

LAB MANUAL

CLASS: S.Y. BTech

SEMESTER: IV

Course Outcomes:

C01	Recognize the characteristics of machine learning that make it useful to real-world problems.
C02	Characterize machine learning algorithms as supervised, semi-supervised, and unsupervised.
C03	Design and implement machine learning solutions to classification, regression, and clustering problems;
C04	Be able to evaluate and interpret the results of the algorithms
C05	Effectively use machine learning toolboxes.

Subject Name: Machine Learning Algorithms (UAIP205)

Sr. No.	Name of Experiment	Co's mapping
1	Supervised learning regression: i) Take a dataset and perform preprocessing steps and Split into Training Data set and Test Data set. ii) Perform linear regression analysis. iii) Plot the graphs for Training and test dataset iv) Find out prediction result and check the accuracy v) Improve the accuracy while changing data points.	CO4,CO5
2	Supervised Learning – Classification: Implement Naïve Bayes Classifier on data set of your choice. Test and Compare for Accuracy and Precision.	CO3,CO5
3	Implement k-neighbours classification using Scikit-learn (sklearn)	CO3,CO4

4	<p>Given the following data, which specify classifications for nine combinations of VAR1 and VAR2 predict a classification for a case where VAR1=0.906 and VAR2=0.606, using the result of k-means clustering with 3 means (i.e., centroids)</p> <table> <thead> <tr> <th>VAR1</th><th>VAR2</th><th>CLASS</th></tr> </thead> <tbody> <tr><td>1.713</td><td>1.586</td><td>0</td></tr> <tr><td>0.180</td><td>1.786</td><td>1</td></tr> <tr><td>0.353</td><td>1.240</td><td>1</td></tr> <tr><td>0.940</td><td>1.566</td><td>0</td></tr> <tr><td>1.486</td><td>0.759</td><td>1</td></tr> <tr><td>1.266</td><td>1.106</td><td>0</td></tr> <tr><td>1.540</td><td>0.419</td><td>1</td></tr> <tr><td>0.459</td><td>1.799</td><td>1</td></tr> <tr><td>0.773</td><td>0.186</td><td>1</td></tr> </tbody> </table>	VAR1	VAR2	CLASS	1.713	1.586	0	0.180	1.786	1	0.353	1.240	1	0.940	1.566	0	1.486	0.759	1	1.266	1.106	0	1.540	0.419	1	0.459	1.799	1	0.773	0.186	1	CO3,CO4,CO5
VAR1	VAR2	CLASS																														
1.713	1.586	0																														
0.180	1.786	1																														
0.353	1.240	1																														
0.940	1.566	0																														
1.486	0.759	1																														
1.266	1.106	0																														
1.540	0.419	1																														
0.459	1.799	1																														
0.773	0.186	1																														
5	Implement Linear Regression using Scikit-learn (sklearn)	CO4,CO5																														
6	Implement Naïve Bayes theorem to classify the English text	CO3,CO4,CO5																														
7	Unsupervised Learning: Implement K-Means Clustering and Hierarchical clustering on proper data set of your choice. Compare their Convergence	CO3,CO5																														
8	Implement Support Vector Machine algorithm using suitable dataset.	CO3,CO4,CO5																														

Assignment No:- 1

Title of Assignment:

Supervised learning - regression

- i) Take a dataset and perform preprocessing steps and Split into Training Data set and Test Data set.
- ii) Perform linear regression analysis.
- iii) Plot the graphs for Training and test dataset
- iv) Find out prediction result and check the accuracy
- v) Improve the accuracy while changing data points.

During assignment students will be able to:

- Learn Supervised Learning.
- Implementation of Linear Regression.

Prerequisites: Knowledge of python programming and machine learning libraries.

Concepts to be used: Linear Regression.

Input: Any dataset

Output: Successful implementation of Linear Regression.

Relevant Theory / Literature Survey:

The linear regression algorithm in machine learning is a supervised learning technique to approximate the mapping function to get the best predictions. In this article, we will learn about linear regression for machine learning. The following topics are discussed in this practical.

- [What is Regression?](#)
- [Types of Regression](#)
- [What is Linear Regression?](#)
- [Linear Regression Terminologies](#)
- [Advantages And Disadvantages Of Linear Regression](#)
- [Linear Regression Use Cases](#)
- [Use case – Linear Regression Implementation](#)

What is Regression?

The main goal of regression is the construction of an efficient model to predict the dependent attributes from a bunch of attribute variables. A regression problem is when the output variable is either real or a continuous value i.e salary, weight, area, etc.

We can also define regression as a statistical means that is used in applications like housing, investing, etc. It is used to predict the relationship between a dependent variable and a bunch of independent variables. Let us take a look at various types of regression techniques.



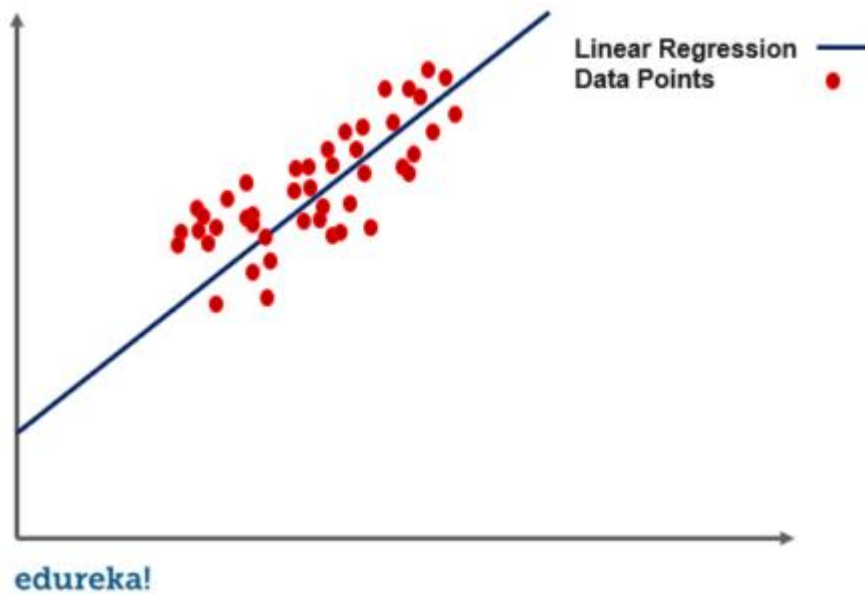
Types Of Regression

The following are types of regression.

1. **Simple Linear Regression**
2. **Polynomial Regression**
3. **Support Vector Regression**
4. **Decision Tree Regression**
5. **Random Forest Regression**

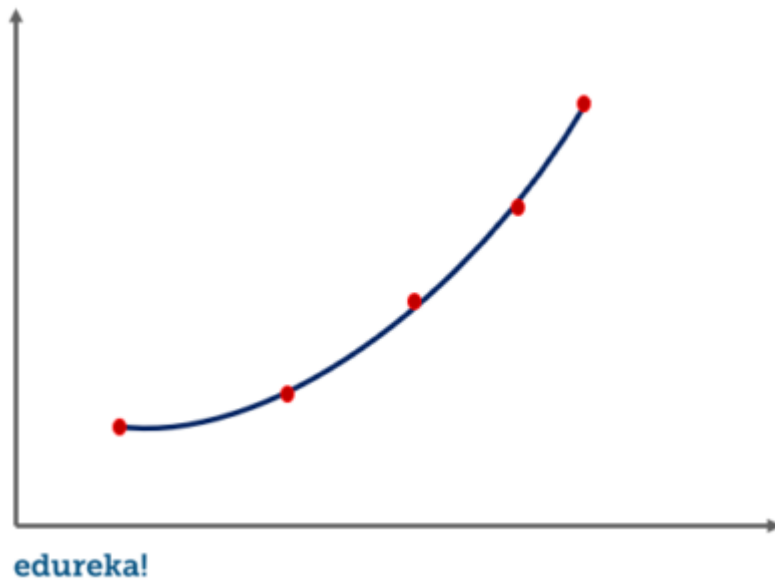
Simple Linear Regression

One of the most interesting and common regression technique is simple linear regression. In this, we predict the outcome of a dependent variable based on the independent variables, the relationship between the variables is linear. Hence, the word linear regression.



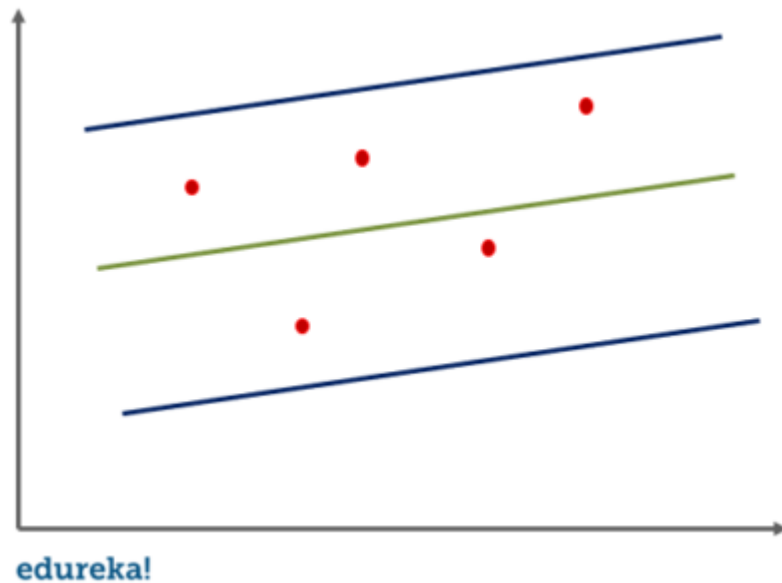
Polynomial Regression

In this regression technique, we transform the original features into polynomial features of a given degree and then perform regression on it.



Support Vector Regression

For support vector machine regression or SVR, we identify a hyperplane with maximum margin such that the maximum number of data points are within those margins. It is quite similar to the support vector machine classification algorithm.



Decision Tree Regression

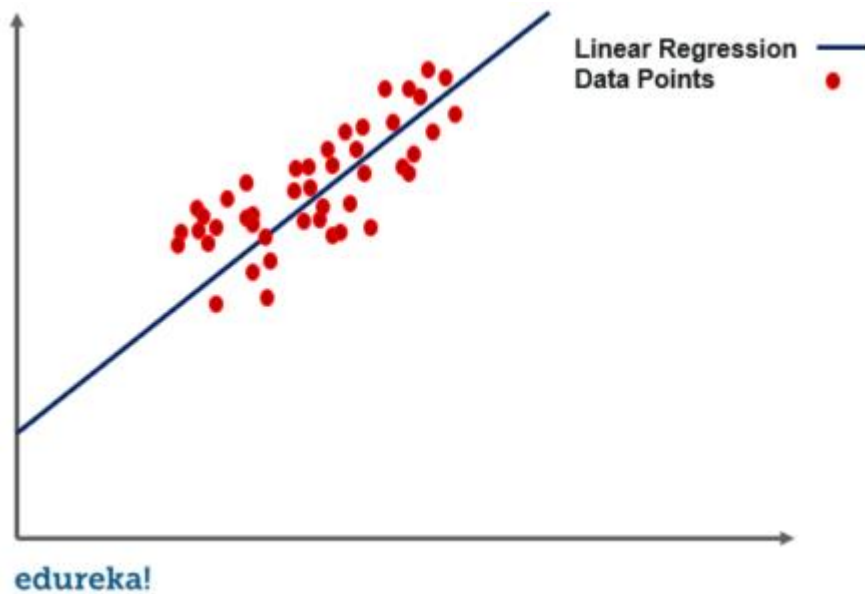
A decision tree can be used for both regression and classification. In the case of regression, we use the ID3 algorithm (Iterative Dichotomiser 3) to identify the splitting node by reducing the standard deviation.

Random Forest Regression

In random forest regression, we ensemble the predictions of several decision tree regressions. Now that we know about different types of regression let us take a look at simple linear regression in detail.

What is Linear Regression?

Simple linear regression is a regression technique in which the independent variable has a linear relationship with the dependent variable. The straight line in the diagram is the best fit line. The main goal of the simple linear regression is to consider the given data points and plot the best fit line to fit the model in the best way possible.



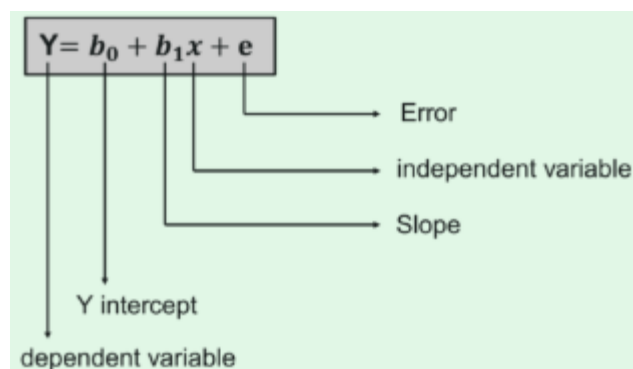
Before moving on to how the linear regression algorithm works, let us take a look at a few important terminologies in simple linear regression.

Linear Regression Terminologies

The following terminologies are important to be familiar with before moving on to the linear regression algorithm.

Cost Function

The best fit line can be based on the linear equation given below.



- The dependent variable that is to be predicted is denoted by Y.
- A line that touches the y-axis is denoted by the intercept b_0 .
- b_1 is the slope of the line, x represents the independent variables that determine the prediction of Y.
- The error in the resultant prediction is denoted by e.

The cost function provides the best possible values for b_0 and b_1 to make the best fit line for the data points. We do it by converting this problem into a minimization problem to get the best values for b_0 and b_1 . The error is minimized in this problem between the actual value and the predicted value.

$$\text{minimize } \frac{1}{n} \sum_{i=1}^n (\text{pred}_i - y_i)^2$$

$$J = \frac{1}{n} \sum_{i=1}^n (\text{pred}_i - y_i)^2$$

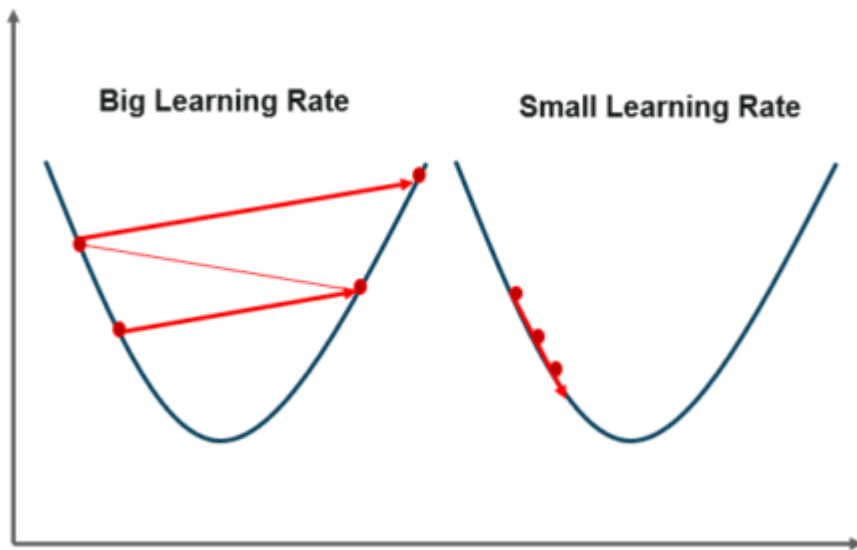
We choose the function above to minimize the error. We square the error difference and sum the error over all data points, the division between the total number of data points. Then, the produced value provides the averaged square error over all data points.

