# 1. Write a short note on Supervised machine learning?

**Ans :** **Supervised Machine Learning** is an algorithm that learns from labeled training data to help you predict outcomes for unforeseen data. In Supervised learning, you train the machine using data that is well "labeled." It means some data is already tagged with correct answers. It can be compared to learning in the presence of a supervisor or a teacher.

Successfully building, scaling, and deploying accurate supervised machine learning models takes time and technical expertise from a team of highly skilled data scientists. Moreover, Data scientist must rebuild models to make sure the insights given remains true until its data changes.

Supervised machine learning uses training data sets to achieve desired results. These data sets contain inputs and the correct output that helps the model to learn faster.
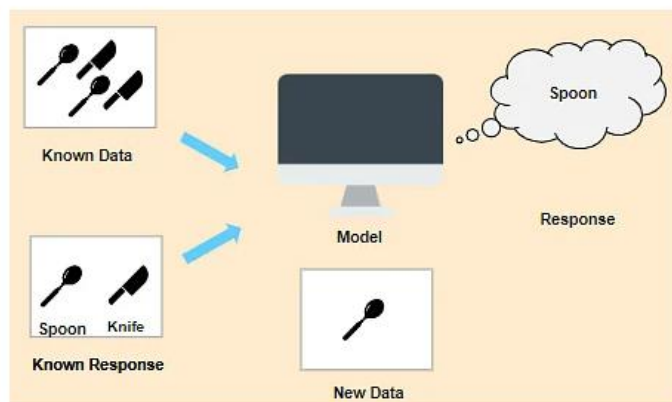
 For example, you want to train a machine to help you predict how long it will take you to drive home from your workplace.

Here, you start by creating a set of labeled data. This data includes:

- Weather conditions
- Time of the day
- Holidays

All these details are your inputs in this Supervised learning example. The output is the amount of time it took to drive back home on that specific day.

In Supervised Learning, the machine learns under supervision. It contains a model that is able to predict with the help of a labeled dataset. A labeled dataset is one where you already know the target answer.

# Module – IV: Supervised Learning –Classification

In this case, we have images that are labeled a spoon or a knife. This known data is fed to the machine, which analyzes and learns the association of these images based on its features such as shape, size, sharpness, etc. Now when a new image is fed to the machine without any label, the machine is able to predict accurately that it is a spoon with the help of the past data.

**Types of Supervised Machine Learning Algorithms**

Following are the types of Supervised Machine Learning algorithms:

**Regression:**

Regression technique predicts a single output value using training data.

**Classification:**

Classification means to group the output inside a class. If the algorithm tries to label input into two distinct classes, it is called binary classification. Selecting between more than two classes is referred to as multiclass classification.

**2. Write about logistic regression?**

Logistic regression is one of the most popular Machine Learning algorithms, which comes under the Supervised Learning technique. It is used for predicting the categorical dependent variable using a given set of independent variables.

Logistic regression predicts the output of a categorical dependent variable. Therefore the outcome must be a categorical or discrete value. It can be either Yes or No, 0 or 1, true or False, etc. but instead of giving the exact value as 0 and 1, **it gives the probabilistic values which lie between 0 and 1**.

Logistic Regression is much similar to the Linear Regression except that how they are used. Linear Regression is used for solving Regression problems, whereas **Logistic regression is used for solving the classification problems**.

In Logistic regression, instead of fitting a regression line, we fit an "S" shaped logistic function, which predicts two maximum values (0 or 1).
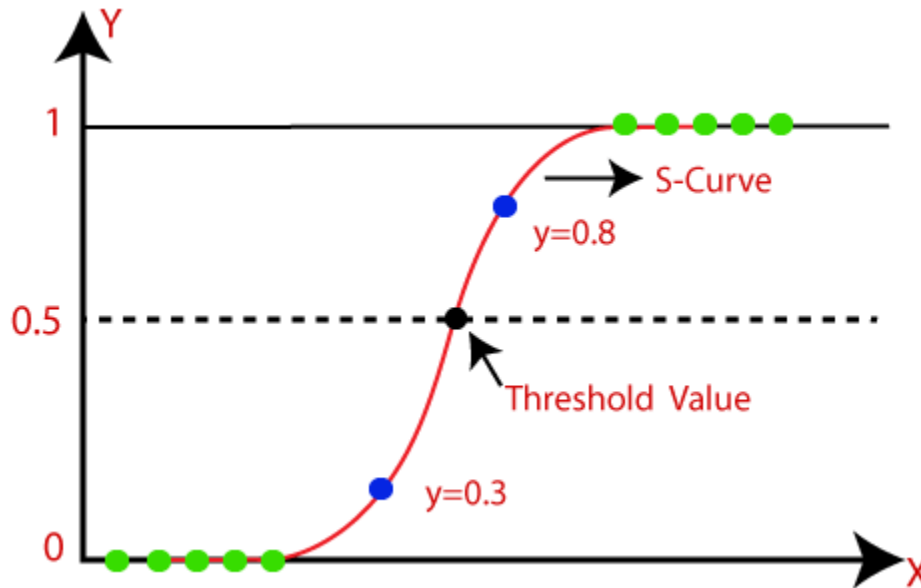
The curve from the logistic function indicates the likelihood of something such as whether the cells are cancerous or not, a mouse is obese or not based on its weight, etc.

# Module – IV: Supervised Learning –Classification

Logistic Regression is a significant machine learning algorithm because it has the ability to provide probabilities and classify new data using continuous and discrete datasets.

Logistic Regression can be used to classify the observations using different types of data and can easily determine the most effective variables used for the classification.

The below image is showing the logistic function:



**Type of Logistic Regression:**

On the basis of the categories, Logistic Regression can be classified into three types:

- **Binomial:** In binomial Logistic regression, there can be only two possible types of the dependent variables, such as 0 or 1, Pass or Fail, etc.
- **Multinomial:** In multinomial Logistic regression, there can be 3 or more possible unordered types of the dependent variable, such as "cat", "dogs", or "sheep"
- **Ordinal:** In ordinal Logistic regression, there can be 3 or more possible ordered types of dependent variables, such as "low", "Medium", or "High".

# Module – IV: Supervised Learning –Classification

### 3. Demonstrate the implementation of logistic regression algorithm using python?

There is a dataset given which contains the information of various users obtained from the social networking sites. There is a car making company that has recently launched a new SUV car. So the company wanted to check how many users from the dataset, wants to purchase the car.

For this problem, we will build a Machine Learning model using the Logistic regression algorithm. The dataset is shown in the below image. In this problem, we will predict the **purchased variable (Dependent Variable)** by using **age and salary (Independent variables)**.

**Steps in Logistic Regression:** To implement the Logistic Regression using Python, we will use the same steps as we have done in previous topics of Regression. Below are the steps:

- o  Data Pre-processing step
- o  Fitting Logistic Regression to the Training set
- o  Predicting the test result
- o  Test accuracy of the result(Creation of Confusion matrix)
- o  Visualizing the test set result.

