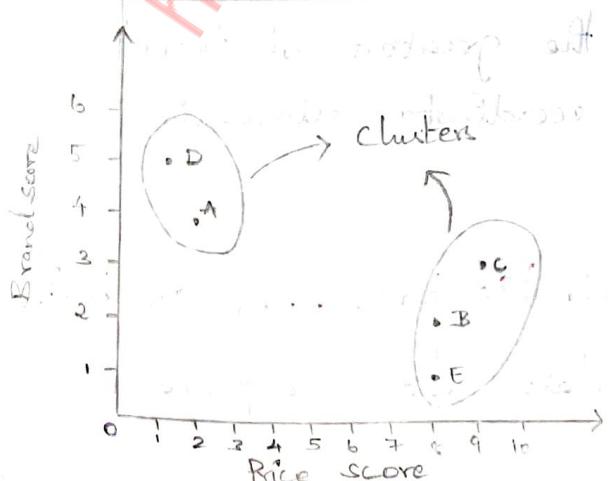
\* clustering is a type of unsupervised learning algorithm. It is a technique of finding similar group of datapoints in the data, called as clusters.

\* It groups data instances that are similar to each other in one cluster and data instances that are very different from each other into different clusters.

\* Cluster is a collection of objects that are homogeneous between them and heterogeneous to objects of other clusters.

Example:

| Individual | Price Score | Brand Score |
|------------|-------------|-------------|
| A          | 2           | 4           |
| B          | 8           | 2           |
| C          | 9           | 3           |
| D          | 1           | 5           |
| E          | 8           | 1           |



Insights:  
From scatter plot, we see that,

- \* customers D & A belong to the cluster of very high brand score
- \* customer C & B & E belongs to the cluster of very high price score

## Distance computation:

practical

In any of the clustering algorithm, distance computation plays a major role in decision making. Most popular distance metrics are

- \* Euclidean distance.
- \* Chebyshov distance.
- \* Manhattan distance.

### Euclidean distance:

Euclidean distance is measured by

drawing a shortest line between two points and plugging the values into the formula

$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + \dots + (z_n - z_1)^2}$$

point

point

point

(x<sub>1</sub>, y<sub>1</sub>)

(x<sub>2</sub>, y<sub>2</sub>)

: adjacent

Euclidean distance is the extension of Pythagoras theorem.

### Chebyshov distance:

In chebyshov distance, the distance between two vectors is the greatest of their differences among all coordinate dimensions.

The formula is,

$$d = \max(|a_1 - b_1|, |a_2 - b_2|, \dots, |a_n - b_n|)$$

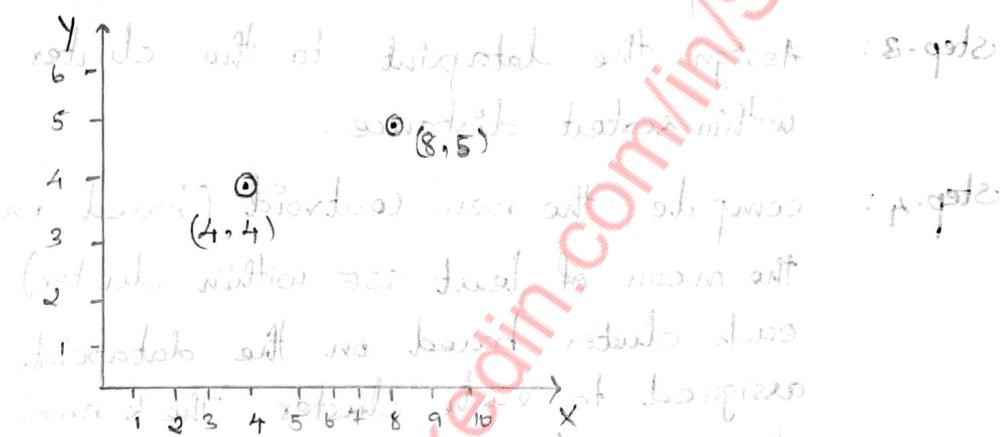
This is used for the data where response is ordinal in nature.

Manhattan distance: ; grid-based search

Manhattan distance is called as the city block distance. It is the sum of perpendicular distances. Formula is,  $d = |a_1 - b_1| + |a_2 - b_2| + |a_3 - b_3| + \dots + |a_n - b_n|$  (add all above) : straight line distance

Example:

Formulate Euclidean, chebychev and manhattan distance for the given two points.



$$\text{Euclidean} = \sqrt{(8-4)^2 + (5-4)^2}$$

$$= \sqrt{16+1} \quad \text{f=40 steps to goal} \quad (8\text{-gate})$$

$$= 4.123 \quad \text{all moves are f=40 steps}$$

$$\text{chebychev} = \max(|8-4|, |5-4|)$$

$$= \max(4, 1)$$

$$= 4$$

$$\text{manhattan} = |8-4| + |5-4|$$

forward diagonal (8) backward diagonal (4)