It is also known as MSE(Mean Squared Error), and we change the values of b_0 and b_1 so that the MSE value is settled at the minimum.

Gradient Descent

The next important terminology to understand linear regression is **gradient descent**. It is a method of updating b_0 and b_1 values to reduce the MSE. The idea behind this is to keep iterating the b_0 and b_1 values until we reduce the MSE to the minimum.

To update b_0 and b_1 , we take gradients from the cost function. To find these gradients, we take partial derivatives with respect to b_0 and b_1 . These partial derivatives are the gradients and are used to update the values of b_0 and b_1 .



edureka!

A smaller learning rate takes closer to the minimum, but it takes more time and in case of a larger learning rate. The time taken is sooner but there is a chance to overshoot the minimum value. Now that we are through with the terminologies in linear regression, let us take a look at a few advantages and disadvantages of linear regression for machine learning.

Advantages And Disadvantages

Advantages	Disadvantages
Linear regression performs exceptionally well for linearly separable data	The assumption of linearity between dependent and independent variables
Easier to implement, interpret and efficient to train	It is often quite prone to noise and overfitting
It handles overfitting pretty well using dimensionally reduction techniques, regularization, and cross-validation	Linear regression is quite sensitive to outliers
One more advantage is the extrapolation beyond a specific data set	It is prone to multicollinearity

Linear Regression Use Cases

- Sales Forecasting
- Risk Analysis
- Housing Applications To Predict the prices and other factors
- Finance Applications To Predict Stock prices, investment evaluation, etc.

The basic idea behind linear regression is to find the relationship between the dependent and independent variables. It is used to get the best fitting line that would predict the outcome with the least error. We can use linear regression in simple real-life situations, like predicting the SAT scores with regard to the number of hours of study and other decisive factors.

With this in mind, let us take a look at a use case.

Use Case – Implementing Linear Regression

The process takes place in the following steps:

1. Loading the Data
2. Exploring the Data
3. Slicing The Data
4. Train and Split Data
5. Generate The Model
6. Evaluate The accuracy

Let us get into the details of each of the steps to implement linear regression.

1. Loading The Data

We can start with the basic diabetes data set that is already present in the sklearn(scikit-learn) data sets module to begin our journey with linear regression.

```
1 from sklearn import datasets
2
3 disease = datasets.load_diabetes()
4 print(disease)
```

Output:

```
Run: usecase x
60., 174., 259., 178., 128., 96., 126., 288., 88., 292., 71.,
197., 186., 25., 84., 96., 195., 53., 217., 172., 131., 214.,
59., 70., 220., 268., 152., 47., 74., 295., 101., 151., 127.,
237., 225., 81., 151., 107., 64., 138., 185., 265., 101., 137.,
143., 141., 79., 292., 178., 91., 116., 86., 122., 72., 129.,
142., 90., 158., 39., 196., 222., 277., 99., 196., 202., 155.,
77., 191., 70., 73., 49., 65., 263., 248., 296., 214., 185.,
78., 93., 252., 150., 77., 208., 77., 108., 160., 53., 220.,
154., 259., 90., 246., 124., 67., 72., 257., 262., 275., 177.,
71., 47., 187., 125., 78., 51., 258., 215., 303., 243., 91.,
150., 310., 153., 346., 63., 89., 50., 39., 103., 308., 116.,
145., 74., 45., 115., 264., 87., 202., 127., 182., 241., 66.,
94., 283., 64., 102., 200., 265., 94., 230., 181., 156., 233.,
60., 219., 80., 68., 332., 248., 84., 200., 55., 85., 89.,
31., 129., 83., 275., 65., 198., 236., 253., 124., 44., 172.,
114., 142., 109., 180., 144., 163., 147., 97., 220., 190., 109.,
191., 122., 230., 242., 248., 249., 192., 131., 237., 78., 135.,
244., 199., 270., 164., 72., 96., 306., 91., 214., 95., 216.,
263., 178., 113., 200., 139., 139., 88., 148., 88., 243., 71.,
77., 109., 272., 60., 54., 221., 90., 311., 281., 182., 321.,
58., 262., 206., 233., 242., 123., 167., 63., 197., 71., 168.,
140., 217., 121., 235., 245., 40., 52., 104., 132., 88., 69.,
219., 72., 201., 110., 51., 277., 63., 118., 69., 273., 258.,
43., 198., 242., 232., 175., 93., 168., 275., 293., 281., 72.,
140., 189., 181., 209., 136., 261., 113., 131., 174., 257., 55.,
84., 42., 146., 212., 233., 91., 111., 152., 120., 67., 310.,
94., 183., 66., 173., 72., 49., 64., 48., 178., 104., 132.,
220., 57.]), 'DESCR': '.._diabetes_dataset:\n\nDiabetes dataset\n-----\n\nTen baseline variables, age, sex, body mass index, average blood\npressure, and six bl

Process finished with exit code 0
```

2. Exploring The Data

After we are done loading the data, we can start exploring by simply checking the labels by using the following code.

```
1 print(disease.keys())
```

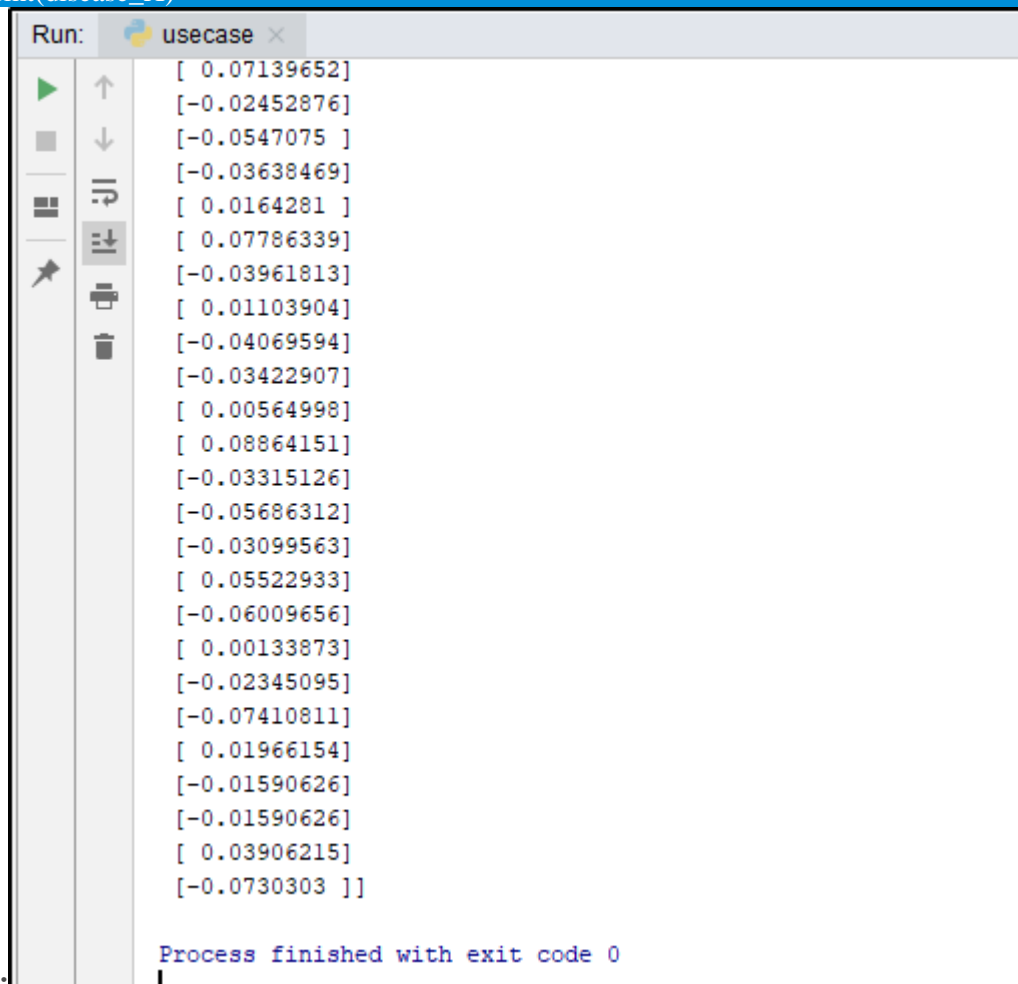
Output:

```
Run: usecase x
C:\Users\Waseem\PycharmProjects\classification\venv\Scripts\python.exe C:/Users/Waseem/PycharmProjects/classification/usecase.py
dict_keys(['data', 'target', 'DESCR', 'feature_names', 'data_filename', 'target_filename'])

Process finished with exit code 0
```

The above code gives all the labels from the data set, after this, we can slice the data so that we can plot the line in the end. We will also use all the data points, for now, we will slice column 2 from the data.

```
1 import numpy as np
2
3 disease_X = disease.data[:, np.newaxis,2]
4 print(disease_X)
```



```
Run: usecase x
[ 0.07139652]
[-0.02452876]
[-0.0547075 ]
[-0.03638469]
[ 0.0164281 ]
[ 0.07786339]
[-0.03961813]
[ 0.01103904]
[-0.04069594]
[-0.03422907]
[ 0.00564998]
[ 0.08864151]
[-0.03315126]
[-0.05686312]
[-0.03099563]
[ 0.05522933]
[-0.06009656]
[ 0.00133873]
[-0.02345095]
[-0.07410811]
[ 0.01966154]
[-0.01590626]
[-0.01590626]
[ 0.03906215]
[-0.0730303 ]

Process finished with exit code 0
```

Output:

After this step, we will split the data into train and test set.

3. Splitting The Data

```
1 disease_X_train = disease_X[:-30]
2 disease_X_test = disease_X[-20:]
3
4 disease_Y_train = disease.target[:-30]
5 disease_Y_test = disease.target[-20:]
```

The next part involves generating the model, which will include importing linear_model from sklearn.

4. Generating the model

```
1 from sklearn import linear_model
2
3 reg = linear_model.LinearRegression()
4 reg.fit(disease_X_train,disease_Y_train)
```

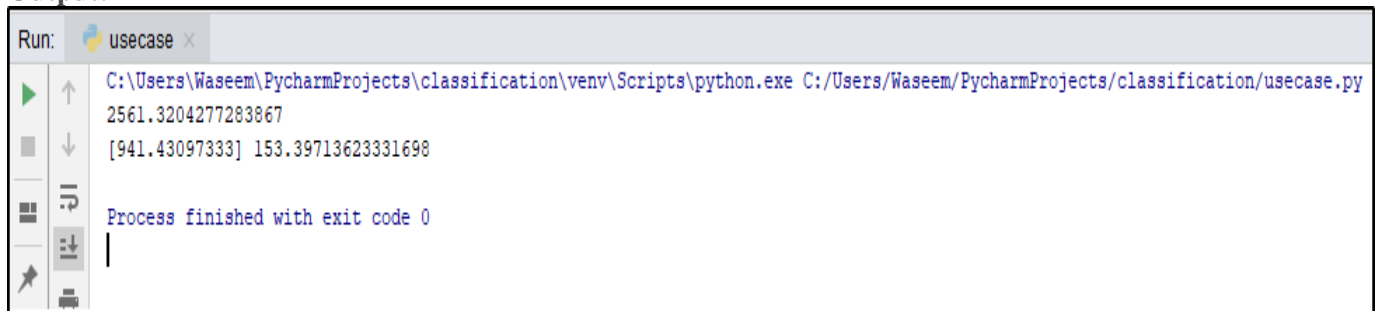
```
5
6 y_predict = reg.predict(disease_X_test)
```

To evaluate the accuracy of the model, we will use the mean squared error from the scikit-learn.

5. Evaluation

```
1 accuracy = mean_squared_error(disease_Y_test,y_predict,)
2
3 print(accuracy)
4
5 weights = reg.coef_
6 intercept = reg.intercept_
7 print(weights,intercept)
```

Output:

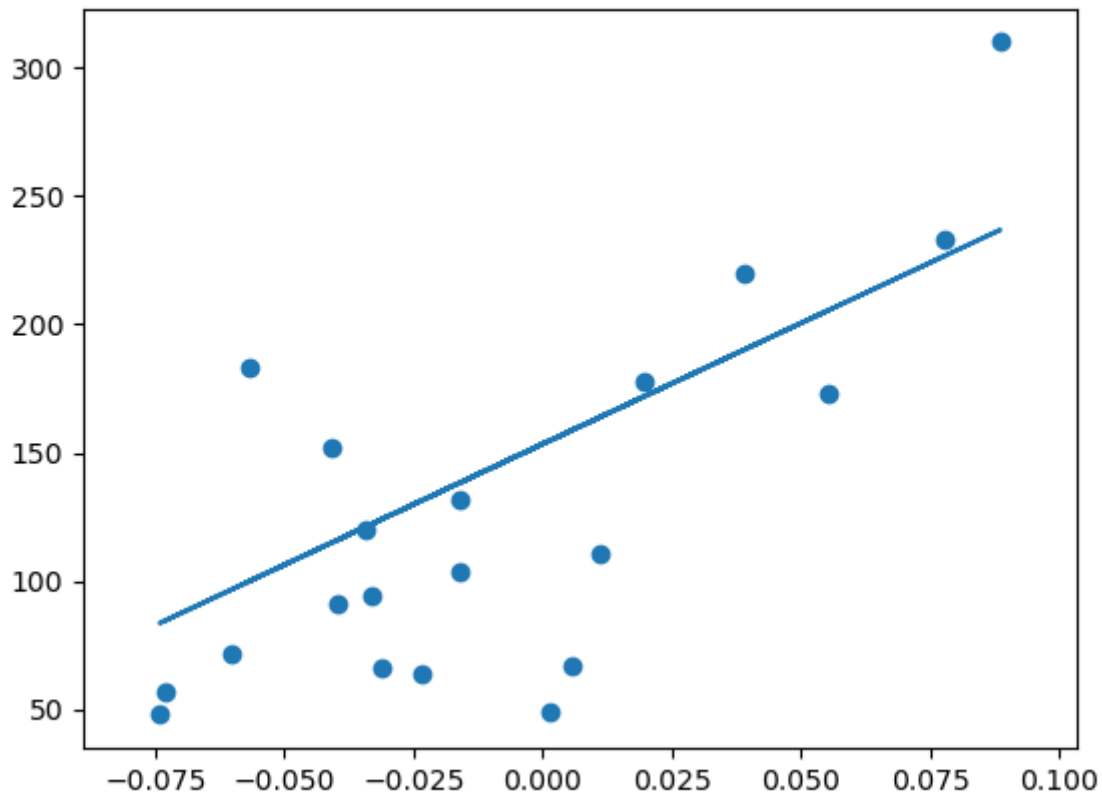


```
Run: usecase x
C:\Users\Waseem\PycharmProjects\classification\venv\Scripts\python.exe C:/Users/Waseem/PycharmProjects/classification/usecase.py
2561.3204277283867
[941.43097333] 153.39713623331698
Process finished with exit code 0
```

To be more clear on how the data points look like on the graph, let us plot the graphs as well.

```
1 import matplotlib.pyplot as plt
2
3 plt.scatter(disease_X_test, disease_Y_test)
4 plt.plot(disease_X_test,y_predict)
5 plt.show()
```

Output:



To get a more accurate model in this scenario, we can use the whole data instead of just the column 2. That would give an accuracy as shown below:

- 1 #make a little change in the code above, and remove the plotting code to avoid errors
- 2 disease_X = disease.data

```
Run: usecase x
C:\Users\Waseem\PycharmProjects\classification\venv\Scripts\python.exe C:/Users/Waseem/PycharmProjects/classification/usecase.py
2004.3086353199665
[ -1.16924976 -237.18461486 518.30606657 309.04865826 -763.14121622
458.90999325 80.62441437 174.32183366 721.49712065 79.19307944] 153.05827988224112
Process finished with exit code 0
```

Output:

FAQ'S

1. What is meant by Regression?
2. What is meant by Linear Regression?
3. What is cost function?
4. What is dependence & independent variable?
5. What is difference between classification and regression?