**1. Data Pre-processing step:** In this step, we will pre-process/prepare the data so that we can use it in our code efficiently. It will be the same as we have done in Data pre-processing topic. The code for this is given below:

```
#Data Pre-procesing Step
# importing libraries
import numpy as nm
import matplotlib.pyplot as mtp
import pandas as pd
 #importing datasets
data_set= pd.read_csv('user_data.csv')
```

Now, we will extract the dependent and independent variables from the given dataset. Below is the code for it:

```
#Extracting Independent and dependent Variable

x= data_set.iloc[:, [2,3]].values
y= data_set.iloc[:, 4].values
```

# Module – IV: Supervised Learning –Classification

In the above code, we have taken [2, 3] for x because our independent variables are age and salary, which are at index 2, 3. And we have taken 4 for y variable because our dependent variable is at index 4. The output will be:

Now we will split the dataset into a training set and test set. Below is the code for it:

# Splitting the dataset into training and test set.
from sklearn.model_selection **import** train_test_split
x_train, x_test, y_train, y_test= train_test_split(x, y, test_size= 0.25, random_state=0)

In logistic regression, we will do feature scaling because we want accurate result of predictions. Here we will only scale the independent variable because dependent variable have only 0 and 1 values. Below is the code for it:

#feature Scaling
from sklearn.preprocessing **import** StandardScaler
st_x= StandardScaler()
x_train= st_x.fit_transform(x_train)
x_test= st_x.transform(x_test)

## 2. Fitting Logistic Regression to the Training set:

We have well prepared our dataset, and now we will train the dataset using the training set. For providing training or fitting the model to the training set, we will import the **LogisticRegression** class of the **sklearn** library.

After importing the class, we will create a classifier object and use it to fit the model to the logistic regression. Below is the code for it:

#Fitting Logistic Regression to the training set
from sklearn.linear_model **import** LogisticRegression
classifier= LogisticRegression(random_state=0)
classifier.fit(x_train, y_train)

**Out[5]:**

LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,

intercept_scaling=1, l1_ratio=None, max_iter=100,

multi_class='warn', n_jobs=None, penalty='l2',

random_state=0, solver='warn', tol=0.0001, verbose=0,

warm_start=False)

Hence our model is well fitted to the training set.

## 3. Predicting the Test Result

Our model is well trained on the training set, so we will now predict the result by using test set data. Below is the code for it:

#Predicting the test set result

y_pred= classifier.predict(x_test)

In the above code, we have created a y_pred vector to predict the test set result.

**Output:** By executing the above code, a new vector (y_pred) will be created under the variable explorer option.

## 4. Test Accuracy of the result

Now we will create the confusion matrix here to check the accuracy of the classification. To create it, we need to import the **confusion_matrix** function of the sklearn library. After importing the function, we will call it using a new variable **cm**. The function takes two parameters, mainly **y_true**( the actual values) and **y_pred** (the targeted value return by the classifier). Below is the code for it:

```
#Creating the Confusion matrix
```

from sklearn.metrics import confusion_matrix

cm= confusion_matrix(y_test,y_pred)

**Output:**

```
array([[68,  0],
       [32,  0]], dtype=int64)
```

By executing the above code, a new confusion matrix will be created

We can find the accuracy of the predicted result by interpreting the confusion matrix. By above output, we can interpret that 68 (Correct Output) and 32(Incorrect Output).
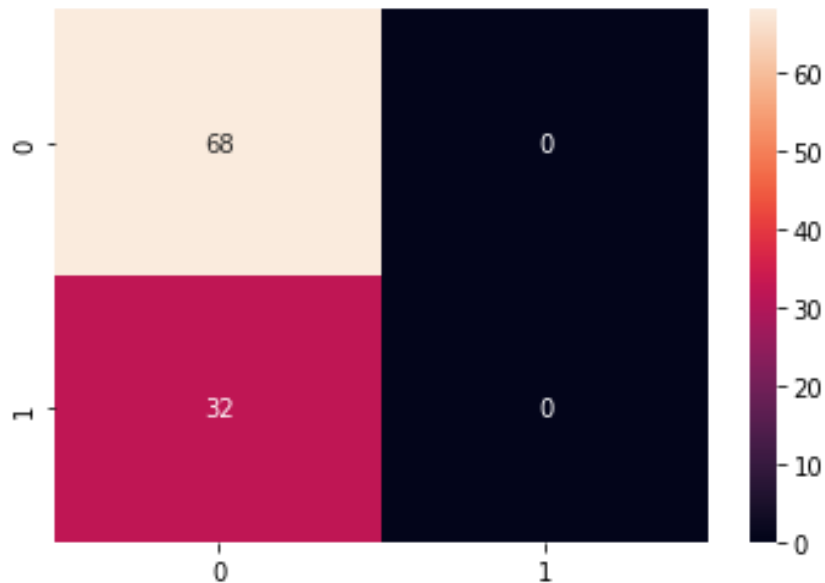
**5. Visualizing the training set result**

Finally, we will visualize the training set result. To visualize the result, we will

```
# Heatmap of Confusion matrix
sns.heatmap(pd.DataFrame(cm), annot=True)
<AxesSubplot:>

fromsklearn.metricsimportaccuracy_score

accuracy =accuracy_score(Y_Test, Y_Pred)
```

```
Out[13]:
0.89
```

## 4. Explain in-detail about K-Nearest Neighbor (KNN) algorithm?

K-Nearest Neighbour is one of the simplest Machine Learning algorithms based on Supervised Learning technique. K-NN algorithm assumes the similarity between the new case/data and available cases and put the new case into the category that is most similar to the available categories.

K-NN algorithm stores all the available data and classifies a new data point based on the similarity. This means when new data appears then it can be easily classified into a well suite category by using K- NN algorithm. K-NN algorithm can be used for Regression as well as for Classification but mostly it is used for the Classification problems.

K-NN is a **non-parametric algorithm**, which means it does not make any assumption on underlying data. It is also called a **lazy learner algorithm** because it does not learn from the training set immediately instead it stores the dataset and at the time of classification, it performs an action on the dataset.

KNN algorithm at the training phase just stores the dataset and when it gets new data, then it classifies that data into a category that is much similar to the new data.

**Example:** Suppose, we have an image of a creature that looks similar to cat and dog, but we want to know either it is a cat or dog. So for this identification, we can use the KNN algorithm, as it works on a similarity measure. Our KNN model will find the similar features of the new data set to the cats and dogs images and based on the most similar features it will put it in either cat or dog category.
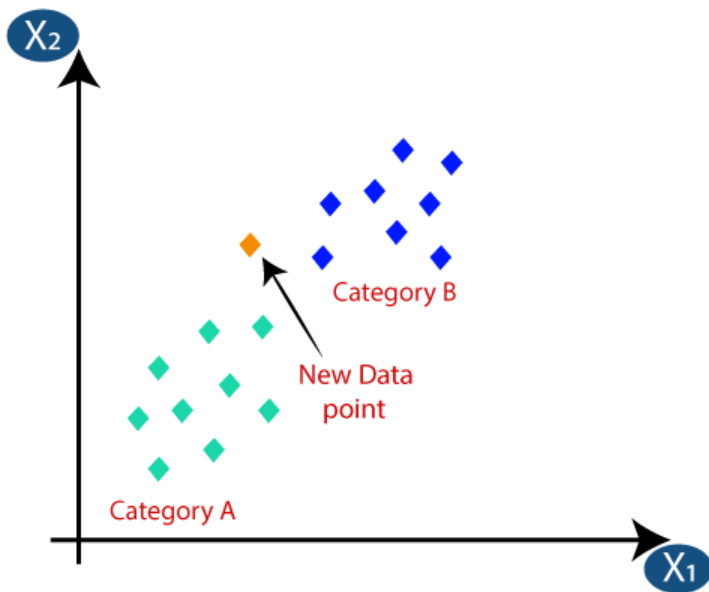
## How does K-NN work?