$$= 4 + 1 \quad \text{f=40 steps to goal}$$

$$= 5$$

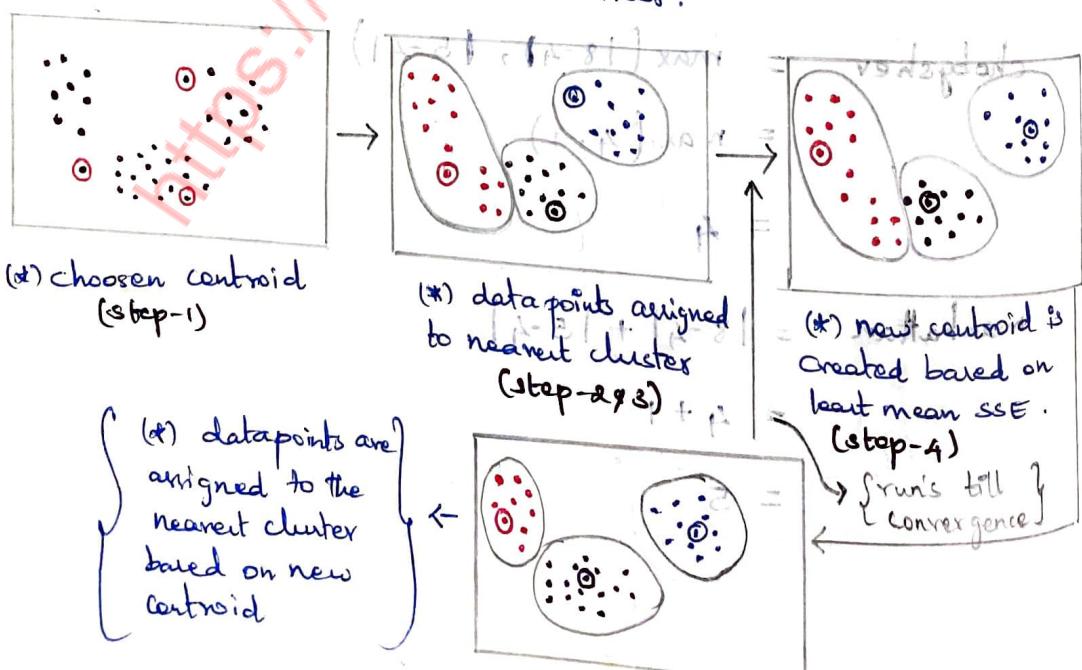
forward diagonal (5) backward diagonal (4)  
left of target  
right of target  
upward to goal  
downward to goal

## K-means clustering:

photo of k-means clustering introduction unsupervised, non-hierarchical clustering algorithm. It is based on within cluster variation. (squared distance from the centre of cluster).

steps in k-means algorithm: (Inside blackbox)

- step-1: Assume k-centroids (for k-clusters)
- step-2: compute Euclidean distance of each datapoint with these centroids.
- step-3: Assign the datapoint to the cluster within shortest distance.
- step-4: compute the new centroid (based on the mean of least SSE within cluster) of each cluster based on the datapoints assigned to each cluster. The k number of means obtained will become the new centroid for each cluster.
- step-5: Repeat steps 2-4 till convergence. That is defined by no movement of datapoints from one cluster to another.



centroids are not the existing datapoints, they are not created as new data point  $\times$

for example, kept centroids not used

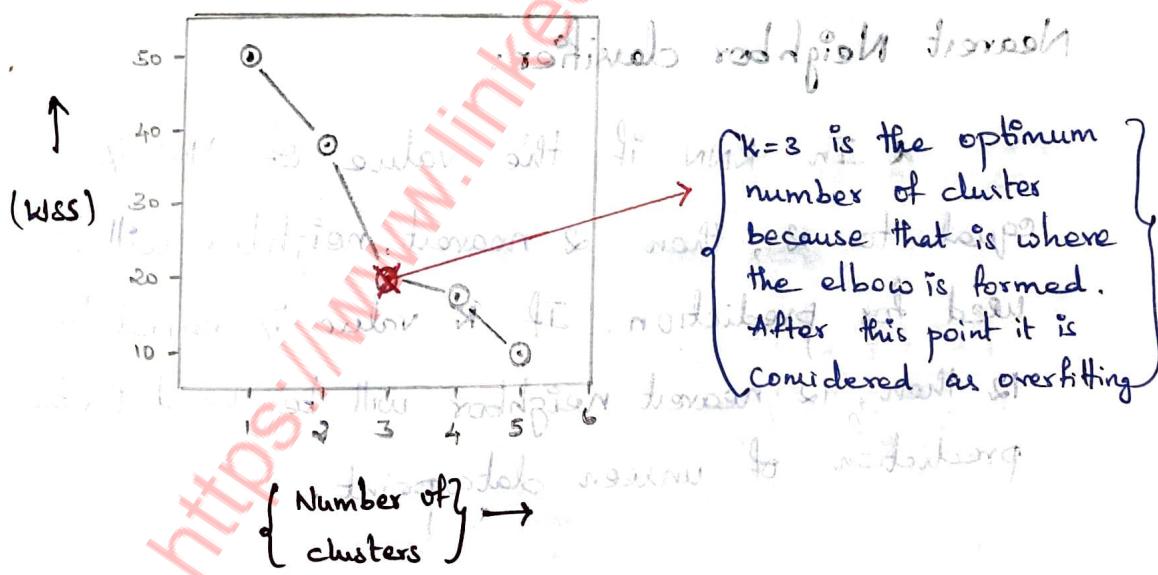
choosing optimum number of clusters:

The one of the biggest challenge of k-means clustering is pre-defining the number of clusters (step -1). To pre-define the optimum number of clusters we have to find it using elbow graph, that is popularly called as WSS plot.

WSS plot:

Within Sum of Square error (WSS) plot is used to find the optimum number of cluster.

This plots the number of clusters and the within group sum of square error of each cluster.



Standard Scalar / preprocessing .scale()

It converts normal distribution to standard normal distribution where  $\mu=0, \sigma=1$  always. This makes all the feature in same range.

Required In - Regression, k-means, KNN, SVM

Not required In - Decision tree, Random forest (As it uses GINI)

## K-Nearest Neighbor classification (KNN)

KNN is a classification type, supervised learning algorithm. It is an instance based classifier, that stores the training record and uses training records to predict the class label of unseen data / classes.