Outcome:

With the completion of this assignment the students will be able to implement linear regression on given dataset.

Conclusion:

Assignment No: - 2

Title of Assignment:

Supervised Learning – Classification:

Implement Naïve Bayes Classifier on data set of your choice. Test and Compare for Accuracy and Precision.

During assignment students will be able to:

- Learn Supervised Learning.
- Implementation of Naïve Bayes Classifier.

Prerequisites: Knowledge of python programming and machine learning libraries.

Concepts to be used: Naïve Bayes Classifier

Input: Any dataset

Output: Successful implementation of Naïve Bayes Classifier.

Relevant Theory / Literature Survey:

Naïve Bayes Classifier Algorithm

- Naïve Bayes algorithm is a supervised learning algorithm, which is based on **Bayes theorem** and used for solving classification problems.
- It is mainly used in *text classification* that includes a high-dimensional training dataset.
- 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.
- **It is a probabilistic classifier, which means it predicts on the basis of the probability of an object.**
- 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:

- **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.

- **Bayes:** It is called Bayes because it depends on the principle of [Bayes' Theorem](#).

Bayes' Theorem:

- 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.
- 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.

Working of Naïve Bayes' Classifier:

Working of Naïve Bayes' Classifier can be understood with the help of the below example:

Suppose we have a dataset of **weather conditions** and corresponding target variable "**Play**". So using this dataset we need to decide that whether we should play or not on a particular day according to the weather conditions. So to solve this problem, we need to follow the below steps:

1. Convert the given dataset into frequency tables.
2. Generate Likelihood table by finding the probabilities of given features.
3. Now, use Bayes theorem to calculate the posterior probability.

Problem: If the weather is sunny, then the Player should play or not?

Solution: To solve this, first consider the below dataset:

	Outlook	Play
0	Rainy	Yes
1	Sunny	Yes
2	Overcast	Yes
3	Overcast	Yes
4	Sunny	No
5	Rainy	Yes
6	Sunny	Yes
7	Overcast	Yes
8	Rainy	No
9	Sunny	No
10	Sunny	Yes
11	Rainy	No
12	Overcast	Yes
13	Overcast	Yes

Frequency table for the Weather Conditions:

Weather	Yes	No
Overcast	5	0
Rainy	2	2
Sunny	3	2
Total	10	5

Likelihood table weather condition:

Weather	No	Yes	
Overcast	0	5	5/14= 0.35
Rainy	2	2	4/14=0.29
Sunny	2	3	5/14=0.35
All	4/14=0.29	10/14=0.71	

Applying Bayes'theorem:

$$P(\text{Yes}|\text{Sunny}) = P(\text{Sunny}|\text{Yes}) * P(\text{Yes}) / P(\text{Sunny})$$

$$P(\text{Sunny}|\text{Yes}) = 3/10 = 0.3$$

$$P(\text{Sunny}) = 0.35$$

$$P(\text{Yes}) = 0.71$$

$$\text{So } P(\text{Yes}|\text{Sunny}) = 0.3 * 0.71 / 0.35 = \mathbf{0.60}$$

$$P(\text{No}|\text{Sunny}) = P(\text{Sunny}|\text{No}) * P(\text{No}) / P(\text{Sunny})$$

$$P(\text{Sunny}|\text{NO}) = 2/4 = 0.5$$

$$P(\text{No}) = 0.29$$

$$P(\text{Sunny}) = 0.35$$

$$\text{So } P(\text{No}|\text{Sunny}) = 0.5 * 0.29 / 0.35 = \mathbf{0.41}$$

So as we can see from the above calculation that **$P(\text{Yes}|\text{Sunny}) > P(\text{No}|\text{Sunny})$**

Hence on a Sunny day, Player can play the game.

Advantages of Naïve Bayes Classifier:

- Naïve Bayes is one of the fast and easy ML algorithms to predict a class of datasets.
- It can be used for Binary as well as Multi-class Classifications.
- It performs well in Multi-class predictions as compared to the other Algorithms.
- It is the most popular choice for **text classification problems**.

Disadvantages of Naïve Bayes Classifier:

- Naive Bayes assumes that all features are independent or unrelated, so it cannot learn the relationship between features.

Applications of Naïve Bayes Classifier:

- It is used for **Credit Scoring**.
- It is used in **medical data classification**.
- It can be used in **real-time predictions** because Naïve Bayes Classifier is an eager learner.
- It is used in Text classification such as **Spam filtering** and **Sentiment analysis**.

Types of Naïve Bayes Model:

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

- **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.
- **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.
- **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.

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:

- Data Pre-processing step
- Fitting Naive Bayes 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:

In this step, we will pre-process/prepare the data so that we can use it efficiently in our code. It is similar as we did in [data-pre-processing](#). The code for this is given below:

1. Importing the libraries
2. **import** numpy as nm
3. **import** matplotlib.pyplot as mtp
4. **import** pandas as pd
- 5.
6. **# Importing the dataset**
7. `dataset = pd.read_csv('user_data.csv')`
8. `x = dataset.iloc[:, [2, 3]].values`
9. `y = dataset.iloc[:, 4].values`
- 10.
11. **# Splitting the dataset into the Training set and Test set**
12. **from** sklearn.model_selection **import** train_test_split
13. `x_train, x_test, y_train, y_test = train_test_split(x, y, test_size = 0.25, random_state = 0)`
- 14.
15. **# Feature Scaling**
16. **from** sklearn.preprocessing **import** StandardScaler
17. `sc = StandardScaler()`
18. `x_train = sc.fit_transform(x_train)`
19. `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.

The output for the dataset is given as:

Index	User ID	Gender	Age	EstimatedSalary	Purchased
0	15624510	Male	19	19000	0
1	15810944	Male	35	20000	0
2	15668575	Female	26	43000	0
3	15603246	Female	27	57000	0
4	15804002	Male	19	76000	0
5	15728773	Male	27	58000	0
6	15598044	Female	27	84000	0
7	15694829	Female	32	150000	1
8	15600575	Male	25	33000	0
9	15727311	Female	35	65000	0
10	15570769	Female	26	80000	0
11	15606274	Female	26	52000	0
12	15746139	Male	20	86000	0
13	15704987	Male	32	18000	0
14	15628972	Male	18	82000	0
15	15697686	Male	29	80000	0
16	15733883	Male	47	25000	1
17	15617482	Male	45	26000	1
18	15704583	Male	46	28000	1
19	15621083	Female	48	29000	1

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:

1. `# Fitting Naive Bayes to the Training set`
2. `from sklearn.naive_bayes import GaussianNB`
3. `classifier = GaussianNB()`
4. `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:

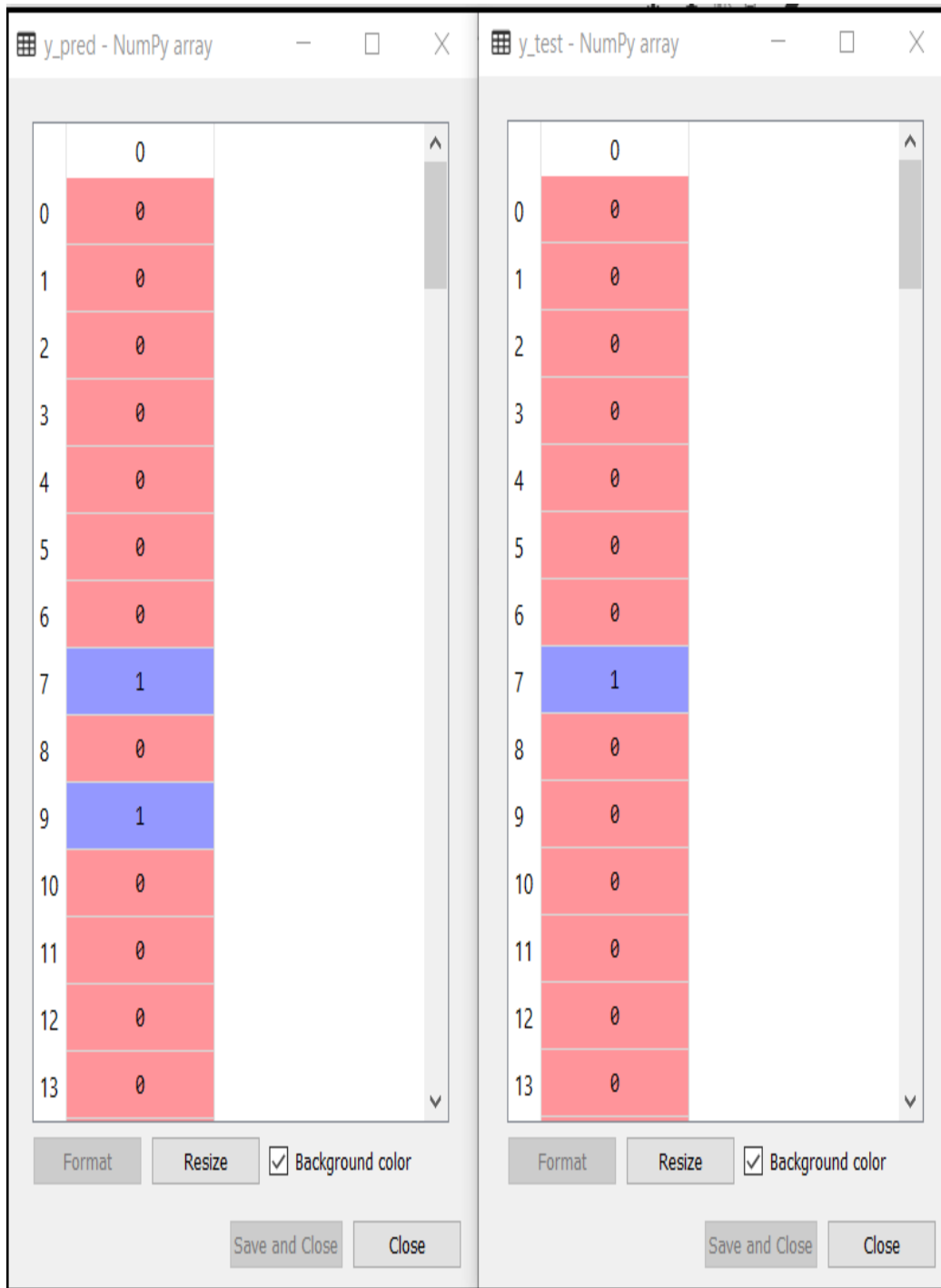
```
Out[6]: GaussianNB(priors=None, var_smoothing=1e-09)
```

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.

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

Output:



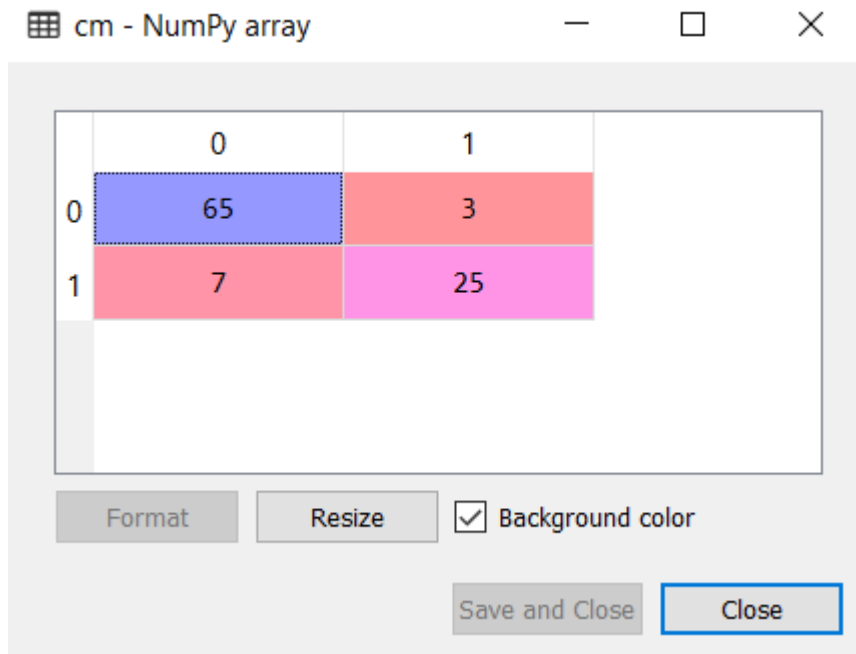
The above output shows the result for prediction vector **y_pred** and real vector y_test. We can see that some predictions are different from the real values, which are the incorrect predictions.

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:

1. `# Making the Confusion Matrix`
2. `from sklearn.metrics import confusion_matrix`
3. `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.

5) Visualizing the training set result:

Next we will visualize the training set result using Naïve Bayes Classifier. Below is the code for it:

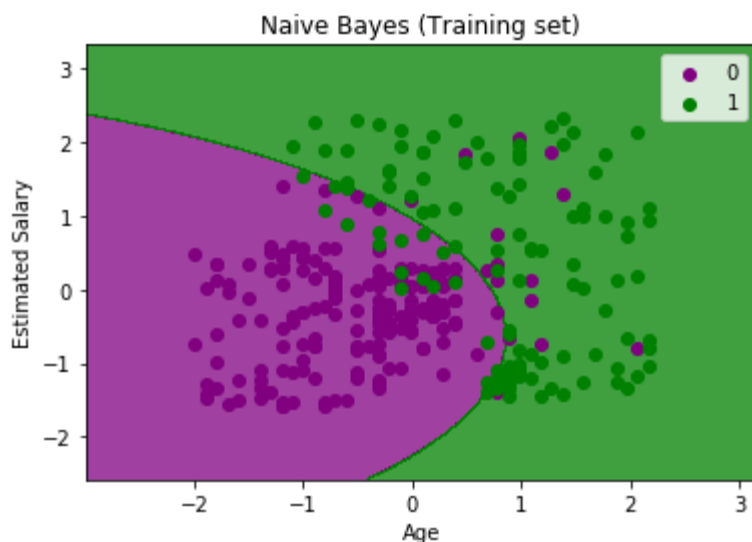
1. `# Visualising the Training set results`
2. `from matplotlib.colors import ListedColormap`
3. `x_set, y_set = x_train, y_train`
4. `X1, X2 = nm.meshgrid(nm.arange(start = x_set[:, 0].min() - 1, stop = x_set[:, 0].max() + 1, step = 0.01),`
5. `nm.arange(start = x_set[:, 1].min() - 1, stop = x_set[:, 1].max() + 1, step = 0.01))`
6. `mtp.contourf(X1, X2, classifier.predict(nm.array([X1.ravel(), X2.ravel()]).T).reshape(X1.shape),`
7. `alpha = 0.75, cmap = ListedColormap(('purple', 'green')))`
8. `mtp.xlim(X1.min(), X1.max())`
9. `mtp.ylim(X2.min(), X2.max())`
10. `for i, j in enumerate(nm.unique(y_set)):`

```

11. mtp.scatter(x_set[y_set == j, 0], x_set[y_set == j, 1],
12.             c = ListedColormap(('purple', 'green'))(i), label = j)
13. mtp.title('Naive Bayes (Training set)')
14. mtp.xlabel('Age')
15. mtp.ylabel('Estimated Salary')
16. mtp.legend()
17. mtp.show()

```

Output:



In the above output we can see that the Naïve Bayes classifier has segregated the data points with the fine boundary. It is Gaussian curve as we have used **GaussianNB** classifier in our code.