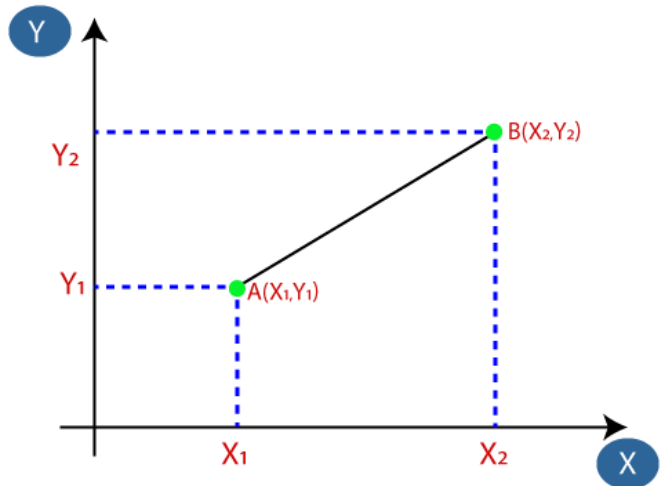The K-NN working can be explained on the basis of the below algorithm:

- o **Step-1:** Select the number K of the neighbors

- o **Step-2:** Calculate the Euclidean distance of **K number of neighbors**

- o **Step-3:** Take the K nearest neighbors as per the calculated Euclidean distance.

- o **Step-4:** Among these k neighbors, count the number of the data points in each category.

- o **Step-5:** Assign the new data points to that category for which the number of the neighbor is maximum.

- o **Step-6:** Our model is ready.

Suppose we have a new data point and we need to put it in the required category. Consider the below image:



- o Firstly, we will choose the number of neighbors, so we will choose the k=5.

- o Next, we will calculate the **Euclidean distance** between the data points. The Euclidean distance is the distance between two points, which we have already studied in geometry. It can be calculated as:

# Module – IV: Supervised Learning –Classification



Euclidean Distance between A₁ and B₂ = $\sqrt{(X_2-X_1)^2+(Y_2-Y_1)^2}$

- o By calculating the Euclidean distance we got the nearest neighbors, as three nearest neighbors in category A and two nearest neighbors in category B. Consider the below image:



- o As we can see the 3 nearest neighbors are from category A, hence this new data point must belong to category A.

## How to select the value of K in the K-NN Algorithm?

Below are some points to remember while selecting the value of K in the K-NN algorithm:

- There is no particular way to determine the best value for "K", so we need to try some values to find the best out of them. The most preferred value for K is 5.

- A very low value for K such as K=1 or K=2, can be noisy and lead to the effects of outliers in the model.

- Large values for K are good, but it may find some difficulties.

## 4. Demonstrate the implementation of K-Nearest Neighbor (KNN) algorithm using python?

**Problem for K-NN Algorithm:** There is a Car manufacturer company that has manufactured a new SUV car. The company wants to give the ads to the users who are interested in buying that SUV. So for this problem, we have a dataset that contains multiple user's information through the social network. The dataset contains lots of information but the **Estimated Salary** and **Age** we will consider for the independent variable and the **Purchased variable** is for the dependent variable. Below is the dataset:

**Steps to implement the K-NN algorithm:**

- Data Pre-processing step
- Fitting the K-NN algorithm to the Training set
- Predicting the test result
- Test accuracy of the result(Creation of Confusion matrix)
- Visualizing the test set result.

**Data Pre-Processing Step:**

The Data Pre-processing step will remain exactly the same as Logistic Regression. Below is the code for it:

```
# importing libraries
import numpy as nm
import matplotlib.pyplot as mtp
import pandas as pd

#importing datasets
data_set= pd.read_csv('user_data.csv')
```

## Module – IV: Supervised Learning –Classification

#Extracting Independent and dependent Variable

x= data_set.iloc[:, [2,3]].values

y= data_set.iloc[:, 4].values

# Splitting the dataset into training and test set.

from sklearn.model_selection import train_test_split

x_train, x_test, y_train, y_test= train_test_split(x, y, test_size= 0.25, random_state=0)

#feature Scaling

from sklearn.preprocessing import StandardScaler

st_x= StandardScaler()

x_train= st_x.fit_transform(x_train)

x_test= st_x.transform(x_test)

From the above output image, we can see that our data is successfully scaled.

**Fitting K-NN classifier to the Training data:**

Now we will fit the K-NN classifier to the training data. To do this we will import the **KNeighborsClassifier** class of **SklearnNeighbors** library. After importing the class, we will create the **Classifier** object of the class. The Parameter of this class will be

o   **n_neighbors:** To define the required neighbors of the algorithm. Usually, it takes 5.

o   **metric='minkowski':** This is the default parameter and it decides the distance between the points.

o   **p=2:** It is equivalent to the standard Euclidean metric.

And then we will fit the classifier to the training data. Below is the code for it:

#Fitting K-NN classifier to the training set

from sklearn.neighbors import KNeighborsClassifier

classifier= KNeighborsClassifier(n_neighbors=5, metric='minkowski', p=2 )

classifier.fit(x_train, y_train)

**Output: By executing the above code, we will get the output as:**

# Module – IV: Supervised Learning –Classification

- o **Predicting the Test Result:**
- o To predict the test set result, we will create a **y_pred** vector as we did in Logistic Regression. Below is the code for it:

#Predicting the test set result
y_pred= classifier.predict(x_test)

**Output:**

The output for the above code will be:

- o **Creating the Confusion Matrix:**
  Now we will create the Confusion Matrix for our K-NN model to see the accuracy of the classifier. Below is the code for it:

#Creating the Confusion matrix
from sklearn.metrics **import** confusion_matrix
cm= confusion_matrix(y_test, y_pred)

In above code, we have imported the confusion_matrix function and called it using the variable cm.

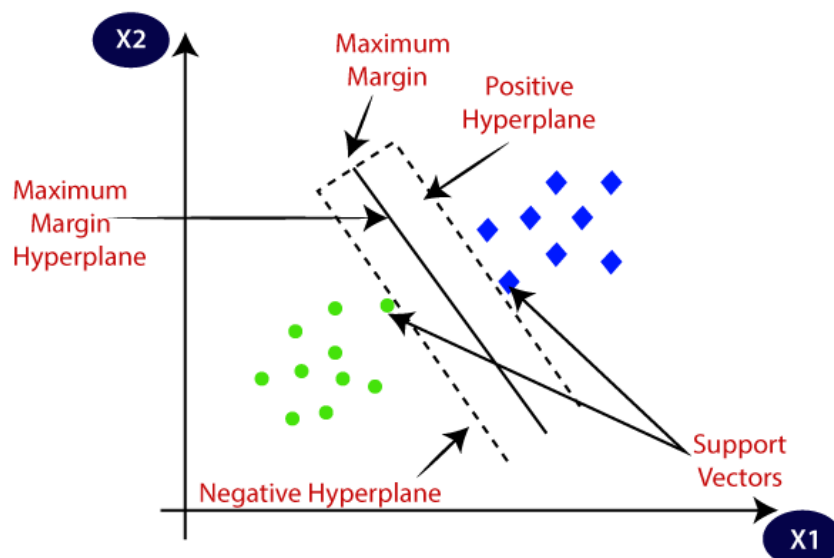**Output:** By executing the above code, we will get the matrix as below:

In the above image, we can see there are 64+29= 93 correct predictions and 3+4= 7 incorrect predictions, whereas, in Logistic Regression, there were 11 incorrect predictions. So we can say that the performance of the model is improved by using the K-NN algorithm.