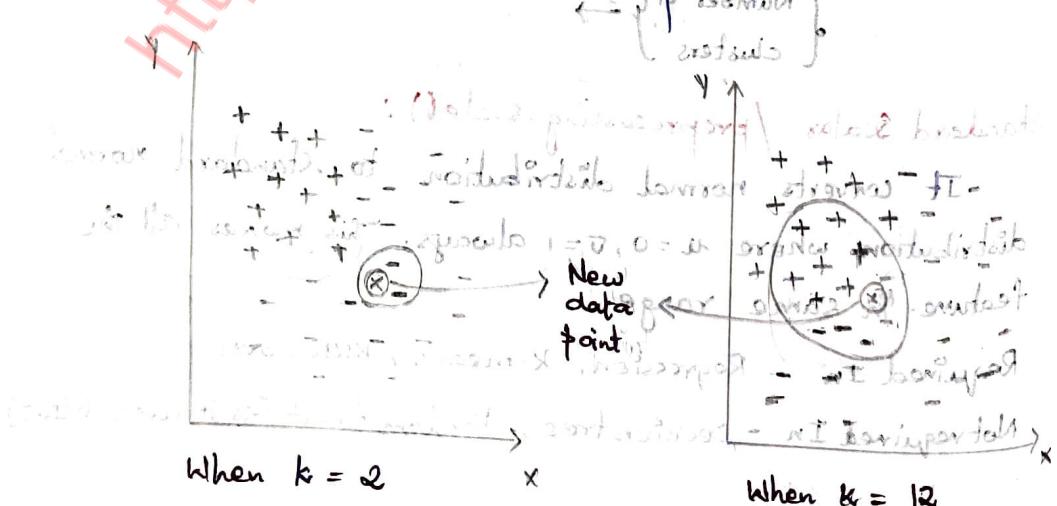
Example of instance based classifier:

\* **Rule learner**: Memorizes entire training data and perform classification only if attributes of record matches other training example.

\* **Nearest Neighbor**: Uses  $k$ ' closest point, that is called as nearest neighbor for performing classification.

Nearest Neighbor classifier:

\* In KNN if the value of ' $k$ ' is equal to 2, then 2 nearest neighbor will be used for prediction. If  $k$  value is equal to 12 then, 12 nearest neighbor will be used for the prediction of unseen data point.



\* The distance between the nearest neighbor and the unseen datapoint is calculated using Euclidean distance.

$$d(p, q) = \sqrt{\sum_i (p_i - q_i)^2}$$

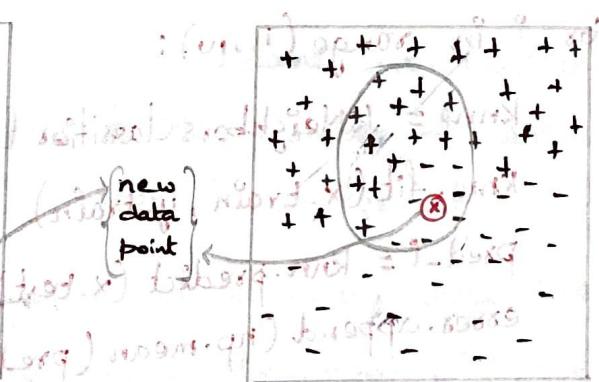
$$\text{similarity} = 1/d(p, q)$$

$\propto 1/d(p, q) \Rightarrow \begin{cases} \text{When the distance between the points increases, the similarity between them decreases & viceversa.} \\ \text{as } d(p, q) \text{ increases} \end{cases}$

choosing the value of 'k'

- \* The value of 'k' should be chosen as odd number to make the correct decision. If it is even number, there is 50-50 chance of both classes and the prediction can fall in any of those classes.
- \* If the value of 'k' is too small, sensitive to noise points.
- \* If the value of 'k' is too large, neighborhood may include points from other class.

|   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|
| + | + | + | - | - | - | - |
| + | + | + | + | + | - | - |
| + | + | + | + | + | - | - |
| - | + | + | + | + | - | - |
| - | + | + | + | + | - | - |
| - | + | + | + | + | - | - |
| + | - | - | - | - | - | - |
| - | - | - | - | - | - | - |
| + | - | - | - | - | - | - |
| + | + | - | - | - | - | - |



(+) When 'k' value is small. Actually the new data point has to be (-) but due to low 'k' value it is (+).

(\*) When 'k' value is large. The new data point has to be (-), but due to high 'k' value it is predicted as (+) by taking points from other class.

\* The optimum value of 'k' has to be chosen by running a loop with a range of 'k' value and plotting it with graph. The value with optimum and less number of error rate (misclassification) has to be chosen.

code for k-NearestClassifier:

```
from sklearn.neighbors import KNeighborsClassifier  
classifier = KNeighborsClassifier(n_neighbors=5)
```

classifier.fit(x\_train, y\_train) # to calculate fit

Code for prediction & Evaluation:

y\_pred = classifier.predict(x\_test)

```
from sklearn.metrics import classification_report,  
confusion_matrix
```

print(confusion\_matrix(y\_test, y\_pred))

print(classification\_report(y\_test, y\_pred))

Code for optimum k-value detection:

error = [] # to store error rate

for i in range(1, 40):

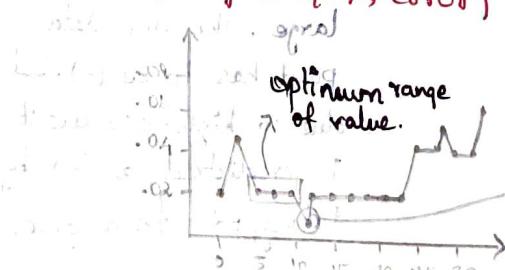
knn = KNeighborsClassifier(n\_neighbors=i)

knn.fit(x\_train, y\_train)

pred\_i = knn.predict(x\_test)

error.append(np.mean(pred\_i != y\_test))

plt.plot(range(1, 40), error)



→ outer 'k' and W(k)

→ plot A. Name

at end bring back k

at end but (1-k)

→ overfitting above it will

\* KNN is a lazy learner, it does not build model explicitly unlike eager learners such as decision tree & random forest. So, this is computationally expensive algorithm.