6) Visualizing the Test set result:

```

1. # Visualising the Test set results
2. from matplotlib.colors import ListedColormap
3. x_set, y_set = x_test, y_test
4. X1, X2 = nm.meshgrid(nm.arange(start = x_set[:, 0].min() - 1, stop = x_set[:, 0].max() + 1, step = 0.01),
5.                      nm.arange(start = x_set[:, 1].min() - 1, stop = x_set[:, 1].max() + 1, step = 0.01))
6. mtp.contourf(X1, X2, classifier.predict(nm.array([X1.ravel(), X2.ravel()]).T).reshape(X1.shape),
7.              alpha = 0.75, cmap = ListedColormap(('purple', 'green')))
8. mtp.xlim(X1.min(), X1.max())
9. mtp.ylim(X2.min(), X2.max())
10. for i, j in enumerate(nm.unique(y_set)):

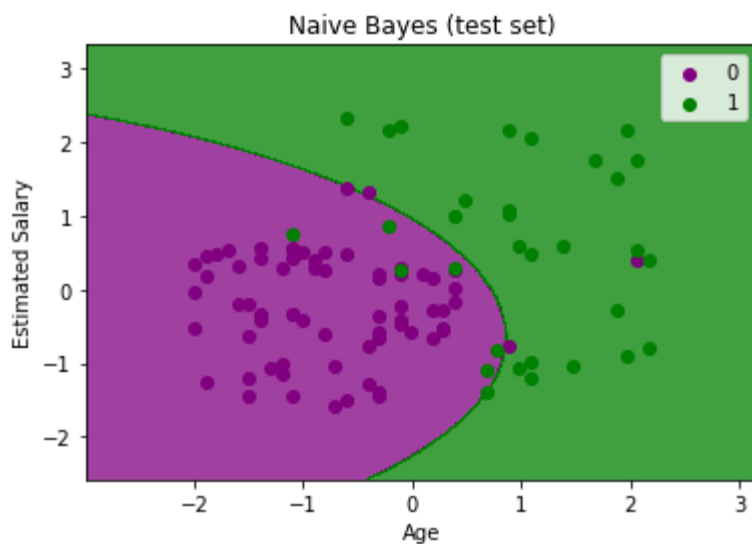
```

```

11. mtp.scatter(x_set[y_set == j, 0], x_set[y_set == j, 1],
12.             c = ListedColormap(('purple', 'green'))(i), label = j)
13. mtp.title('Naive Bayes (test set)')
14. mtp.xlabel('Age')
15. mtp.ylabel('Estimated Salary')
16. mtp.legend()
17. mtp.show()

```

Output:



The above output is final output for test set data. As we can see the classifier has created a Gaussian curve to divide the "purchased" and "not purchased" variables. There are some wrong predictions which we have calculated in Confusion matrix. But still it is pretty good classifier.

FAQ'S

1. What is meant by Regression?
2. What is meant by Naïve Bias classification?
3. What is classification?
4. What is dependence & independent variable?
5. What is difference between classification and regression?

Outcome:

With the completion of this assignment the students will be able to implement Naïve bias on given dataset.

Conclusion:

Assignment No:- 3

Title of Assignment:

Implement K Nearest Neighbor(KNN) classification using Scikit-learn (sklearn)

During assignment students will be able to:

- Learn Classification.
- Implementation of K Nearest Neighbor(KNN) classification.

Prerequisites: Knowledge of python programming and machine learning libraries.

Concepts to be used: Classification.

Input: Any dataset

Output: Successful implementation of K Nearest Neighbor(KNN) classification.

Relevant Theory / Literature Survey:

K Nearest Neighbor(KNN) is a very simple, easy to understand, versatile and one of the topmost machine learning algorithms. KNN used in the variety of applications such as finance, healthcare, political science, handwriting detection, image recognition and video recognition. In Credit ratings, financial institutes will predict the credit rating of customers. In loan disbursement, banking institutes will predict whether the loan is safe or risky. In political science, classifying potential voters in two classes will vote or won't vote. KNN algorithm used for both classification and regression problems. KNN algorithm based on feature similarity approach.

In this tutorial, you are going to cover the following topics:

- K-Nearest Neighbor Algorithm
- How does the KNN algorithm work?
- Eager Vs Lazy learners
- How do you decide the number of neighbors in KNN?
- Curse of Dimensionality
- Classifier Building in Scikit-learn

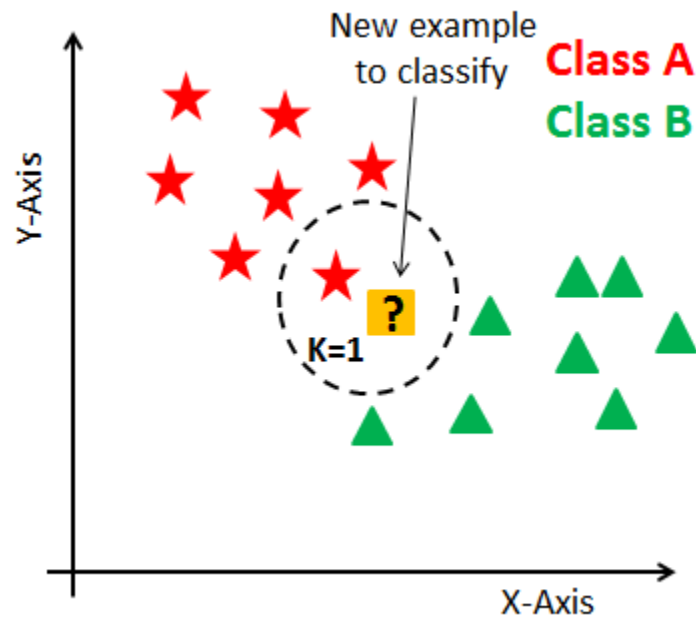
- Pros and Cons
- How to improve KNN performance?
- Conclusion

K-Nearest Neighbors

KNN is a non-parametric and lazy learning algorithm. Non-parametric means there is no assumption for underlying data distribution. In other words, the model structure determined from the dataset. This will be very helpful in practice where most of the real world datasets do not follow mathematical theoretical assumptions. Lazy algorithm means it does not need any training data points for model generation. All training data used in the testing phase. This makes training faster and testing phase slower and costlier. Costly testing phase means time and memory. In the worst case, KNN needs more time to scan all data points and scanning all data points will require more memory for storing training data.

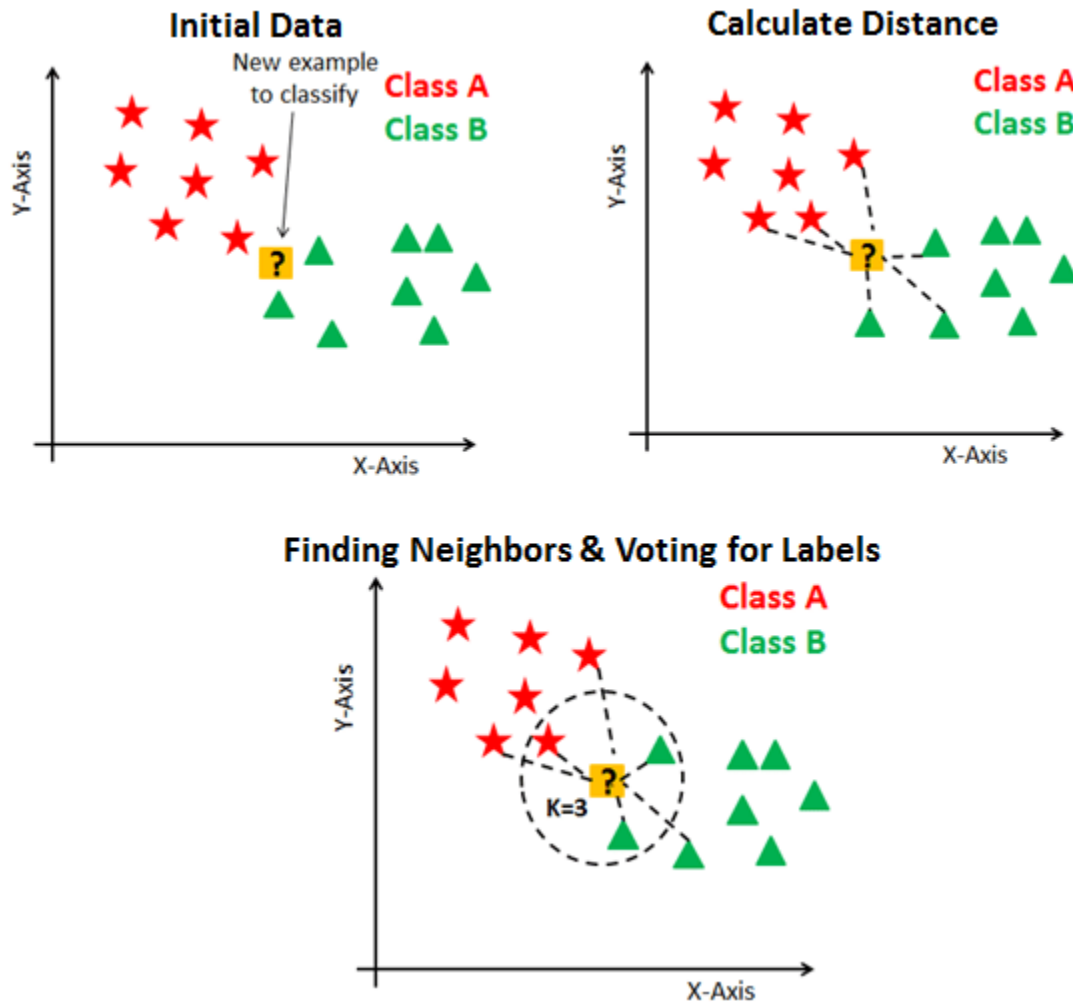
How does the KNN algorithm work?

In KNN, K is the number of nearest neighbors. The number of neighbors is the core deciding factor. K is generally an odd number if the number of classes is 2. When K=1, then the algorithm is known as the nearest neighbor algorithm. This is the simplest case. Suppose P1 is the point, for which label needs to predict. First, you find the one closest point to P1 and then the label of the nearest point assigned to P1.



Suppose P1 is the point, for which label needs to predict. First, you find the k closest point to P1 and then classify points by majority vote of its k neighbors. Each object votes for their class and the class with the most votes is taken as the prediction. For finding closest similar points, you find the distance between points using distance measures such as Euclidean distance, Hamming distance, Manhattan distance and Minkowski distance. KNN has the following basic steps:

1. Calculate distance
2. Find closest neighbors
3. Vote for labels



Eager Vs. Lazy Learners

Eager learners mean when given training points will construct a generalized model before performing prediction on given new points to classify. You can think of such learners as being ready, active and eager to classify unobserved data points.

Lazy Learning means there is no need for learning or training of the model and all of the data points used at the time of prediction. Lazy learners wait until the last minute before classifying any data point. Lazy learner stores merely the training dataset and waits until classification needs to perform. Only when it sees the test tuple does it perform generalization to classify the tuple based on its similarity to the stored training tuples. Unlike eager learning methods, lazy learners do less work in the training phase and more work in the testing phase to make a classification. Lazy learners are also known as instance-based learners because lazy learners store the training points or instances, and all learning is based on instances.

Curse of Dimensionality

KNN performs better with a lower number of features than a large number of features. You can say that when the number of features increases than it requires more data. Increase in dimension also leads to the problem of overfitting. To avoid overfitting, the needed data will need to grow exponentially as you increase the number of dimensions. This problem of higher dimension is known as the Curse of Dimensionality.

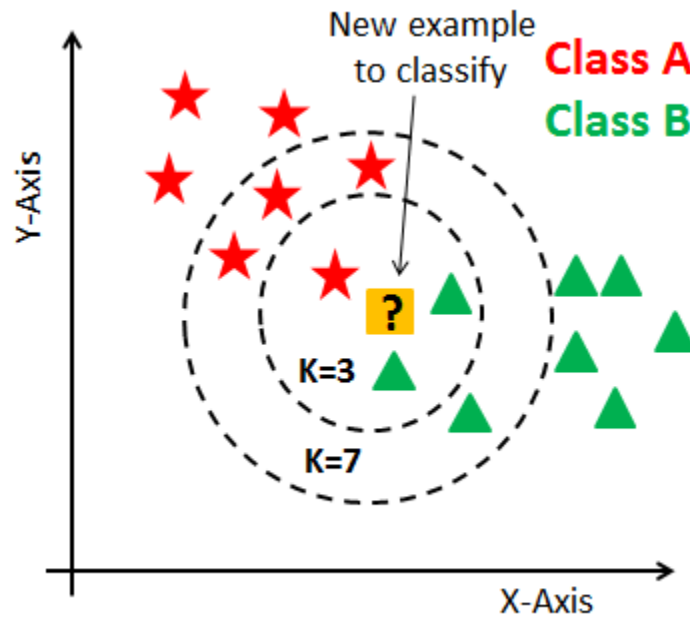
To deal with the problem of the curse of dimensionality, you need to perform principal component analysis before applying any machine learning algorithm, or you can also use feature selection approach. Research has shown that in large dimension Euclidean distance is not useful anymore. Therefore, you can prefer other measures such as cosine similarity, which get decidedly less affected by high dimension.

How do you decide the number of neighbors in KNN?

Now, you understand the KNN algorithm working mechanism. At this point, the question arises that How to choose the optimal number of neighbors? And what are its effects on the classifier? The number of neighbors(K) in KNN is a hyperparameter that you need choose at the time of model building. You can think of K as a controlling variable for the prediction model.

Research has shown that no optimal number of neighbors suits all kind of data sets. Each dataset has it's own requirements. In the case of a small number of neighbors, the noise will have a higher influence on the result, and a large number of neighbors make it computationally expensive. Research has also shown that a small amount of neighbors are most flexible fit which will have low bias but high variance and a large number of neighbors will have a smoother decision boundary which means lower variance but higher bias.

Generally, Data scientists choose as an odd number if the number of classes is even. You can also check by generating the model on different values of k and check their performance. You can also try Elbow method [here](#).



Classifier Building in Scikit-learn

KNN Classifier

Defining dataset

Let's first create your own dataset. Here you need two kinds of attributes or columns in your data: Feature and label. The reason for two type of column is "supervised nature of KNN algorithm".

```
# Assigning features and label variables
```

```
# First Feature
```

```
weather=['Sunny','Sunny','Overcast','Rainy','Rainy','Rainy','Overcast','Sun  
ny','Sunny',
```

```
'Rainy','Sunny','Overcast','Overcast','Rainy']
```

```
# Second Feature
```

```
temp=['Hot','Hot','Hot','Mild','Cool','Cool','Cool','Mild','Cool','Mild','M  
ild','Mild','Hot','Mild']
```

```
# Label or target variable
```

```
play=['No', 'No', 'Yes', 'Yes', 'Yes', 'No', 'Yes', 'No', 'Yes', 'Yes', 'Yes', 'Yes', 'Yes', 'Yes', 'No']
```

In this dataset, you have two features (weather and temperature) and one label(play).