## 6. Discuss in-detail about Support Vector Machine (SVM) algorithm?

Support Vector Machine or SVM is one of the most popular Supervised Learning algorithms, which is used for Classification as well as Regression problems. However, primarily, it is used for Classification problems in Machine Learning.

The goal of the SVM algorithm is to create the best line or decision boundary that can segregate n-dimensional space into classes so that we can easily put the new data point in the correct category in the future. This best decision boundary is called a hyperplane.

SVM chooses the extreme points/vectors that help in creating the hyperplane. These extreme cases are called as support vectors, and hence algorithm is termed as Support Vector Machine. Consider the below diagram in which there are two different categories that are classified using a decision boundary or hyperplane:



VM algorithm can be used for **Face detection, image classification, text categorization,** etc.

**Types of SVM**

**SVM can be of two types:**

- o **Linear SVM:** Linear SVM is used for linearly separable data, which means if a dataset can be classified into two classes by using a single straight line, then such data is termed as linearly separable data, and classifier is used called as Linear SVM classifier.

o **Non-linear SVM:** Non-Linear SVM is used for non-linearly separated data, which means if a dataset cannot be classified by using a straight line, then such data is termed as non-linear data and classifier used is called as Non-linear SVM classifier.

## Hyperplane and Support Vectors in the SVM algorithm:

**Hyperplane:** There can be multiple lines/decision boundaries to segregate the classes in n-dimensional space, but we need to find out the best decision boundary that helps to classify the data points. This best boundary is known as the hyperplane of SVM.

The dimensions of the hyperplane depend on the features present in the dataset, which means if there are 2 features (as shown in image), then hyperplane will be a straight line. And if there are 3 features, then hyperplane will be a 2-dimension plane.

We always create a hyperplane that has a maximum margin, which means the maximum distance between the data points.
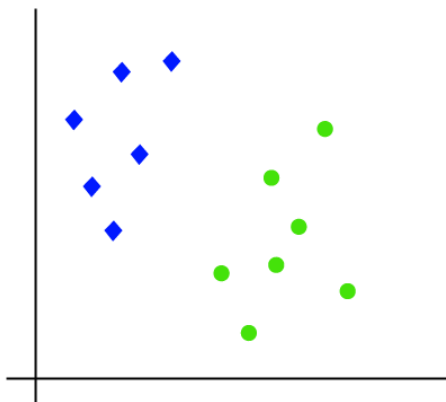
**Support Vectors:**

The data points or vectors that are the closest to the hyperplane and which affect the position of the hyperplane are termed as Support Vector. Since these vectors support the hyperplane, hence called a Support vector.
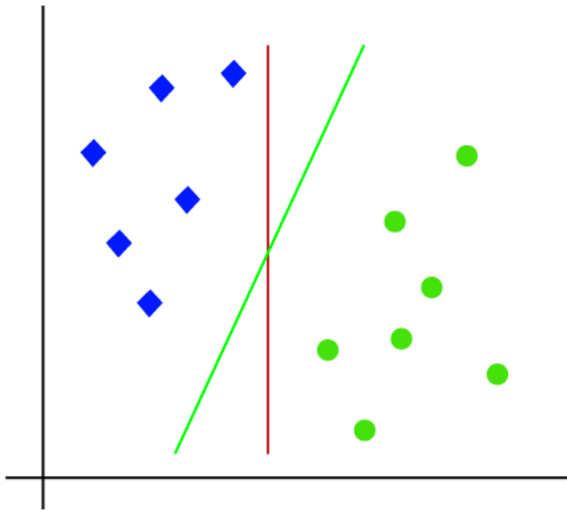
## How does SVM works?

**Linear SVM:**

The working of the SVM algorithm can be understood by using an example. Suppose we have a dataset that has two tags (green and blue), and the dataset has two features x1 and x2. We want a classifier that can classify the pair(x1, x2) of coordinates in either green or blue. Consider the below image:

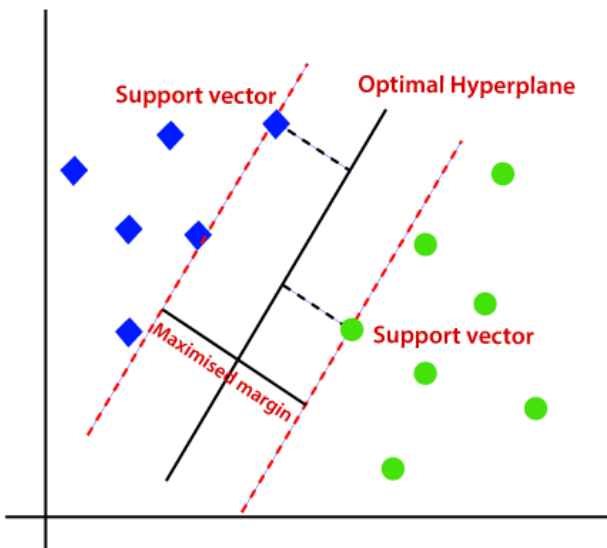So as it is 2-d space so by just using a straight line, we can easily separate these two classes. But there can be multiple lines that can separate these classes. Consider the below image:



Hence, the SVM algorithm helps to find the best line or decision boundary; this best boundary or region is called as a **hyperplane**. SVM algorithm finds the closest point of the lines from both the classes. These points are called support vectors. The distance between the vectors and the hyperplane is called as **margin**. And the goal of SVM is to maximize this margin. The **hyperplane** with maximum margin is called the **optimal hyperplane**.
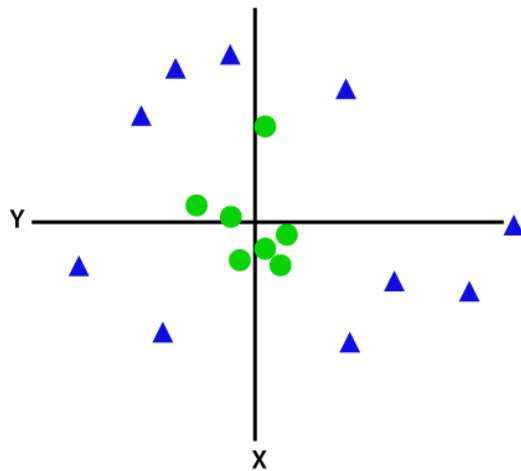


**Non-Linear SVM:**

If data is linearly arranged, then we can separate it by using a straight line, but for non-linear data, we cannot draw a single straight line. Consider the below image:

## Module – IV: Supervised Learning –Classification



So to separate these data points, we need to add one more dimension. For linear data, we have used two dimensions x and y, so for non-linear data, we will add a third dimension z. It can be calculated as:

$$z=x^2 +y^2$$

By adding the third dimension, the sample space will become as below image:



So now, SVM will divide the datasets into classes in the following way. Consider the below image:

**Module – IV: Supervised Learning –Classification**



Since we are in 3-d Space, hence it is looking like a plane parallel to the x-axis. If we convert it in 2d space with z=1, then it will become as:



Hence we get a circumference of radius 1 in case of non-linear data.