### Difference between clustering and classification?

- (\*) Clustering is unsupervised ML algorithm, the clusters are homogeneous within themselves and heterogeneous to each other.
- (\*) Classification is supervised ML algorithm, the classifiers can be heterogeneous or homogeneous to each other.

### Difference between KNN and K-means?

- (\*) KNN is a classification type supervised learning algorithm. The 'k' value has to be predicted by using misclassification rate (error) by looping over a range of values and the value of 'k' represents the number of neighbors to select near the point.
- (\*) K-means is a clustering type unsupervised learning algorithm. The value of 'k' in K-means is found using within sum of square error (WSS) plot. The value of 'k' represents the number of centroids and the number of clusters to be formed for clustering.  
**similarity:** Both KNN and K-means use the distance metric (Euclidean distance) to calculate the nearest points.

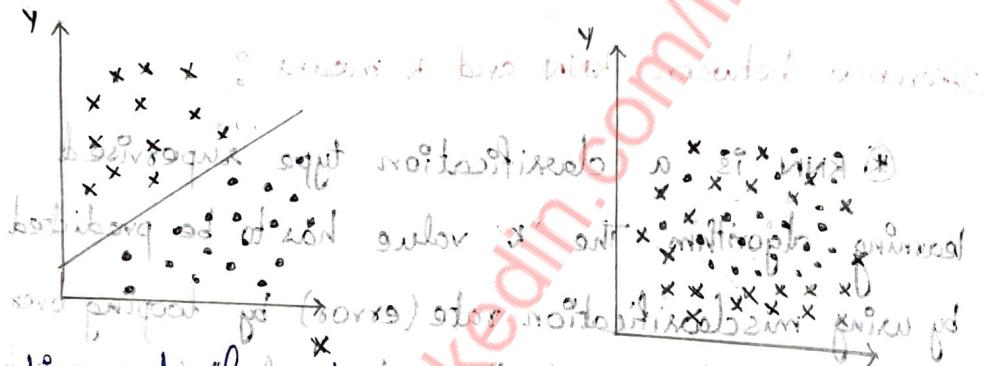
{  
annals - scikit  
allie - pyglet  
fjern - fennel}

of these result with

# Support Vector Machines (SVM)

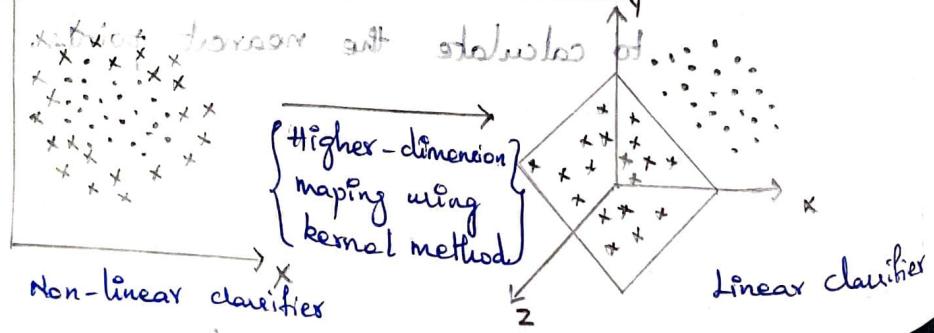
- \* SVM is a classification type, supervised learning algorithm. It is related to statistical learning theory developed in 1992. It has a greater success in handwritten digit recognition.
- \* It is an example of kernel learning method / kernel method. It is one of the key area in machine learning.

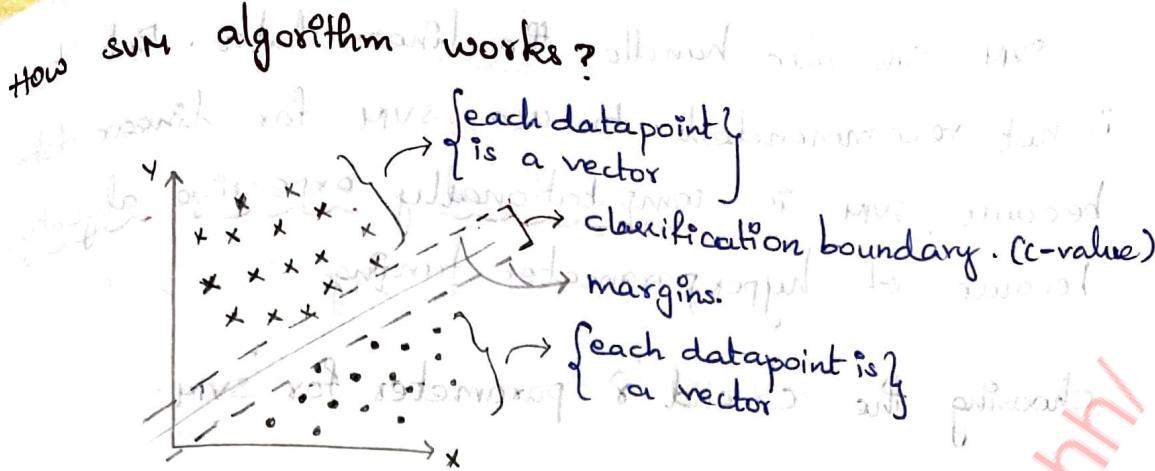
Types of classifier: Linear & non-linear classifier.



Linear classifier, where all the data points can be separated by a simple straight line.  
Non-linear classifier, where the data points can't be separated by a simple straight line.