Encoding data columns

Various machine learning algorithms require numerical input data, so you need to represent categorical columns in a numerical column.

In order to encode this data, you could map each value to a number. e.g. Overcast:0, Rainy:1, and Sunny:2.

This process is known as label encoding, and sklearn conveniently will do this for you using Label Encoder.

```
# Import LabelEncoder
from sklearn import preprocessing
#creating labelEncoder
le = preprocessing.LabelEncoder()
# Converting string labels into numbers.
weather_encoded=le.fit_transform(weather)
print(weather_encoded)
[2 2 0 1 1 1 0 2 2 1 2 0 0 1]
```

Here, you imported preprocessing module and created Label Encoder object. Using this LabelEncoder object, you can fit and transform "weather" column into the numeric column.

Similarly, you can encode temperature and label into numeric columns.

```
# converting string labels into numbers
temp_encoded=le.fit_transform(temp)
label=le.fit_transform(play)
```

Combining Features

Here, you will combine multiple columns or features into a single set of data using "zip" function

```
#combinig weather and temp into single listof tuples
```

```
features=list(zip(weather_encoded,temp_encoded))
```

Generating Model

Let's build KNN classifier model.

First, import the KNeighborsClassifier module and create KNN classifier object by passing argument number of neighbors in KNeighborsClassifier() function.

Then, fit your model on the train set using fit() and perform prediction on the test set using predict().

```
from sklearn.neighbors import KNeighborsClassifier
```

```
model = KNeighborsClassifier(n_neighbors=3)
```

```
# Train the model using the training sets
```

```
model.fit(features,label)
```

```
#Predict Output
```

```
predicted= model.predict([[0,2]]) # 0:Overcast, 2:Mild
```

```
print(predicted)
```

```
[1]
```

In the above example, you have given input [0,2], where 0 means Overcast weather and 2 means Mild temperature. Model predicts [1], which means play.

KNN with Multiple Labels

Till now, you have learned How to create KNN classifier for two in python using scikit-learn. Now you will learn about KNN with multiple classes.

In the model the building part, you can use the wine dataset, which is a very famous multi-class classification problem. This data is the result of a chemical analysis of wines grown in the same region in Italy using three

different cultivars. The analysis determined the quantities of 13 constituents found in each of the three types of wines.

The dataset comprises 13 features ('alcohol', 'malic_acid', 'ash', 'alcalinity_of_ash', 'magnesium', 'total_phenols', 'flavanoids', 'nonflavanoid_phenols', 'proanthocyanins', 'color_intensity', 'hue', 'od280/od315_of_diluted_wines', 'proline') and a target (type of cultivars).

This data has three types of cultivar classes: 'class_0', 'class_1', and 'class_2'. Here, you can build a model to classify the type of cultivar. The dataset is available in the scikit-learn library, or you can also download it from the UCI Machine Learning Library.

Loading Data

Let's first load the required wine dataset from scikit-learn datasets.

```
#Import scikit-learn dataset library  
  
from sklearn import datasets
```

```
#Load dataset  
  
wine = datasets.load_wine()
```

Exploring Data

After you have loaded the dataset, you might want to know a little bit more about it. You can check feature and target names.

```
# print the names of the features  
  
print(wine.feature_names)  
  
['alcohol', 'malic acid', 'ash', 'alcalinity of ash', 'magnesium',  
'total_phenols', 'flavanoids', 'nonflavanoid_phenols', 'proanthocyanins',  
'color_intensity', 'hue', 'od280/od315 of diluted wines', 'proline']  
  
# print the label species(class 0, class 1, class 2)  
  
print(wine.target_names)  
  
['class 0' 'class 1' 'class 2']
```

Let's check top 5 records of the feature set.

```
# print the wine data (top 5 records)

print(wine.data[0:5])

[[ 1.42300000e+01  1.71000000e+00  2.43000000e+00  1.56000000e+01
 1.27000000e+02  2.80000000e+00  3.06000000e+00  2.80000000e-01
 2.29000000e+00  5.64000000e+00  1.04000000e+00  3.92000000e+00
 1.06500000e+03]

[ 1.32000000e+01  1.78000000e+00  2.14000000e+00  1.12000000e+01
 1.00000000e+02  2.65000000e+00  2.76000000e+00  2.60000000e-01
 1.28000000e+00  4.38000000e+00  1.05000000e+00  3.40000000e+00
 1.05000000e+03]

[ 1.31600000e+01  2.36000000e+00  2.67000000e+00  1.86000000e+01
 1.01000000e+02  2.80000000e+00  3.24000000e+00  3.00000000e-01
 2.81000000e+00  5.68000000e+00  1.03000000e+00  3.17000000e+00
 1.18500000e+03]

[ 1.43700000e+01  1.95000000e+00  2.50000000e+00  1.68000000e+01
 1.13000000e+02  3.85000000e+00  3.49000000e+00  2.40000000e-01
 2.18000000e+00  7.80000000e+00  8.60000000e-01  3.45000000e+00
 1.48000000e+03]

[ 1.32400000e+01  2.59000000e+00  2.87000000e+00  2.10000000e+01
 1.18000000e+02  2.80000000e+00  2.69000000e+00  3.90000000e-01
 1.82000000e+00  4.32000000e+00  1.04000000e+00  2.93000000e+00
 7.35000000e+02]]
```

Let's check records of the target set.

[illegible]

[illegible]

Let's explore it for a bit more. You can also check the shape of the dataset using `shape`.

```
# print data(feature).shape

print(wine.data.shape)

(178, 13)

# print target(or label).shape

print(wine.target.shape)

(178,)
```

Splitting Data

To understand model performance, dividing the dataset into a training set and a test set is a good strategy.

Let's split dataset by using function `train_test_split()`. You need to pass 3 parameters features, target, and test_set size. Additionally, you can use `random_state` to select records randomly.

```
# Import train_test_split function

from sklearn.model_selection import train_test_split

# Split dataset into training set and test set

X_train, X_test, y_train, y_test = train_test_split(wine.data, wine.target,
test size=0.3) # 70% training and 30% test
```

Generating Model for K=5

Let's build KNN classifier model for $k=5$.

```
#Import knearest neighbors Classifier model
```

```
from sklearn.neighbors import KNeighborsClassifier
```

```
#Create KNN Classifier

knn = KNeighborsClassifier(n_neighbors=5)

#Train the model using the training sets

knn.fit(X_train, y_train)

#Predict the response for test dataset

y_pred = knn.predict(X_test)
```

Model Evaluation for k=5

Let's estimate, how accurately the classifier or model can predict the type of cultivars.

Accuracy can be computed by comparing actual test set values and predicted values.

```
#Import scikit-learn metrics module for accuracy calculation

from sklearn import metrics

# Model Accuracy, how often is the classifier correct?

print("Accuracy:",metrics.accuracy_score(y_test, y_pred))

Accuracy: 0.685185185185
```

Well, you got a classification rate of 68.51%, considered as good accuracy.

For further evaluation, you can also create a model for a different number of neighbors.

Re-generating Model for K=7

Let's build KNN classifier model for k=7.

```
#Import knearest neighbors Classifier model

from sklearn.neighbors import KNeighborsClassifier

#Create KNN Classifier

knn = KNeighborsClassifier(n_neighbors=7)
```



```
#Train the model using the training sets
```

```
knn.fit(X_train, y_train)
```

```
#Predict the response for test dataset
```

```
y_pred = knn.predict(X_test)
```

Model Evaluation for k=7

Let's again estimate, how accurately the classifier or model can predict the type of cultivars for k=7.

```
#Import scikit-learn metrics module for accuracy calculation
```

```
from sklearn import metrics
```

```
# Model Accuracy, how often is the classifier correct?
```

```
print("Accuracy:",metrics.accuracy_score(y_test, y_pred))
```

```
Accuracy: 0.777777777778
```

Well, you got a classification rate of 77.77%, considered as good accuracy.

Here, you have increased the number of neighbors in the model and accuracy got increased. But, this is not necessary for each case that an increase in many neighbors increases the accuracy. For a more detailed understanding of it, you can refer section "How to decide the number of neighbors?" of this tutorial.

Pros

The training phase of K-nearest neighbor classification is much faster compared to other classification algorithms. There is no need to train a model for generalization, That is why KNN is known as the simple and instance-based learning algorithm. KNN can be useful in case of nonlinear data. It can be used with the regression problem. Output value for the object is computed by the average of k closest neighbors value.

Cons

The testing phase of K-nearest neighbor classification is slower and costlier in terms of time and memory. It requires large memory for storing the entire training dataset for prediction. KNN requires scaling of data because KNN uses the Euclidean distance between two data points to find nearest neighbors. Euclidean distance is

sensitive to magnitudes. The features with high magnitudes will weight more than features with low magnitudes. KNN also not suitable for large dimensional data.

How to improve KNN?

For better results, normalizing data on the same scale is highly recommended. Generally, the normalization range considered between 0 and 1. KNN is not suitable for the large dimensional data. In such cases, dimension needs to reduce to improve the performance. Also, handling missing values will help us in improving results.

FAQ'S

1. What is meant by classification?
2. What is KNN?
3. What is the pros and cons of KNN?
4. What is dependence & independent variable?
5. What is difference between classification and regression?

Outcome:

With the completion of this assignment the students will be able to implement KNN on given dataset.

Conclusion:

Assignment No:- 4

Title of Assignment:

Given the following data, which specify classifications for nine combinations of VAR1 and VAR2 predict a classification for a case where VAR1=0.906 and VAR2=0.606, using the result of k-means clustering with 3 means (i.e., centroids)

VAR1	VAR2	CLASS
1.713	1.586	0
0.180	1.786	1
0.353	1.240	1
0.940	1.566	0
1.486	0.759	1
1.266	1.106	0
1.540	0.419	1
0.459	1.799	1
0.773	0.186	1

During assignment students will be able to:

- Learn Unsupervised Learning & Clustering.
- Implementation of k-means clustering.

Prerequisites: Knowledge of python programming and machine learning libraries.

Concepts to be used: Unsupervised Learning & Clustering.

Input: given dataset

Output: Successful implementation of k-means clustering.

Relevant Theory / Literature Survey:

K-Means Clustering Algorithm

K-Means Clustering is an unsupervised learning algorithm that is used to solve the clustering problems in machine learning or data science. In this topic, we will learn what is K-means clustering algorithm, how the algorithm works, along with the Python implementation of k-means clustering.

What is K-Means Algorithm?

K-Means Clustering is an [Unsupervised Learning algorithm](#)

, which groups the unlabeled dataset into different clusters. Here K defines the number of pre-defined clusters that need to be created in the process, as if $K=2$, there will be two clusters, and for $K=3$, there will be three clusters, and so on.

It is an iterative algorithm that divides the unlabeled dataset into k different clusters in such a way that each dataset belongs only one group that has similar properties.

It allows us to cluster the data into different groups and a convenient way to discover the categories of groups in the unlabeled dataset on its own without the need for any training.

It is a centroid-based algorithm, where each cluster is associated with a centroid. The main aim of this algorithm is to minimize the sum of distances between the data point and their corresponding clusters.

The algorithm takes the unlabeled dataset as input, divides the dataset into k -number of clusters, and repeats the process until it does not find the best clusters. The value of k should be predetermined in this algorithm.

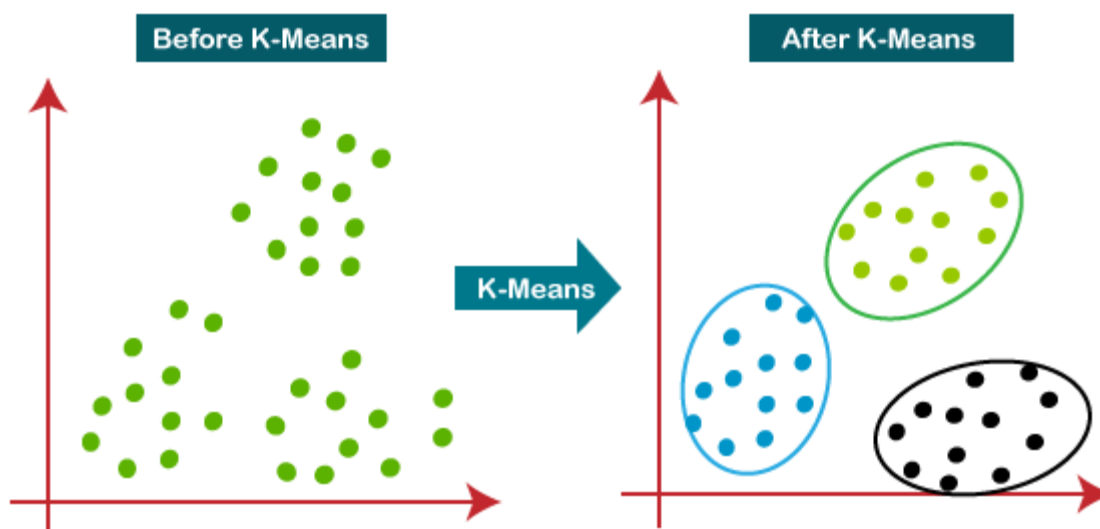
The k -means [clustering](#)

algorithm mainly performs two tasks:

- Determines the best value for K center points or centroids by an iterative process.
- Assigns each data point to its closest k -center. Those data points which are near to the particular k -center, create a cluster.

Hence each cluster has datapoints with some commonalities, and it is away from other clusters.

The below diagram explains the working of the K -means Clustering Algorithm:



How does the K -Means Algorithm Work?

The working of the K-Means algorithm is explained in the below steps:

Step-1: Select the number K to decide the number of clusters.

Step-2: Select random K points or centroids. (It can be other from the input dataset).

Step-3: Assign each data point to their closest centroid, which will form the predefined K clusters.

Step-4: Calculate the variance and place a new centroid of each cluster.

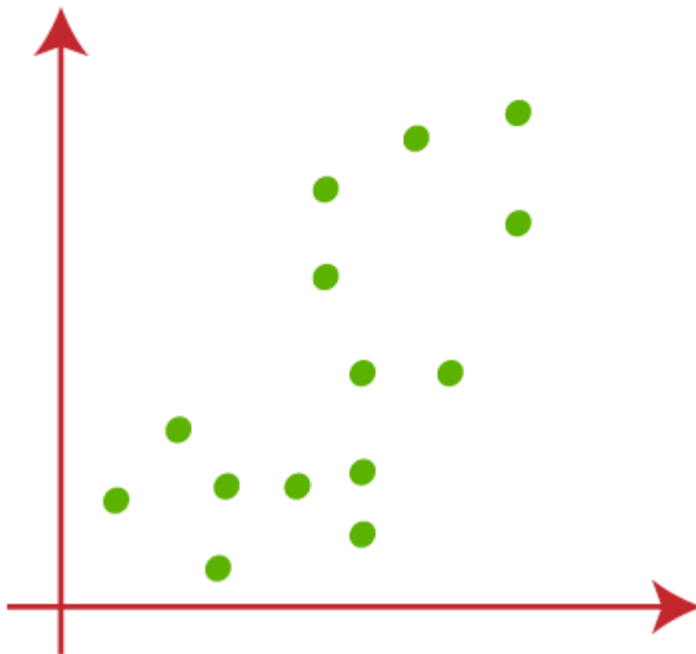
Step-5: Repeat the third steps, which means reassign each datapoint to the new closest centroid of each cluster.