## 7. Demonstrate the implementation of Support Vector Machine (SVM) algorithm using python?

### Ans :

Now we will implement the SVM algorithm using Python. Here we will use the same dataset **user_data**, which we have used in Logistic regression and KNN classification.
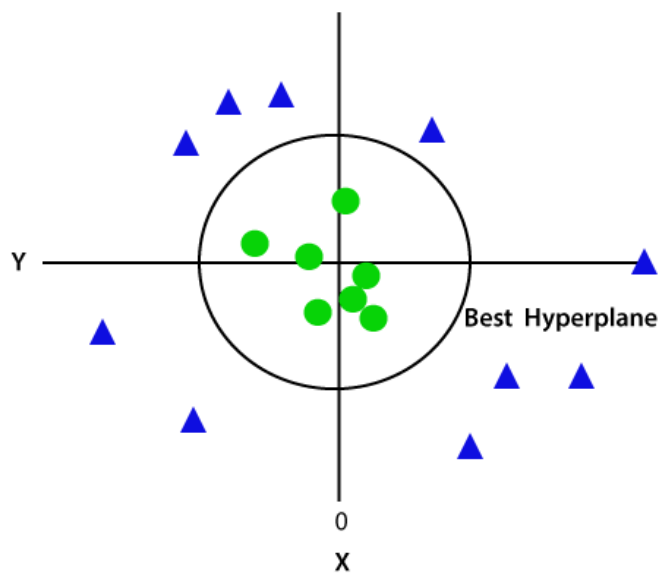
- o **Data Pre-processing step**

Till the Data pre-processing step, the code will remain the same. Below is the code:

```
#Data Pre-processing Step
# importing libraries
import numpy as nm
import matplotlib.pyplot as mtp
import pandas as pd

#importing datasets
data_set= pd.read_csv('user_data.csv')

#Extracting Independent and dependent Variable
x= data_set.iloc[:, [2,3]].values
y= data_set.iloc[:, 4].values

# Splitting the dataset into training and test set.
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test= train_test_split(x, y, test_size= 0.25, random_state=0)
#feature Scaling
from sklearn.preprocessing import StandardScaler
st_x= StandardScaler()
x_train= st_x.fit_transform(x_train)
x_test= st_x.transform(x_test)
```

**Fitting the SVM classifier to the training set:**

## Module – IV: Supervised Learning –Classification

Now the training set will be fitted to the SVM classifier. To create the SVM classifier, we will import **SVC** class from **Sklearn.svm** library. Below is the code for it:

from sklearn.svm **import** SVC # "Support vector classifier"

classifier = SVC(kernel='linear', random_state=0)

classifier.fit(x_train, y_train)

In the above code, we have used **kernel='linear'**, as here we are creating SVM for linearly separable data. However, we can change it for non-linear data. And then we fitted the classifier to the training dataset(x_train, y_train)

**Output:**

```
Out[8]:
SVC(C=1.0, cache_size=200, class_weight=None, coef0=0.0,
decision_function_shape='ovr', degree=3, gamma='auto_deprecated',
    kernel='linear', max_iter=-1, probability=False, random_state=0,
    shrinking=True, tol=0.001, verbose=False)
```

The model performance can be altered by changing the value of **C(Regularization factor), gamma, and kernel**.

**Predicting the test set result:**

Now, we will predict the output for test set. For this, we will create a new vector y_pred. Below is the code for it:

#Predicting the test set result

y_pred= classifier.predict(x_test)

After getting the y_pred vector, we can compare the result of **y_pred** and **y_test** to check the difference between the actual value and predicted value.

**Creating the confusion matrix:**

Now we will see the performance of the SVM classifier that how many incorrect predictions are there as compared to the Logistic regression classifier. To create the confusion matrix, we need to import the **confusion_matrix** function of the sklearn library. After importing the function, we will call it using a new variable **cm**. The function takes two parameters, mainly **y_true**( the actual values) and **y_pred** (the targeted value return by the classifier). Below is the code for it:

#Creating the Confusion matrix

from sklearn.metrics **import** confusion_matrix

cm= confusion_matrix(y_test, y_pred)

As we can see in the above output image, there are 66+24= 90 correct predictions and 8+2= 10 correct predictions. Therefore we can say that our SVM model improved as compared to the Logistic regression model.

## 8. What is Bayes' Theorem? Discuss about Naïve BayesClassifier algorithm?

**Ans :**

Naïve Bayes Classifier Algorithm

- o Naïve Bayes algorithm is a supervised learning algorithm, which is based on **Bayes theorem** and used for solving classification problems.

- o It is mainly used in *text classification* that includes a high-dimensional training dataset.

- o Naïve Bayes Classifier is one of the simple and most effective Classification algorithms which helps in building the fast machine learning models that can make quick predictions.

- o **It is a probabilistic classifier, which means it predicts on the basis of the probability of an object**.

- o Some popular examples of Naïve Bayes Algorithm are **spam filtration, Sentimental analysis, and classifying articles**.

**Why is it called Naïve Bayes?**

The Naïve Bayes algorithm is comprised of two words Naïve and Bayes, Which can be described as:

- o **Naïve**: It is called Naïve because it assumes that the occurrence of a certain feature is independent of the occurrence of other features. Such as if the fruit is identified on the bases of color, shape, and taste, then red, spherical, and sweet fruit is recognized as an apple. Hence each feature individually contributes to identify that it is an apple without depending on each other.

- o **Bayes**: It is called Bayes because it depends on the principle of <u>Bayes' Theorem</u>

## Module – IV: Supervised Learning –Classification

**Bayes' Theorem:**

- o Bayes' theorem is also known as **Bayes' Rule** or **Bayes' law**, which is used to determine the probability of a hypothesis with prior knowledge. It depends on the conditional probability.

- o The formula for Bayes' theorem is given as:

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

**Where,**

**P(A|B) is Posterior probability**: Probability of hypothesis A on the observed event B.
**P(B|A) is Likelihood probability**: Probability of the evidence given that the probability of a hypothesis is true.
**P(A) is Prior Probability**: Probability of hypothesis before observing the evidence.
**P(B) is Marginal Probability**: Probability of Evidence.

**Advantages of Naïve Bayes Classifier:**

- o Naïve Bayes is one of the fast and easy ML algorithms to predict a class of datasets.

- o It can be used for Binary as well as Multi-class Classifications.

- o It performs well in Multi-class predictions as compared to the other Algorithms.

- o It is the most popular choice for **text classification problems**.

**Types of Naïve Bayes Model:**
There are three types of Naive Bayes Model, which are given below:

- o **Gaussian**: The Gaussian model assumes that features follow a normal distribution. This means if predictors take continuous values instead of discrete, then the model assumes that these values are sampled from the Gaussian distribution.

- o **Multinomial**: The Multinomial Naïve Bayes classifier is used when the data is multinomial distributed. It is primarily used for document classification problems, it means a particular document belongs to which category such as Sports, Politics, education, etc. The classifier uses the frequency of words for the predictors.