SVM is capable of handling both linear and non-linear data by using its own kernel method. Few of the SVM kernels are polynomial kernel, Gaussian kernel, Gramian radial basis function, Laplace RBF, Hyperbolic tangent kernel, sigmoid kernel etc.,





\* SVM creates the classification boundary as a threshold that touches the extreme datapoints in both the cases. Below or above which the classes are predicted by the threshold.

\* The datapoints are called as the vectors. The margin of boundary (classification boundary) is supported by three vectors and so it's called as SVM.

### Choosing the kernel function:

i) Polynomial kernel. → For non-linear data.

ii) Radical or RBF kernel. → For non-linear data,

iii) Linear kernel. → For linear data.

### How to find the data is linear/non-linear?

\* We can't simply say by seeing the data, As there is no way of visually finding the linearity of data. As in real time it has many dimensions. Number of dimensions will be the number of features so it is not possible to find visually by seeing data.

\* When the basic classification algorithm gives low accuracy even after all EDA (65% - 75%), then the data may be non-linear type. Then try SVM and if accuracy increases then its non-linear data.

SVM can also handle the linear data. But it is not recommended to use SVM for linear data because SVM is computationally expensive algorithm because of hyperparameter tuning.

choosing the C and  $\gamma$  parameter for SVM:

C parameter : It is a regularisation parameter that controls the tradeoff between achieving a low training error and low testing error. It is a ability to generalise the classifier to unseen data.

C represents the width of the margin.

for classification purpose

Higher C = less width  $\Rightarrow$  not a good model

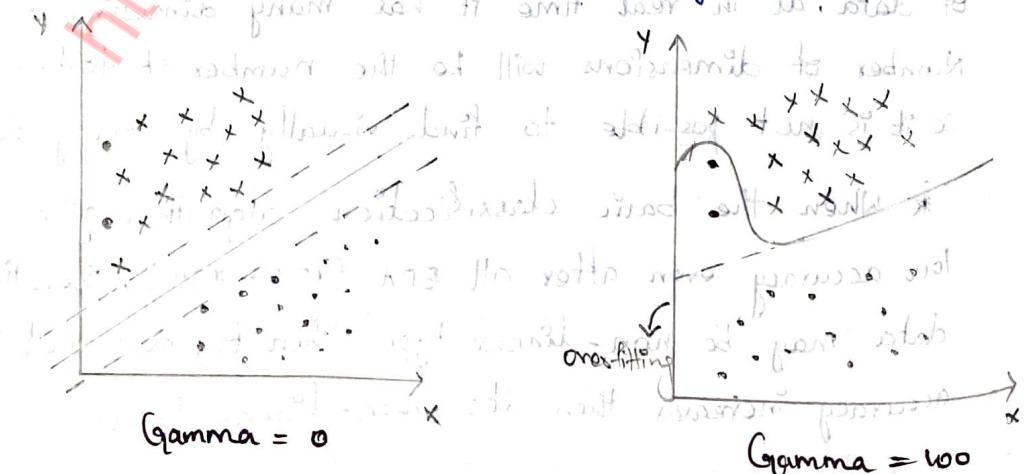
Lower C = Higher width  $\Rightarrow$  good model

Optimum value of C has to be calculated based on model accuracy.

Gamma parameter : It defines how far the influence of single training example reaches.

Low value means far and high value means close.

Optimum value of  $\gamma$  has to be chosen otherwise model will run into overfitting.



How to choose the optimum value for c, gamma and the kernel method?

The optimum values of c, gamma and the right kernel can be chosen with the help of GridSearchCV. The GridSearchCV was designed originally for SVM because the hyperparameter tuning is difficult in SVM.

\* SVM should not be used unless until the accuracy can't be increased and the data can't be linearly separable. This is defined by the law of parsimony, it states that "work of sword needle should not be done by sword". We have to try other non expensive algorithm for classification and prefer SVM and KNN at worst case.

\* Random forest is a space consuming algorithm when we have a constrain of space we can use decision tree. Because random forest creates data with using bagging, so it creates more trees.

When we know the data is non-linear which algorithm can be used?

SVM can be used because it can use kernel method to transform the non-linear data into linear data. In this way, SVM is superior in handling non-linear data.

# Logistic Regression with Example

Short Note by [REDACTED]

\* Logistic Regression is a supervised learning, classification algorithm. It is used to predict for binary outcome based on a set of independent variable. It is probabilistic model.

\* like other classification algorithm, logistic regression can be used with independent variables that falls in following categories, and are called as,

- i) Binary logistic Regression - spam / not-spam
- ii) Multinomial logistic Regression - veg / Non-veg / vegan
- iii) Ordinal logistic Regression - Movie rating (1 to 5)

Sigmoid Function:

Sigmoid function can convert any number between  $-\infty$  to  $+\infty$  into the range of 0 to 1. The formula for sigmoid function is,

$$\textcircled{1} \quad -S(x) = \frac{1}{1+e^{-x}} = \frac{1}{1+e^{-y}}$$

$$\textcircled{2} \quad -y = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + \dots + w_n x_n$$

Substitute  $\textcircled{2}$  in  $\textcircled{1}$ , we get,

$$S(x) = \frac{1}{1+e^{-(w_0 + w_1 x_1 + \dots + w_n x_n)}}$$

→ Sigmoid function gives a a sigmoid curve that is also called as 'S-curve'.

Non-parallel sigmoid not blurriness

is transformed with activation of sigmoid

from  $-\infty$  to  $\infty$  for output 0 to mapping into  $+\infty$

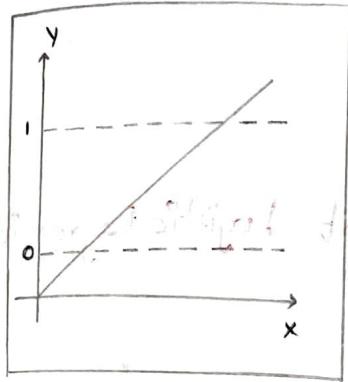
mapping the values between  $-\infty$  to  $+\infty$  into range of 0 to 1.