Step-6: If any reassignment occurs, then go to step-4 else go to FINISH.

Step-7: The model is ready.

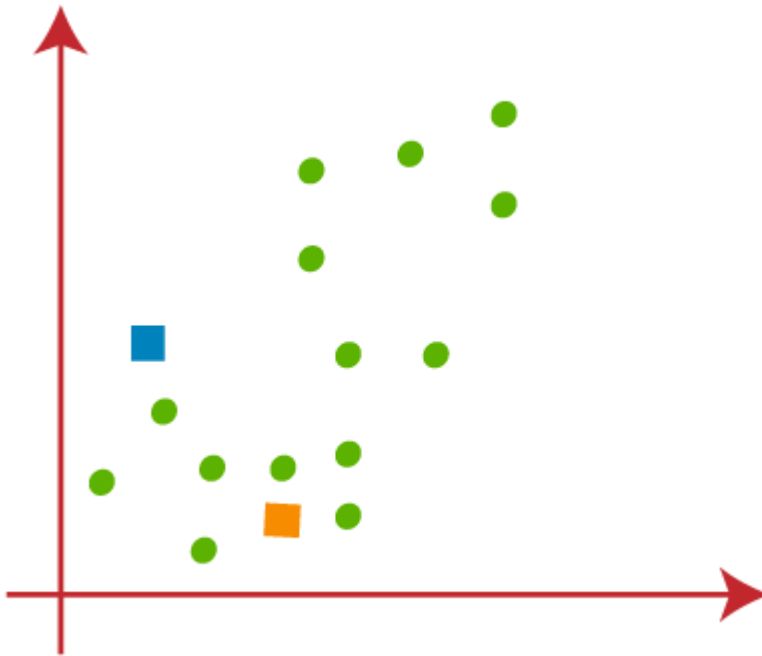
Let's understand the above steps by considering the visual plots:

Suppose we have two variables M1 and M2. The x-y axis scatter plot of these two variables is given below:



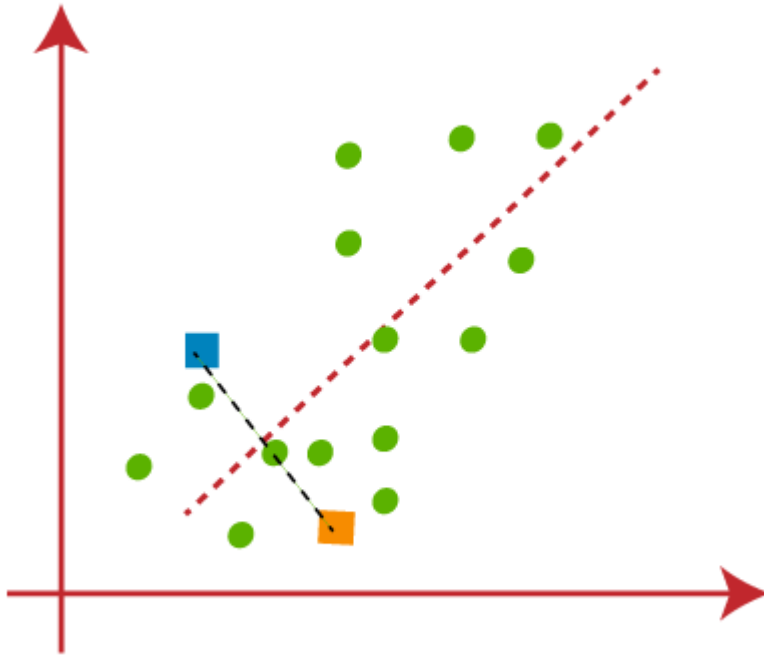
- Let's take number k of clusters, i.e., $K=2$, to identify the dataset and to put them into different clusters. It means here we will try to group these datasets into two different clusters.

- We need to choose some random k points or centroid to form the cluster. These points can be either the points from the dataset or any other point. So, here we are selecting the below two points as k points, which are not the part of our dataset. Consider the below image:

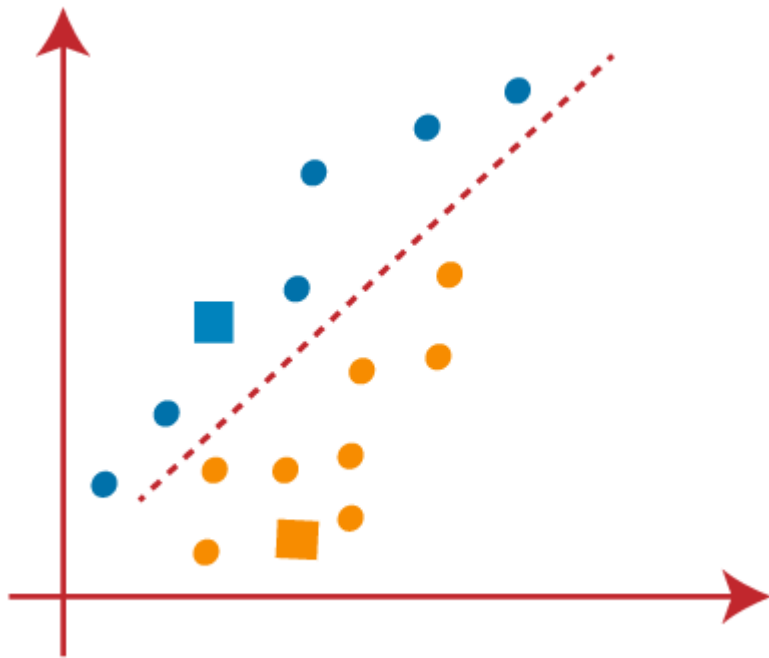


- Now we will assign each data point of the scatter plot to its closest K-point or centroid. We will compute it by applying some mathematics that we have studied to calculate the distance between two points. So, we will draw a median between both the centroids. Consider the below

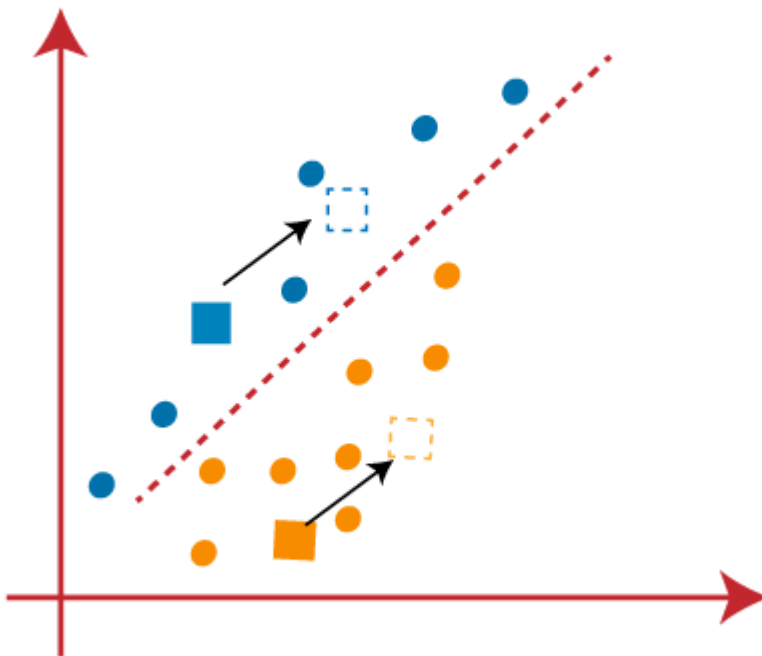
image:



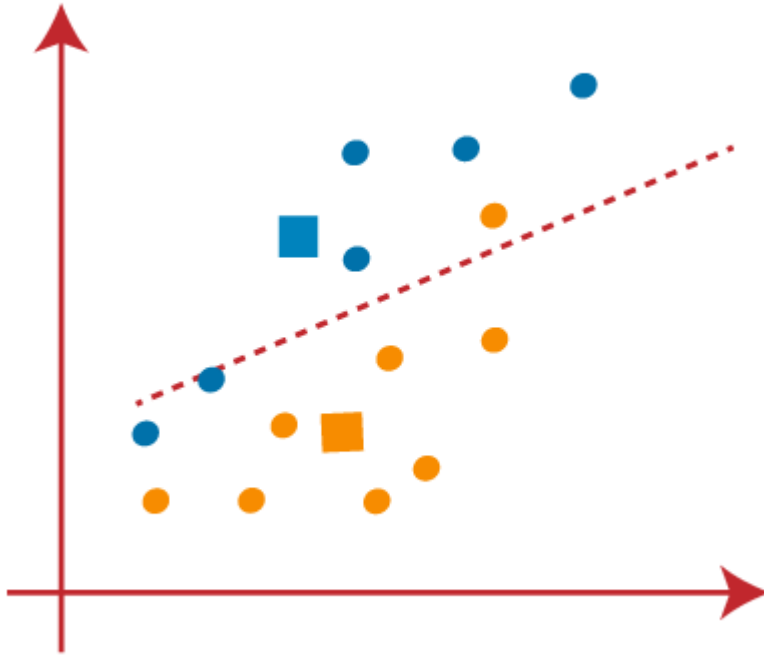
From the above image, it is clear that points left side of the line is near to the K1 or blue centroid, and points to the right of the line are close to the yellow centroid. Let's color them as blue and yellow for clear visualization.



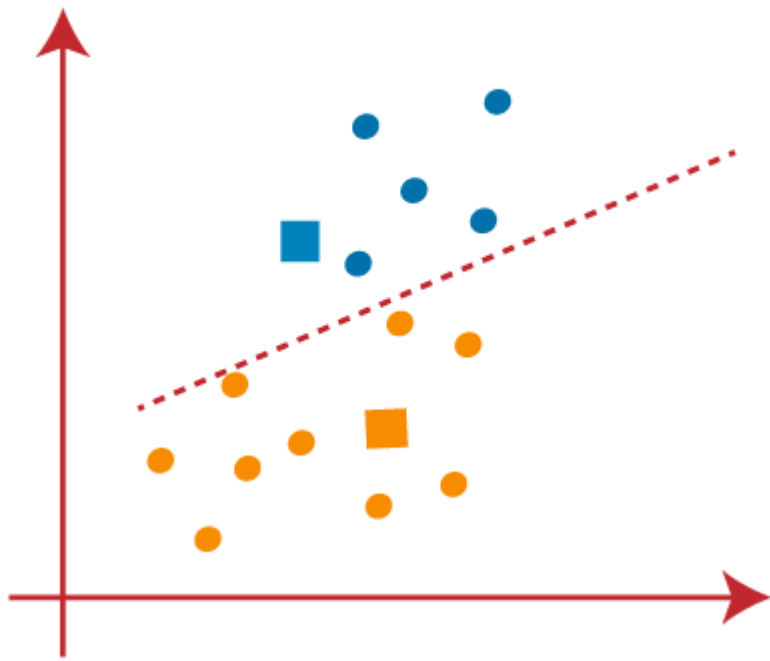
- As we need to find the closest cluster, so we will repeat the process by choosing a **new centroid**. To choose the new centroids, we will compute the center of gravity of these centroids, and will find new centroids as below:



- Next, we will reassign each datapoint to the new centroid. For this, we will repeat the same process of finding a median line. The median will be like below image:

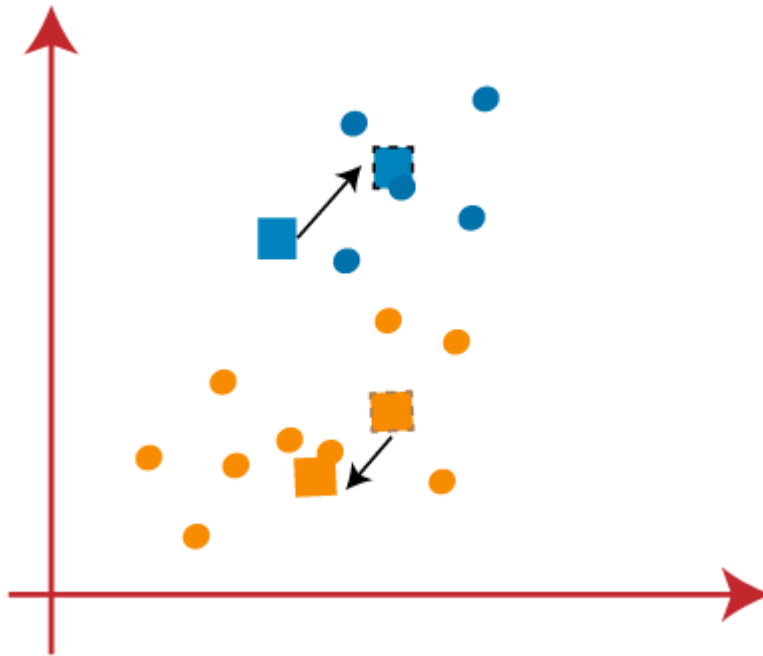


From the above image, we can see, one yellow point is on the left side of the line, and two blue points are right to the line. So, these three points will be assigned to new centroids.

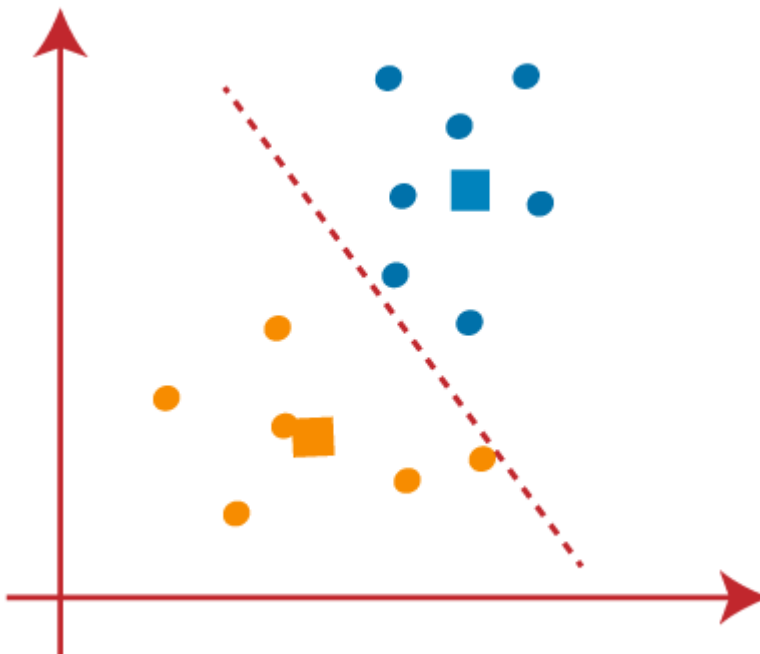


As reassignment has taken place, so we will again go to the step-4, which is finding new centroids or K-points.