- o **Bernoulli**: The Bernoulli classifier works similar to the Multinomial classifier, but the predictor variables are the independent Booleans variables. Such as if a particular word is present or not in a document. This model is also famous for document classification tasks.

## 9. Demonstrate the implementation of Naïve Bayes Classifier algorithm using python?

**Python Implementation of the Naïve Bayes algorithm:**

Now we will implement a Naive Bayes Algorithm using Python. So for this, we will use the "**user_data**" **dataset**, which we have used in our other classification model. Therefore we can easily compare the Naive Bayes model with the other models.

**Steps to implement:**

- o   Data Pre-processing step
- o   Fitting Naive Bayes to the Training set
- o   Predicting the test result
- o   Test accuracy of the result(Creation of Confusion matrix)
- o   Visualizing the test set result.

**1) Data Pre-processing step:**
The code for this is given below:

Importing the libraries

**import** numpy as nm

**import** matplotlib.pyplot as mtp

**import** pandas as pd

 # Importing the dataset

dataset = pd.read_csv('user_data.csv')

x = dataset.iloc[:, [2, 3]].values

y = dataset.iloc[:, 4].values

 # Splitting the dataset into the Training set and Test set

**from** sklearn.model_selection **import** train_test_split

x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.25, random_state = 0)

 # Feature Scaling

**from** sklearn.preprocessing **import** StandardScaler

sc = StandardScaler()

x_train = sc.fit_transform(x_train)

x_test = sc.transform(x_test)

In the above code, we have loaded the dataset into our program using "**dataset = pd.read_csv('user_data.csv')**. The loaded dataset is divided into training and test set, and then we have scaled the feature variable.

**2) Fitting Naive Bayes to the Training Set:**

After the pre-processing step, now we will fit the Naive Bayes model to the Training set. Below is the code for it:

```
# Fitting Naive Bayes to the Training set
from sklearn.naive_bayes import GaussianNB
classifier = GaussianNB()
classifier.fit(x_train, y_train)
```

In the above code, we have used the **GaussianNB classifier** to fit it to the training dataset. We can also use other classifiers as per our requirement.

**Output:**
```
Out[6]: GaussianNB(priors=None, var_smoothing=1e-09)
```

**3) Prediction of the test set result:**

Now we will predict the test set result. For this, we will create a new predictor variable **y_pred**, and will use the predict function to make the predictions.

```
# Predicting the Test set results
y_pred = classifier.predict(x_test)
```

The above output shows the result for prediction vector **y_pred** and real vector y_test. We can see that some predications are different from the real values, which are the incorrect predictions.

**4) Creating Confusion Matrix:**

Now we will check the accuracy of the Naive Bayes classifier using the Confusion matrix. Below is the code for it:

```
# Making the Confusion Matrix
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
```
**Output:**

As we can see in the above confusion matrix output, there are 7+3= 10 incorrect predictions, and 65+25=90 correct predictions.

## 10. Discuss in-detail about Decision Tree Classification algorithm?

Decision Tree is a **Supervised learning technique** that can be used for both classification and Regression problems, but mostly it is preferred for solving Classification problems. It is a tree-structured classifier, where **internal nodes represent the features of a dataset, branches represent the decision rules** and **each leaf node represents the outcome.**
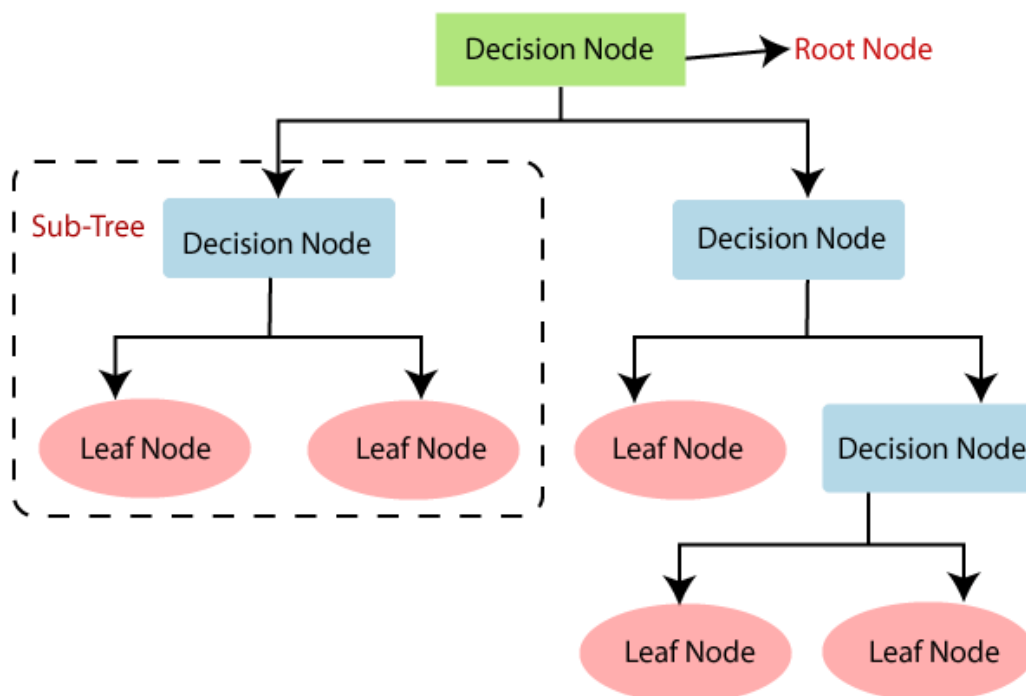
In a Decision tree, there are two nodes, which are the **Decision Node** and **Leaf Node.** Decision nodes are used to make any decision and have multiple branches, whereas Leaf nodes are the output of those decisions and do not contain any further branches.

*It is a graphical representation for getting all the possible solutions to a problem/decision based on given conditions.* It is called a decision tree because, similar to a tree, it starts with the root node, which expands on further branches and constructs a tree-like structure.

In order to build a tree, we use the **CART algorithm,** which stands for **Classification and Regression Tree algorithm.** A decision tree simply asks a question, and based on the answer (Yes/No), it further split the tree into subtrees.

Below diagram explains the general structure of a decision tree:

# Module – IV: Supervised Learning –Classification



## Why use Decision Trees?

There are various algorithms in Machine learning, so choosing the best algorithm for the given dataset and problem is the main point to remember while creating a machine learning model. Below are the two reasons for using the Decision tree:

- o Decision Trees usually mimic human thinking ability while making a decision, so it is easy to understand.

- o The logic behind the decision tree can be easily understood because it shows a tree-like structure.

## Decision Tree Terminologies

**Root Node:** Root node is from where the decision tree starts. It represents the entire dataset, which further gets divided into two or more homogeneous sets.
**Leaf Node:** Leaf nodes are the final output node, and the tree cannot be segregated further after getting a leaf node.
**Splitting:** Splitting is the process of dividing the decision node/root node into sub-nodes according to the given conditions.
**Branch/Sub Tree:** A tree formed by splitting the tree.
**Pruning:** Pruning is the process of removing the unwanted branches from the tree.
**Parent/Child node:** The root node of the tree is called the parent node, and other nodes are called the child nodes.