Horizontal axis of linear function present

obtain threshold with  $y = 0.5$

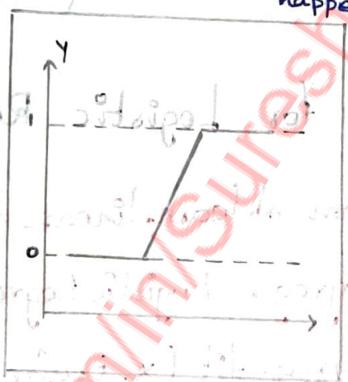
$[(\text{val} < 0.5) = 0, (\text{val} > 0.5) = 1]$

This is how the classification happen in logreg



Transforming a line into s-curve

$$S(n) = \frac{e^n}{1+e^{-n}}$$



(not x) different input, a mapping

softmax function:

$$P = \frac{e^{w_1 x_1 + w_2 x_2 + \dots + w_n x_n}}{1 + e^{w_1 x_1 + w_2 x_2 + \dots + w_n x_n}}$$

ODDS RATIO:

The weight factors ( $w_1 x_1 + w_2 x_2 + \dots + w_n x_n$ ), the value of  $w_n$  will be calculated by odds ratio in logistic regression. The odds ratio is defined as the ratio of odds of A in the presence of B, A in the absence of B or equivalently odd of B in the presence of A, B in the absence of A.

It quantifies the strength of association between two events A and B.

$$\text{odd} = \frac{\text{Probability}}{1 - \text{probability}}$$

## Threshold for Logistic Regression:

Threshold for achieving the highest accuracy

in the logistic regression depends on the output

of the training data and it is internally

assigned by the sklearn module.

$$y = \begin{cases} 0, & w_i x_i < \text{threshold} \\ 1, & w_i x_i \geq \text{threshold} \end{cases}$$

## Code for Logistic Regression:

```
from sklearn.linear_model import LogisticRegression  
logreg = LogisticRegression()  
logreg.fit(x-train, y-train)  
y-pred = logreg.predict(x-test)
```

## Code for accuracy and performance:

```
cmt-matrix = metrics.confusion_matrix(y-test, y-pred)
```

cmt-matrix

```
print('Accuracy', metrics.accuracy_score(y-test, y-pred))
```

```
print('precision', metrics.precision_score(y-test, y-pred))
```

```
print('Recall', metrics.recall_score(y-test, y-pred))
```

As the accuracy is 1.0 due to other

the precision and recall is 1.0 due to same reason

As the accuracy is 1.0 due to same reason

As the accuracy is 1.0 due to same reason

As the accuracy is 1.0 due to same reason

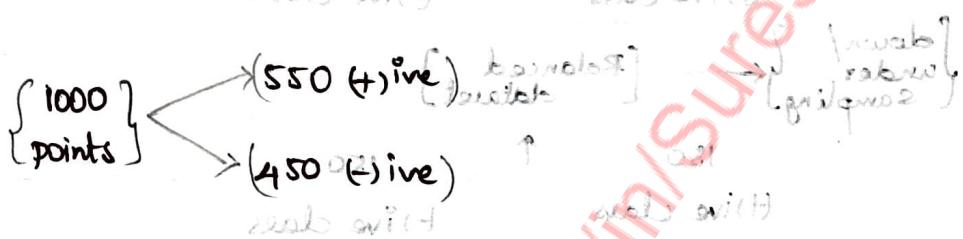
$$\text{precision} = 1.0$$

How to handle a Imbalance dataset? X

Balanced dataset : dataset having equal numbers of samples from each class

Balanced dataset is one which contains equal or approximately equal number of samples from each available class.

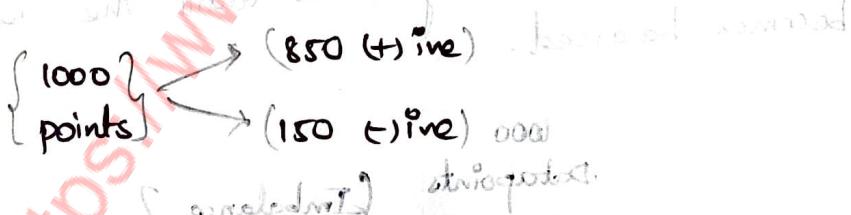
Ex : out of 1000 datapoints available, if 550 belongs to +ive class and 450 belongs to -ive class then its balanced dataset.



Imbalance dataset :

Imbalance dataset is one which does not contain equal or approximately equal number of samples for each available class.

Ex : out of 1000 datapoints available, if 850 belong to (+)ive and 150 belong to (-)ive class then its imbalanced dataset.



Ways to handle imbalance dataset :

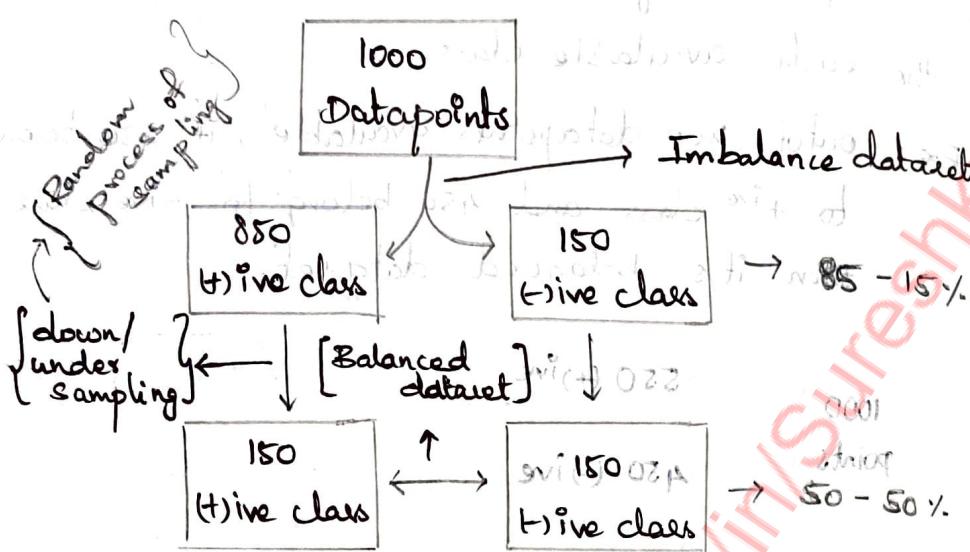
Imbalance dataset may lead to low recall scores and larger misclassification in minority class. This can be prevented by converting imbalance dataset to balanced one. This is done by,