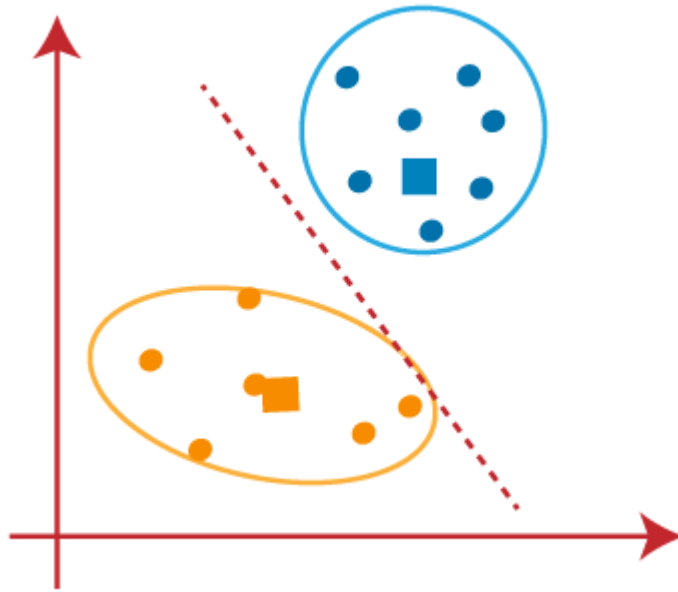
- We will repeat the process by finding the center of gravity of centroids, so the new centroids will be as shown in the below image:



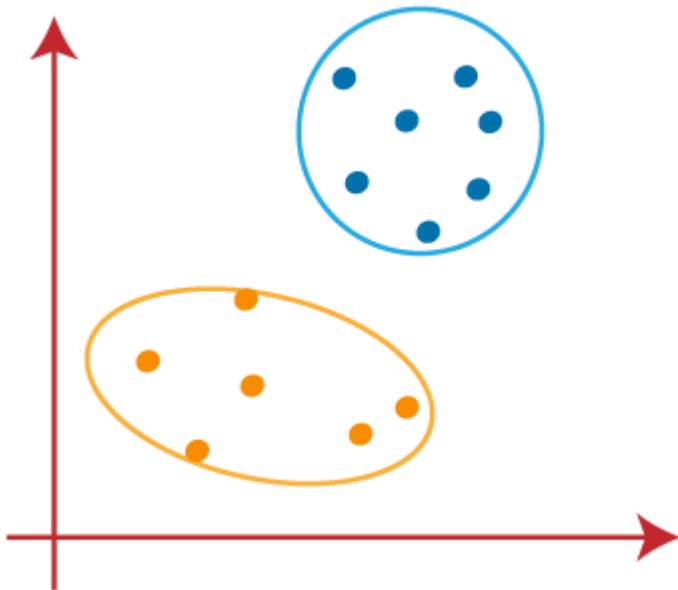
- As we got the new centroids so again will draw the median line and reassign the data points. So, the image will be:



- We can see in the above image; there are no dissimilar data points on either side of the line, which means our model is formed. Consider the below image:



As our model is ready, so we can now remove the assumed centroids, and the two final clusters will be as shown in the below image:



How to choose the value of "K number of clusters" in K-means Clustering?

The performance of the K-means clustering algorithm depends upon highly efficient clusters that it forms. But choosing the optimal number of clusters is a big task. There are some different ways to find the optimal number of clusters, but here we are discussing the most appropriate method to find the number of clusters or value of K. The method is given below:

Elbow Method

The Elbow method is one of the most popular ways to find the optimal number of clusters. This method uses the concept of WCSS value. **WCSS** stands for **Within Cluster Sum of Squares**, which defines the total variations within a cluster. The formula to calculate the value of WCSS (for 3 clusters) is given below:

$$WCSS = \sum_{P_i \text{ in Cluster1}} \text{distance}(P_i C_1)^2 + \sum_{P_i \text{ in Cluster2}} \text{distance}(P_i C_2)^2 + \sum_{P_i \text{ in Cluster3}} \text{distance}(P_i C_3)^2$$

In the above formula of WCSS,

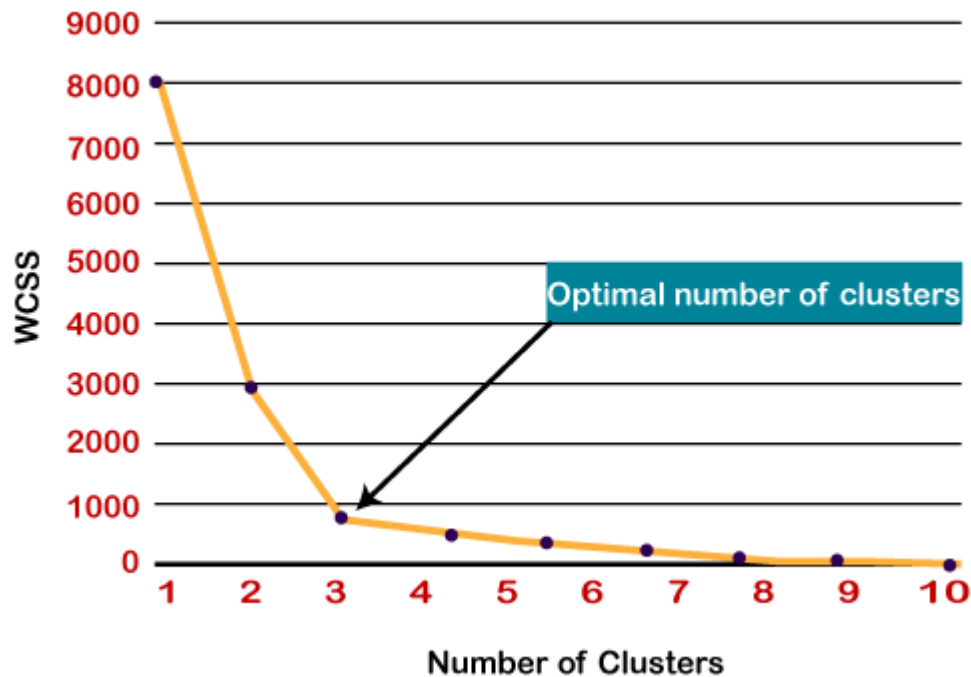
$\sum_{P_i \text{ in Cluster1}} \text{distance}(P_i C_1)^2$: It is the sum of the square of the distances between each data point and its centroid within a cluster1 and the same for the other two terms.

To measure the distance between data points and centroid, we can use any method such as Euclidean distance or Manhattan distance.

To find the optimal value of clusters, the elbow method follows the below steps:

- It executes the K-means clustering on a given dataset for different K values (ranges from 1-10).
- For each value of K, calculates the WCSS value.
- Plots a curve between calculated WCSS values and the number of clusters K.
- The sharp point of bend or a point of the plot looks like an arm, then that point is considered as the best value of K.

Since the graph shows the sharp bend, which looks like an elbow, hence it is known as the elbow method. The graph for the elbow method looks like the below image:



Note: We can choose the number of clusters equal to the given data points. If we choose the number of clusters equal to the data points, then the value of WCSS becomes zero, and that will be the endpoint of the plot.

Python Implementation of K-means Clustering Algorithm

In the above section, we have discussed the K-means algorithm, now let's see how it can be implemented using [Python](#)

.

Before implementation, let's understand what type of problem we will solve here. So, we have a dataset of **Mall_Customers**, which is the data of customers who visit the mall and spend there.

In the given dataset, we have **Customer_Id**, **Gender**, **Age**, **Annual Income (\$)**, and **Spending Score** (which is the calculated value of how much a customer has spent in the mall, the more the value, the more he has spent). From this dataset, we need to calculate some patterns, as it is an unsupervised method, so we don't know what to calculate exactly.

The steps to be followed for the implementation are given below:

- **Data Pre-processing**
- **Finding the optimal number of clusters using the elbow method**
- **Training the K-means algorithm on the training dataset**
- **Visualizing the clusters**

Step-1: Data pre-processing Step

The first step will be the data pre-processing, as we did in our earlier topics of Regression and Classification. But for the clustering problem, it will be different from other models. Let's discuss it:

- **Importing Libraries**

As we did in previous topics, firstly, we will import the libraries for our model, which is part of data pre-processing. The code is given below:

1. # importing libraries
2. **import** numpy as nm
3. **import** matplotlib.pyplot as mtp
4. **import** pandas as pd

In the above code, the **numpy** we have imported for the performing mathematics calculation, **matplotlib** is for plotting the graph, and **pandas** are for managing the dataset.

- **Importing the Dataset:**

Next, we will import the dataset that we need to use. So here, we are using the Mall_Customer_data.csv dataset. It can be imported using the below code:

1. # Importing the dataset
2. dataset = pd.read_csv('Mall_Customers_data.csv')

By executing the above lines of code, we will get our dataset in the Spyder IDE. The dataset looks like the below image:

Index	CustomerID	Genre	Age	Annual Income (k\$)	Spending Score (1-100)
0	1	Male	19	15	39
1	2	Male	21	15	81
2	3	Female	20	16	6
3	4	Female	23	16	77
4	5	Female	31	17	40
5	6	Female	22	17	76
6	7	Female	35	18	6
7	8	Female	23	18	94
8	9	Male	64	19	3
9	10	Female	30	19	72
10	11	Male	67	19	14
11	12	Female	35	19	99
12	13	Female	58	20	15
13	14	Female	24	20	77
14	15	Male	37	20	13
15	16	Male	22	20	79

From the above dataset, we need to find some patterns in it.

- **Extracting Independent Variables**

Here we don't need any dependent variable for data pre-processing step as it is a clustering problem, and we have no idea about what to determine. So we will just add a line of code for the matrix of features.

1. `x = dataset.iloc[:, [3, 4]].values`

As we can see, we are extracting only 3rd and 4th feature. It is because we need a 2d plot to visualize the model, and some features are not required, such as customer_id.

Step-2: Finding the optimal number of clusters using the elbow method

In the second step, we will try to find the optimal number of clusters for our clustering problem. So, as discussed above, here we are going to use the elbow method for this purpose.

As we know, the elbow method uses the WCSS concept to draw the plot by plotting WCSS values on the Y-axis and the number of clusters on the X-axis. So we are going to calculate the value for WCSS for different k values ranging from 1 to 10. Below is the code for it:

```
1. #finding optimal number of clusters using the elbow method
2. from sklearn.cluster import KMeans
3. wcss_list= [] #Initializing the list for the values of WCSS
4.
5. #Using for loop for iterations from 1 to 10.
6. for i in range(1, 11):
7.     kmeans = KMeans(n_clusters=i, init='k-means++', random_state= 42)
8.     kmeans.fit(x)
9.     wcss_list.append(kmeans.inertia_)
10. mtp.plot(range(1, 11), wcss_list)
11. mtp.title('The Elbow Method Graph')
12. mtp.xlabel('Number of clusters(k)')
13. mtp.ylabel('wcss_list')
14. mtp.show()
```

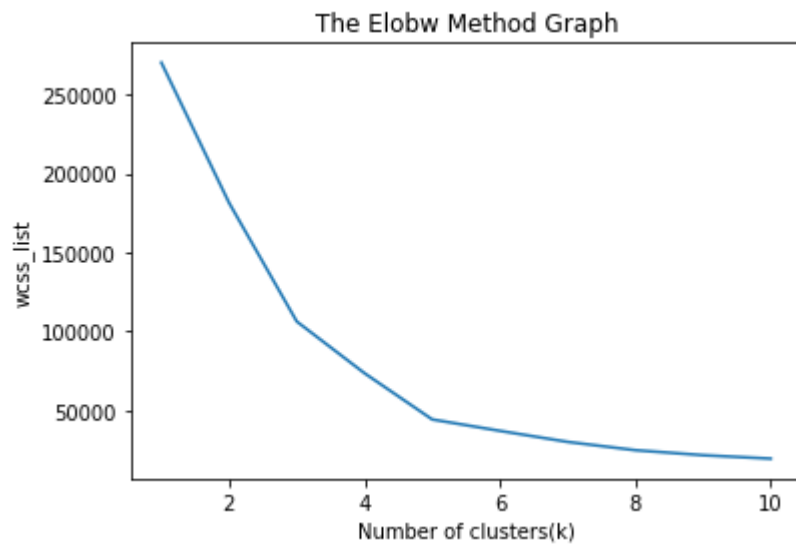
As we can see in the above code, we have used **the KMeans** class of sklearn. cluster library to form the clusters.

Next, we have created the **wcss_list** variable to initialize an empty list, which is used to contain the value of wcss computed for different values of k ranging from 1 to 10.

After that, we have initialized the for loop for the iteration on a different value of k ranging from 1 to 10; since for loop in Python, exclude the outbound limit, so it is taken as 11 to include 10th value.

The rest part of the code is similar as we did in earlier topics, as we have fitted the model on a matrix of features and then plotted the graph between the number of clusters and WCSS.

Output: After executing the above code, we will get the below output:



From the above plot, we can see the elbow point is at **5**. So the number of clusters here will be **5**.

wcss_list - List (10 elements)			
Index	Type	Size	Value
0	float64	1	269981.28
1	float64	1	181363.59595959596
2	float64	1	106348.37306211118
3	float64	1	73679.78903948834
4	float64	1	44448.45544793371
5	float64	1	37233.81451071001
6	float64	1	30259.65720728547
7	float64	1	25011.83934915659
8	float64	1	21850.165282585633
9	float64	1	19672.07284901432

Save and Close
Close

Step- 3: Training the K-means algorithm on the training dataset

As we have got the number of clusters, so we can now train the model on the dataset.

To train the model, we will use the same two lines of code as we have used in the above section, but here instead of using `i`, we will use `5`, as we know there are 5 clusters that need to be formed. The code is given below:

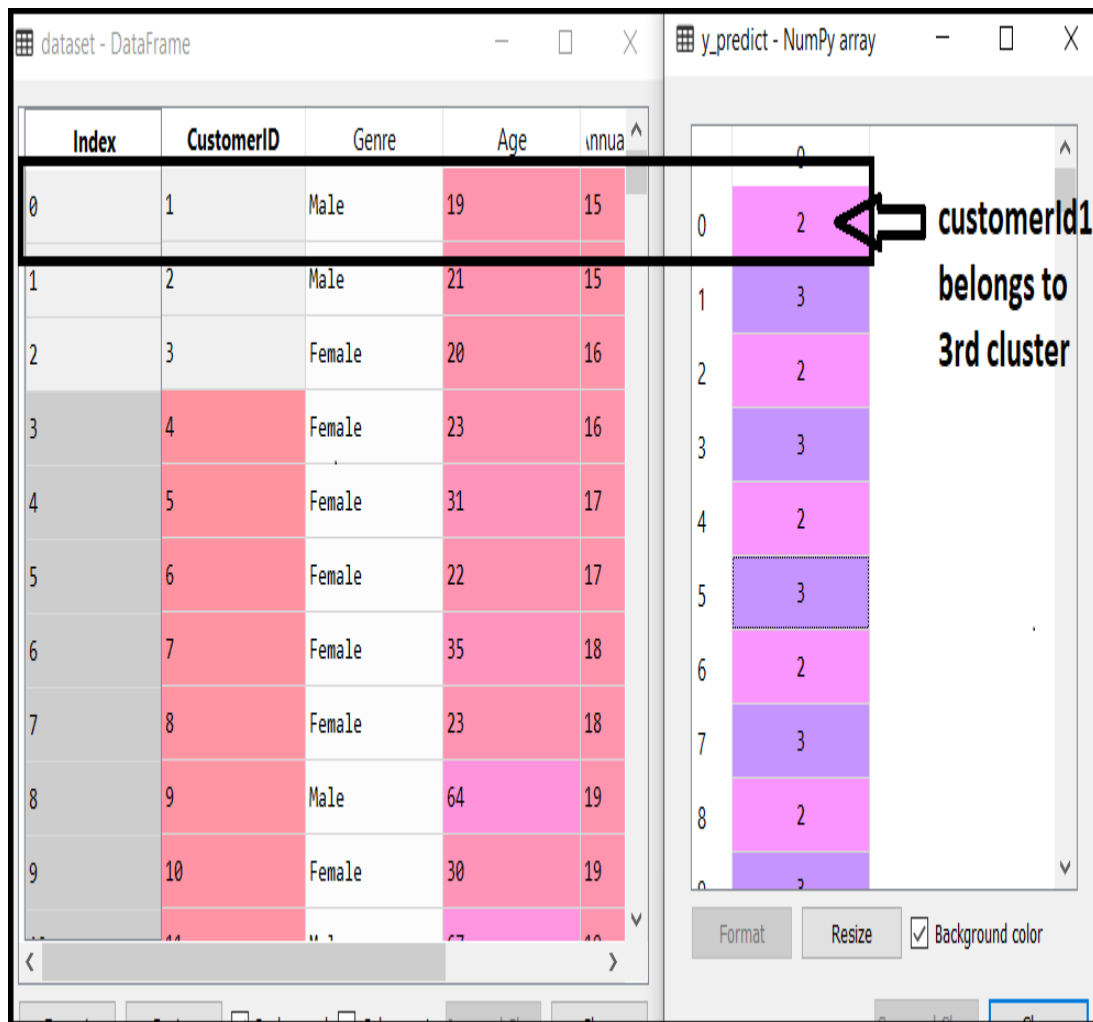
1. #training the K-means model on a dataset
2. `kmeans = KMeans(n_clusters=5, init='k-means++', random_state= 42)`

3. `y_predict= kmeans.fit_predict(x)`

The first line is the same as above for creating the object of KMeans class.