# Module – IV: Supervised Learning –Classification

**How does the Decision Tree algorithm Work?**

In a decision tree, for predicting the class of the given dataset, the algorithm starts from the root node of the tree. This algorithm compares the values of root attribute with the record (real dataset) attribute and, based on the comparison, follows the branch and jumps to the next node.

For the next node, the algorithm again compares the attribute value with the other sub-nodes and move further. It continues the process until it reaches the leaf node of the tree. The complete process can be better understood using the below algorithm:

- o **Step-1:** Begin the tree with the root node, says S, which contains the complete dataset.

- o **Step-2:** Find the best attribute in the dataset using **Attribute Selection Measure (ASM).**

- o **Step-3:** Divide the S into subsets that contains possible values for the best attributes.

- o **Step-4:** Generate the decision tree node, which contains the best attribute.

- o **Step-5:** Recursively make new decision trees using the subsets of the dataset created in step -3. Continue this process until a stage is reached where you cannot further classify the nodes and called the final node as a leaf node.

**Example:** Suppose there is a candidate who has a job offer and wants to decide whether he should accept the offer or Not. So, to solve this problem, the decision tree starts with the root node (Salary attribute by ASM). The root node splits further into the next decision node (distance from the office) and one leaf node based on the corresponding labels. The next decision node further gets split into one decision node (Cab facility) and one leaf node. Finally, the decision node splits into two leaf nodes (Accepted offers and Declined offer). Consider the below diagram:

## Attribute Selection Measures

While implementing a Decision tree, the main issue arises that how to select the best attribute for the root node and for sub-nodes. So, to solve such problems there is a technique which is called as **Attribute selection measure or ASM.** By this measurement, we can easily select the best attribute for the nodes of the tree. There are two popular techniques for ASM, which are:

- o **Information Gain**
- o **Gini Index**

## 1. Information Gain:

- o Information gain is the measurement of changes in entropy after the segmentation of a dataset based on an attribute.

- o It calculates how much information a feature provides us about a class.

- o According to the value of information gain, we split the node and build the decision tree.

- o A decision tree algorithm always tries to maximize the value of information gain, and a node/attribute having the highest information gain is split first. It can be calculated using the below formula:

*Information Gain= Entropy(S)- [(Weighted Avg) *Entropy(each feature)*

**Entropy:** Entropy is a metric to measure the impurity in a given attribute. It specifies randomness in data. Entropy can be calculated as:

Entropy(s)= -P(yes)log2 P(yes)- P(no) log2 P(no)

**Where,**

- o **S= Total number of samples**
- o **P(yes)= probability of yes**
- o **P(no)= probability of no**

## 2. Gini Index:

- o Gini index is a measure of impurity or purity used while creating a decision tree in the CART(Classification and Regression Tree) algorithm.

- o An attribute with the low Gini index should be preferred as compared to the high Gini index.

- o It only creates binary splits, and the CART algorithm uses the Gini index to create binary splits.

- o Gini index can be calculated using the below formula:

**Gini Index= 1- $\sum_j P_j^2$**

## Pruning: Getting an Optimal Decision tree

## Module – IV: Supervised Learning –Classification

*Pruning is a process of deleting the unnecessary nodes from a tree in order to get the optimal decision tree.*

A too-large tree increases the risk of overfitting, and a small tree may not capture all the important features of the dataset. Therefore, a technique that decreases the size of the learning tree without reducing accuracy is known as Pruning. There are mainly two types of tree **pruning** technology used:

- **Cost Complexity Pruning**
- **Reduced Error Pruning.**

## 11. Demonstrate the implementation of Decision Tree Classification algorithm using python?

Python Implementation of Decision Tree

Now we will implement the Decision tree using Python. For this, we will use the dataset "**user_data.csv**," which we have used in previous classification models. By using the same dataset, we can compare the Decision tree classifier with other classification models such as KNN,SVM,LogisticRegression, etc.

Steps will also remain the same, which are given below:

- **Data Pre-processing step**
- **Fitting a Decision-Tree algorithm to the Training set**
- **Predicting the test result**
- **Test accuracy of the result(Creation of Confusion matrix)**
- **Visualizing the test set result.**

1. Data Pre-Processing Step:

```
# importing libraries
import numpy as nm
import matplotlib.pyplot as mtp
import pandas as pd

#importing datasets
data_set= pd.read_csv('user_data.csv')
```

## Module – IV: Supervised Learning –Classification

#Extracting Independent and dependent Variable

x= data_set.iloc[:, [2,3]].values

y= data_set.iloc[:, 4].values

# Splitting the dataset into training and test set.

from sklearn.model_selection import train_test_split

x_train, x_test, y_train, y_test= train_test_split(x, y, test_size= 0.25, random_state=0)

#feature Scaling

from sklearn.preprocessing import StandardScaler

st_x= StandardScaler()

x_train= st_x.fit_transform(x_train)

x_test= st_x.transform(x_test)

### 2. Fitting a Decision-Tree algorithm to the Training set

Now we will fit the model to the training set. For this, we will import the **DecisionTreeClassifier** class from **sklearn.tree** library. Below is the code for it:

#Fitting Decision Tree classifier to the training set

From sklearn.tree import DecisionTreeClassifier

classifier= DecisionTreeClassifier(criterion='entropy', random_state=0)

classifier.fit(x_train, y_train)

In the above code, we have created a classifier object, in which we have passed two main parameters;

- o **"criterion='entropy':** Criterion is used to measure the quality of split, which is calculated by information gain given by entropy.
- o **random_state=0":** For generating the random states.

Below is the output for this:

```
Out[8]:
DecisionTreeClassifier(class_weight=None, criterion='entropy', max_depth=None,
```

```
max_features=None, max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, presort=False,
random_state=0, splitter='best')
```

## 3. Predicting the test result

Now we will predict the test set result. We will create a new prediction vector **y_pred.** Below is the code for it:

#Predicting the test set result

y_pred= classifier.predict(x_test)

**Output:**

In the below output image, the predicted output and real test output are given. We can clearly see that there are some values in the prediction vector, which are different from the real vector values. These are prediction errors.

## 4. Test accuracy of the result (Creation of Confusion matrix)

In the above output, we have seen that there were some incorrect predictions, so if we want to know the number of correct and incorrect predictions, we need to use the confusion matrix. Below is the code for it:

#Creating the Confusion matrix

from sklearn.metrics **import** confusion_matrix

cm= confusion_matrix(y_test, y_pred)

**Output:**

In the above output image, we can see the confusion matrix, which has **6+3= 9 incorrect predictions** and**62+29=91 correct predictions. Therefore, we can say that compared to other classification models, the Decision Tree classifier made a good prediction.**