\* Down sampling (undersampling)

\* Up sampling (oversampling)

## \* Down sampling / under sampling

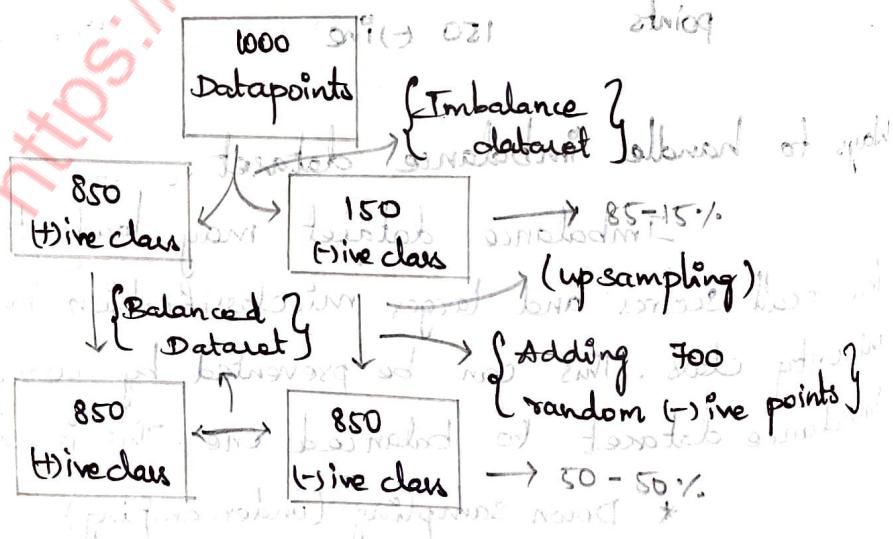
Under sampling simply means that samples from majority class until the dataset becomes balanced.



→ This leads to artificially reducing the size of dataset, by which some valuable information might be missed.

## \* Up sampling / over sampling (simple)

Over sampling simply means duplicating samples from minority class until the dataset becomes balanced.

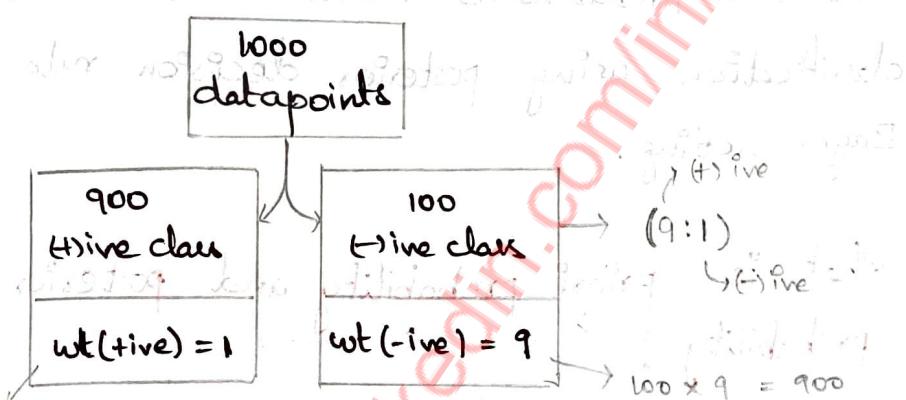


\* upsampling / oversampling by extrapolation:

This is called as synthetic minority over-sampling technique (SMOTE). In this method of oversampling artificial datapoints/synthetic datapoints for minority class is generated till the dataset is balanced.

\* upsampling / oversampling by class-weight:

In this method (the) weight of the minority class will be equal to the proportion of the majority class.



$900 \times 1 = 900$   
→ Now each data point of both classes has same weights.  
→ Synthetic was created with cost reduction.

→ This is known as oversampling.  
→ oversampling may lead to overfitting of the machine learning models because it creates the exact copies / duplicates of the existing minority class data points.

→ increased variance

→ much slower speed with more memory

## Naive Bayes classification

It is a supervised learning classification.

type probabilistic machine learning algorithm.

It is based of Bayes theorem for calculating probabilities and conditional probabilities.

Bayes Rule:

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

It is a probabilistic model that makes classification using posterior decision rule in Bayes setting.

What is prior probability and posterior probability?

Prior probability what is originally believed before new evidence is introduced. Posterior probability takes this information into account. Simply, whatever is known is called as prior probability (coin, dice, card, etc..). Whatever has to be calculated by data is called as posterior probability.

What are all the probabilistic machine learning models?

- \* Logistic Regression.
- \* Naive Bayes classification.