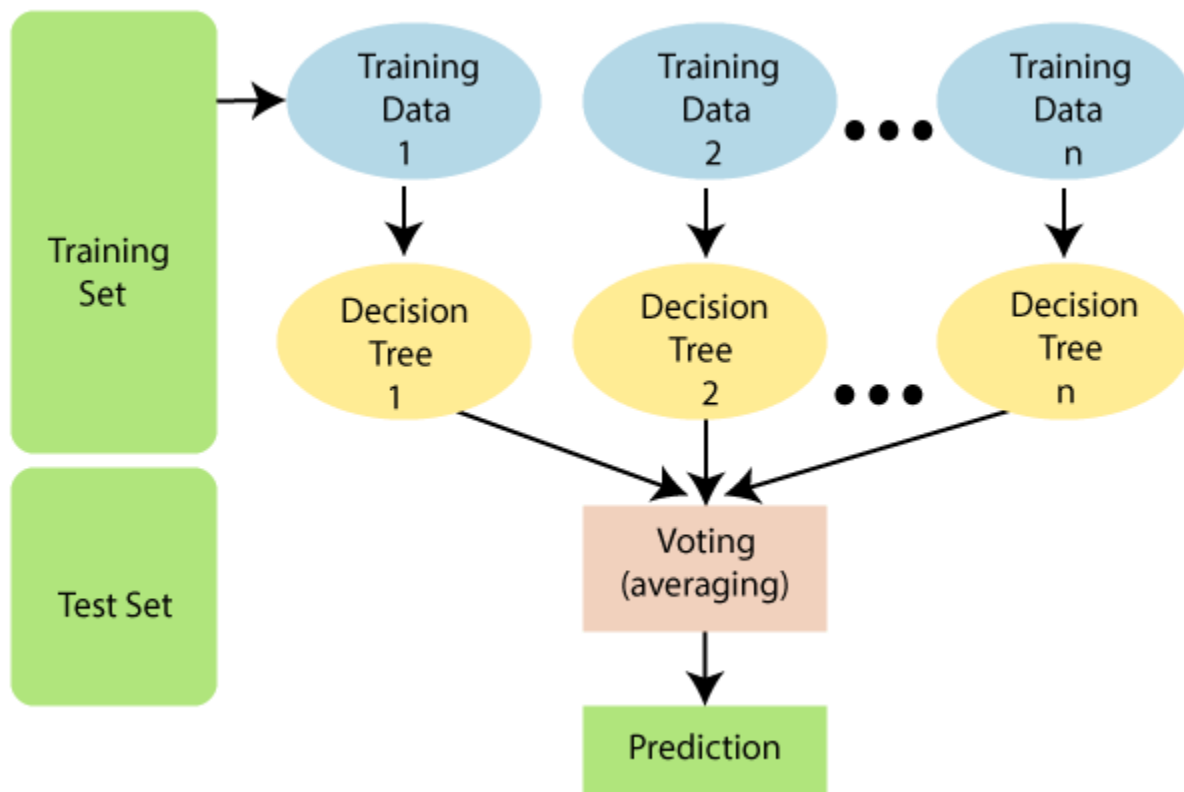## 12. Explain about Random Forest Classification algorithm with an example?

Random Forest Algorithm

Random Forest is a popular machine learning algorithm that belongs to the supervised learning technique. It can be used for both Classification and Regression problems in ML. It is based on the concept of **ensemble learning,** which is a process of *combining multiple classifiers to solve a complex problem and to improve the performance of the model.*

As the name suggests, ***"Random Forest is a classifier that contains a number of decision trees on various subsets of the given dataset and takes the average to improve the predictive accuracy of that dataset."*** Instead of relying on one decision tree, the random forest takes the prediction from each tree and based on the majority votes of predictions, and it predicts the final output.

**The greater number of trees in the forest leads to higher accuracy and prevents the problem of overfitting.**

The below diagram explains the working of the Random Forest algorithm:



Assumptions for Random Forest

# Module – IV: Supervised Learning –Classification

Since the random forest combines multiple trees to predict the class of the dataset, it is possible that some decision trees may predict the correct output, while others may not. But together, all the trees predict the correct output. Therefore, below are two assumptions for a better Random forest classifier:

- There should be some actual values in the feature variable of the dataset so that the classifier can predict accurate results rather than a guessed result.
- The predictions from each tree must have very low correlations.

## How does Random Forest algorithm work?

Random Forest works in two-phase first is to create the random forest by combining N decision tree, and second is to make predictions for each tree created in the first phase.

The Working process can be explained in the below steps and diagram:

**Step-1:** Select random K data points from the training set.

**Step-2:** Build the decision trees associated with the selected data points (Subsets).

**Step-3:** Choose the number N for decision trees that you want to build.

**Step-4:** Repeat Step 1 & 2.

**Step-5:** For new data points, find the predictions of each decision tree, and assign the new data points to the category that wins the majority votes.

The working of the algorithm can be better understood by the below example:

## Python Implementation of Random Forest Algorithm

Now we will implement the Random Forest Algorithm tree using Python. For this, we will use the same dataset "user_data.csv", which we have used in previous classification models. By using the same dataset, we can compare the Random Forest classifier with other classification models such as Decision tree Classifier,KNN,SVM,Logistic Regression,etc.

Implementation Steps are given below:

- Data Pre-processing step
- Fitting the Random forest algorithm to the Training set
- Predicting the test result

# Module – IV: Supervised Learning –Classification

    o   Test accuracy of the result (Creation of Confusion matrix)

    o   Visualizing the test set result.

## 1. Data Pre-Processing Step:

Below is the code for the pre-processing step:

```
# importing libraries
import numpy as nm
import matplotlib.pyplot as mtp
import pandas as pd
#importing datasets
data_set= pd.read_csv('user_data.csv')
#Extracting Independent and dependent Variable
x= data_set.iloc[:, [2,3]].values
y= data_set.iloc[:, 4].values
# Splitting the dataset into training and test set.
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test= train_test_split(x, y, test_size= 0.25, random_state=0)
#feature Scaling
from sklearn.preprocessing import StandardScaler
st_x= StandardScaler()
x_train= st_x.fit_transform(x_train)
x_test= st_x.transform(x_test)
```

## 2. Fitting the Random Forest algorithm to the training set:

Now we will fit the Random forest algorithm to the training set. To fit it, we will import the **RandomForestClassifier** class from the **sklearn.ensemble** library. The code is given below:

```
#Fitting Decision Tree classifier to the training set
from sklearn.ensemble import RandomForestClassifier
classifier= RandomForestClassifier(n_estimators= 10, criterion="entropy")
classifier.fit(x_train, y_train)
```

In the above code, the classifier object takes below parameters:

- o **n_estimators=** The required number of trees in the Random Forest. The default value is 10. We can choose any number but need to take care of the overfitting issue.

- o **criterion=** It is a function to analyze the accuracy of the split. Here we have taken "entropy" for the information gain.

**Output:**

```
RandomForestClassifier(bootstrap=True, class_weight=None, criterion='entropy',
max_depth=None, max_features='auto', max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, n_estimators=10,
n_jobs=None, oob_score=False, random_state=None,
            verbose=0, warm_start=False)
```

## 3. Predicting the Test Set result

Since our model is fitted to the training set, so now we can predict the test result. For prediction, we will create a new prediction vector y_pred. Below is the code for it:

#Predicting the test set result
y_pred= classifier.predict(x_test)

By checking the above prediction vector and test set real vector, we can determine the incorrect predictions done by the classifier.

## 4. Creating the Confusion Matrix

Now we will create the confusion matrix to determine the correct and incorrect predictions. Below is the code for it:

#Creating the Confusion matrix
from sklearn.metrics **import** confusion_matrix
cm= confusion_matrix(y_test, y_pred)

**Output:**

As we can see in the above matrix, there are **4+4= 8 incorrect predictions** and **64+28= 92 correct predictions.**