# Bayes classification work?

How

| Age | Gender | Occupation | skill | Target |
|-----|--------|------------|-------|--------|
| 26  | M      | self       | ML    | 0      |
| 25  | F      | Govt       | C++   | 1      |
| 24  | M      | Govt       | ML    | 0      |
| 22  | M      | Priv       | ML    | 1      |
| 26  | F      | sal        | C++   | 1      |
| 24  | M      | self       | C++   | 0      |

$$FP_0 = (0.000)(-1)$$

$$FP_1 = (0.000)(1)$$

Now a new datapoint arrives for prediction:

$$FP_0 = (0.000)(-1)$$

$$FP_1 = (0.000)(1)$$

→ Age - 26, Gender - M, occupation - self, skill - ML

Now the prediction happens using Bayes rule, with available training data.

$$P(Y=0) \Rightarrow P(Y=0 | \text{Age}) \times P(Y=0 | \text{Gender}) \times P(Y=0 | \text{skill}) \\ \times P(Y=0 | \text{occupation})$$

$$P(Y=1) \Rightarrow P(Y=1 | \text{Age}) \times P(Y=1 | \text{Gender}) \times P(Y=1 | \text{skill}) \times \\ P(Y=1 | \text{occupation})$$

decision : If  $P(Y=0)$  is (say - 0.88) greater than

the  $P(Y=1)$  (say - 0.12), then the new dataset will fall into '0' class.

If  $P(Y=1)$  is (say - 0.7) greater than  $P(Y=0)$  (say - 0.3), then the new dataset will fall into '1' class.

Find the probability of one given another?

A patient takes the lab test and results come as positive. The test returns correct positive results in 98% of cases in which the disease is actually present and correct negative results in 97% where disease is not present. Furthermore, 0.008 of entire population has cancer.

probabilities:

$$P(\text{Cancer}) = 0.008 \quad P(-\text{cancer}) = 0.992$$

$$P(+|\text{cancer}) = 0.98 \quad P(-|\text{cancer}) = 0.02$$

$$P(+|-\text{cancer}) = 0.03 \quad P(-|-\text{cancer}) = 0.97$$

Why Bayes classifier is called as 'Naive' Bayes classifier?

It is called as Naive because the

character of one feature does not affect the character of other features. Each feature has a direct probabilistic relation with the target variable. It is not affected by multicollinearity, it does not require scaling / Normalization.

In simple words, it is so naive that it does not care about other variables.

# Principal Component Analysis (PCA)

Principal component analysis (PCA) is a feature engineering technique used for dimension reduction. It is invented by Pearson and Hotelling.

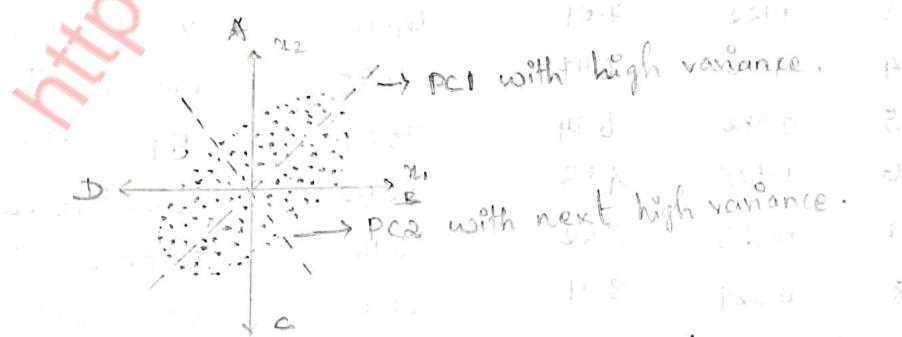
PCA takes a data matrix of  $n$  objects by  $p$  variables, which may be correlated and summarizes it by uncorrelated axes (principal component / principal axes) that are linear combination of original  $p$  variables.

Example:

|     | $x_1$  | $x_2$  | $x_3$  | $x_4$ | $x_5$ | $x_6$ | $x_7$ | $x_8$ |        |
|-----|--------|--------|--------|-------|-------|-------|-------|-------|--------|
| Age | Height | Height | Weight | skill | skill | Qual  | Exp   | Sal   | Target |
| PC1 | PC2    |        |        |       |       |       |       |       | 1<br>8 |

Now PC1, PC2 can be further used for upcoming process such as training & testing.

PC's has to be created in a way that PC1 should have high variance, PC2 should have next high variance and so on.



All PC's should be orthogonal to each other. Because of which there is no possibility of covariance between the PC's.

# The Algebra of PCA

\* PCA uses Euclidean distance calculated from the p variables as the measure of dissimilarity among the n objects. So it is mandatory to standardize / normalize the data before PCA.

\* In practice nobody uses PCA with two variables. Only when the dataset has more number of features PCA can be used (say, 300 variables).

\* Eigen value and Eigen vectors are used behind the PCA for creating PC's.

Eigen Value - It gives the percentage (%) variance covered in particular axis.

Eigen vector - It gives the direction of the particular PC.

PCA is a blackbox technique

| Eigen Value |             |               |                          |
|-------------|-------------|---------------|--------------------------|
| Axis / PC's | Eigen value | % of Variance | Cumulative % of Variance |
| 1           | 5.855       | 36.60         | 36.60                    |
| 2           | 3.420       | 21.38         | 57.97                    |
| 3           | 1.122       | 7.01          | 64.98                    |
| 4           | 1.116       | 6.97          | 71.95                    |
| 5           | 0.982       | 6.14          | 78.09                    |
| 6           | 0.725       | 4.53          | 82.62                    |
| 7           | 0.563       | 3.52          | 86.14                    |
| 8           | 0.529       | 3.31          | 89.45                    |
| 9           | 0.476       | 2.98          | 92.42                    |
| 10          | 0.375       | 2.35          | 94.77                    |
| 11          | 0.312       | 1.98          | 96.75                    |

Optimum number of PCs can be chosen based on cumulative % of variance. We can see that upto 5 PC's the cumulative sum increases drastically between the PC and after which it has decreased. So 5 PC's is considered as optimum.