In the second line of code, we have created the dependent variable **y_predict** to train the model.

By executing the above lines of code, we will get the `y_predict` variable. We can check it under **the variable explorer** option in the Spyder IDE. We can now compare the values of `y_predict` with our original dataset. Consider the below image:



From the above image, we can now relate that the CustomerID 1 belongs to a cluster

3(as index starts from 0, hence 2 will be considered as 3), and 2 belongs to cluster 4, and so on.

Step-4: Visualizing the Clusters

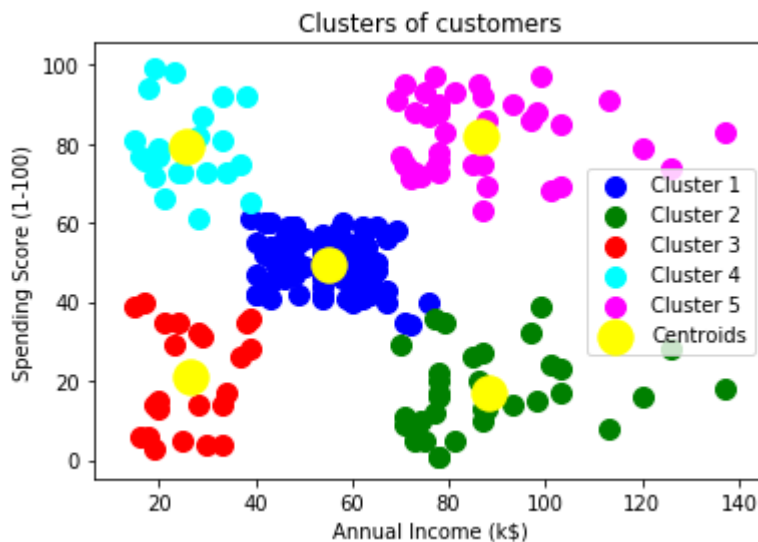
The last step is to visualize the clusters. As we have 5 clusters for our model, so we will visualize each cluster one by one.

To visualize the clusters will use scatter plot using `mtp.scatter()` function of `matplotlib`.

1. `#visualizing the clusters`
2. `mtp.scatter(x[y_predict == 0, 0], x[y_predict == 0, 1], s = 100, c = 'blue', label = 'Cluster 1') #for first cluster`
3. `mtp.scatter(x[y_predict == 1, 0], x[y_predict == 1, 1], s = 100, c = 'green', label = 'Cluster 2') #for second cluster`
4. `mtp.scatter(x[y_predict == 2, 0], x[y_predict == 2, 1], s = 100, c = 'red', label = 'Cluster 3') #for third cluster`
5. `mtp.scatter(x[y_predict == 3, 0], x[y_predict == 3, 1], s = 100, c = 'cyan', label = 'Cluster 4') #for fourth cluster`
6. `mtp.scatter(x[y_predict == 4, 0], x[y_predict == 4, 1], s = 100, c = 'magenta', label = 'Cluster 5') #for fifth cluster`
7. `mtp.scatter(kmeans.cluster_centers[:, 0], kmeans.cluster_centers[:, 1], s = 300, c = 'yellow', label = 'Centroid')`
8. `mtp.title('Clusters of customers')`
9. `mtp.xlabel('Annual Income (k$)')`
10. `mtp.ylabel('Spending Score (1-100)')`
11. `mtp.legend()`
12. `mtp.show()`

In above lines of code, we have written code for each clusters, ranging from 1 to 5. The first coordinate of the `mtp.scatter`, i.e., `x[y_predict == 0, 0]` containing the x value for the showing the matrix of features values, and the `y_predict` is ranging from 0 to 1.

Output:



The output image is clearly showing the five different clusters with different colors. The clusters are formed between two parameters of the dataset; Annual income of customer and Spending. We can change the colors and labels as per the requirement or choice. We can also observe some points from the above patterns, which are given below:

- **Cluster1** shows the customers with average salary and average spending so we can categorize these customers as
- Cluster2 shows the customer has a high income but low spending, so we can categorize them as **careful**.
- Cluster3 shows the low income and also low spending so they can be categorized as sensible.
- Cluster4 shows the customers with low income with very high spending so they can be categorized as **careless**.
- Cluster5 shows the customers with high income and high spending so they can be categorized as target, and these customers can be the most profitable customers for the mall owner.

FAQ'S

1. What is meant by Unsupervised learning?
2. What is meant by clustering?
3. What is k-means clustering?

Outcome:

With the completion of this assignment the students will be able to implement k-means clustering on given dataset.

Conclusion:

Assignment No:- 5

Title of Assignment:

Implement Linear Regression using Scikit-learn (sklearn)

During assignment students will be able to:

- Learn Supervised Learning and regression.
- Implementation of Linear Regression.

Prerequisites: Knowledge of python programming and machine learning libraries.

Concepts to be used: Regression.

Input: Any dataset

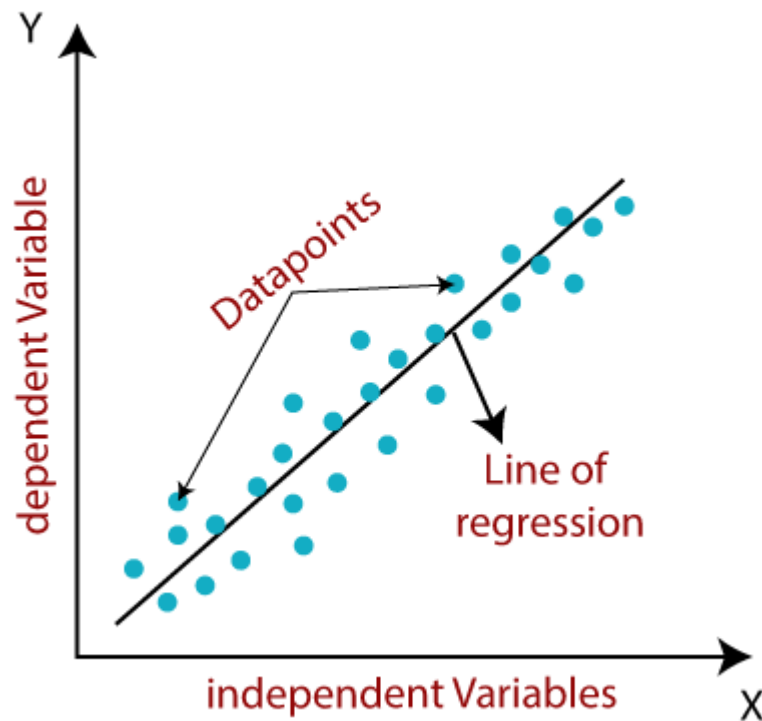
Output: Successful implementation of Linear Regression.

Relevant Theory / Literature Survey:

Linear regression is one of the easiest and most popular Machine Learning algorithms. It is a statistical method that is used for predictive analysis. Linear regression makes predictions for continuous/real or numeric variables such as **sales, salary, age, product price**, etc.

Linear regression algorithm shows a linear relationship between a dependent (y) and one or more independent (x) variables, hence called as linear regression. Since linear regression shows the linear relationship, which means it finds how the value of the dependent variable is changing according to the value of the independent variable.

The linear regression model provides a sloped straight line representing the relationship between the variables. Consider the below image:



Mathematically, we can represent a linear regression as:

$$y = a_0 + a_1x + \epsilon$$

Here,

Y= Dependent Variable (Target Variable)

X= Independent Variable (predictor Variable)

a_0 = intercept of the line (Gives an additional degree of freedom)

a_1 = Linear regression coefficient (scale factor to each input value).

ϵ = random error

The values for x and y variables are training datasets for Linear Regression model representation.

Types of Linear Regression

Linear regression can be further divided into two types of the algorithm:

- **Simple Linear Regression:**

If a single independent variable is used to predict the value of a numerical dependent variable, then such a Linear Regression algorithm is called Simple Linear Regression.

- **Multiple Linear regression:**

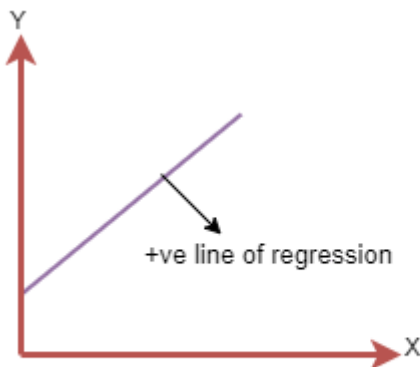
If more than one independent variable is used to predict the value of a numerical dependent variable, then such a Linear Regression algorithm is called Multiple Linear Regression.

Linear Regression Line

A linear line showing the relationship between the dependent and independent variables is called a **regression line**. A regression line can show two types of relationship:

- **Positive Linear Relationship:**

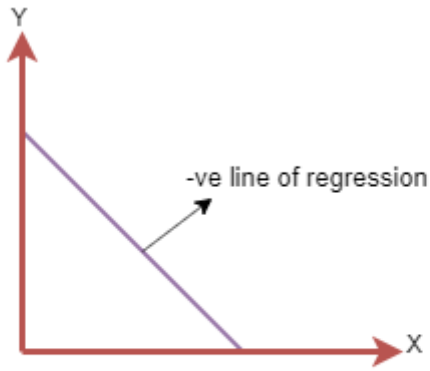
If the dependent variable increases on the Y-axis and independent variable increases on X-axis, then such a relationship is termed as a Positive linear relationship.



The line equation will be: $Y = a_0 + a_1X$

- **Negative Linear Relationship:**

If the dependent variable decreases on the Y-axis and independent variable increases on the X-axis, then such a relationship is called a negative linear relationship.



The line of equation will be: $Y = -a_0 + a_1X$

Finding the best fit line:

When working with linear regression, our main goal is to find the best fit line that means the error between predicted values and actual values should be minimized. The best fit line will have the least error.

The different values for weights or the coefficient of lines (a_0 , a_1) gives a different line of regression, so we need to calculate the best values for a_0 and a_1 to find the best fit line, so to calculate this we use cost function.

Cost function-

- The different values for weights or coefficient of lines (a_0 , a_1) gives the different line of regression, and the cost function is used to estimate the values of the coefficient for the best fit line.
- Cost function optimizes the regression coefficients or weights. It measures how a linear regression model is performing.
- We can use the cost function to find the accuracy of the **mapping function**, which maps the input variable to the output variable. This mapping function is also known as **Hypothesis function**.

For Linear Regression, we use the **Mean Squared Error (MSE)** cost function, which is the average of squared error occurred between the predicted values and actual values. It can be written as:

For the above linear equation, MSE can be calculated as:

$$MSE = \frac{1}{N} \sum_{i=1}^n (y_i - (a_1x_i + a_0))^2$$

Where,

N=Total number of observation

Y_i = Actual value

$(a_1x_i + a_0)$ = Predicted value.

Residuals: The distance between the actual value and predicted values is called residual. If the observed points are far from the regression line, then the residual will be high, and so cost function will high. If the scatter points are close to the regression line, then the residual will be small and hence the cost function.

Gradient Descent:

- Gradient descent is used to minimize the MSE by calculating the gradient of the cost function.
- A regression model uses gradient descent to update the coefficients of the line by reducing the cost function.
- It is done by a random selection of values of coefficient and then iteratively update the values to reach the minimum cost function.

Model Performance:

The Goodness of fit determines how the line of regression fits the set of observations. The process of finding the best model out of various models is called **optimization**. It can be achieved by below method:

1. R-squared method:

- R-squared is a statistical method that determines the goodness of fit.
- It measures the strength of the relationship between the dependent and independent variables on a scale of 0-100%.
- The high value of R-square determines the less difference between the predicted values and actual values and hence represents a good model.
- It is also called a **coefficient of determination**, or **coefficient of multiple determination** for multiple regression.
- It can be calculated from the below formula:

$$\text{R-squared} = \frac{\text{Explained variation}}{\text{Total Variation}}$$

Assumptions of Linear Regression

Below are some important assumptions of Linear Regression. These are some formal checks while building a Linear Regression model, which ensures to get the best possible result from the given dataset.

- **Linear relationship between the features and target:**
Linear regression assumes the linear relationship between the dependent and independent variables.
- **Small or no multicollinearity between the features:**
Multicollinearity means high-correlation between the independent variables. Due to multicollinearity, it may difficult to find the true relationship between the predictors and target variables. Or we can say, it is difficult to determine which predictor variable is affecting the target variable and which is not. So, the model assumes either little or no multicollinearity between the features or independent variables.
- **Homoscedasticity Assumption:**
Homoscedasticity is a situation when the error term is the same for all the values of independent variables. With homoscedasticity, there should be no clear pattern distribution of data in the scatter plot.
- **Normal distribution of error terms:**
Linear regression assumes that the error term should follow the normal distribution pattern. If error terms are not normally distributed, then confidence intervals will become either too wide or too narrow, which may cause difficulties in finding coefficients.
It can be checked using the **q-q plot**. If the plot shows a straight line without any deviation, which means the error is normally distributed.
- **No autocorrelations:**
The linear regression model assumes no autocorrelation in error terms. If there will be any correlation in the error term, then it will drastically reduce the accuracy of the model.
Autocorrelation usually occurs if there is a dependency between residual errors.

FAQ'S

1. What is meant by Regression?
2. What is meant by Linear Regression?
3. What is cost function?
4. What is standard deviation and normalization?
5. What is difference between classification and regression?

Outcome:

With the completion of this assignment the students will be able to implement linear regression on given dataset.

Conclusion:

Assignment No:- 6

Title of Assignment:

Implement Naïve Bayes theorem to classify the English text

During assignment students will be able to:

- Learn Supervised Learning.
- Implementation of Naïve Bayes Classifier.

Prerequisites: Knowledge of python programming and machine learning libraries.

Concepts to be used: Naïve Bayes Classifier

Input: Text dataset

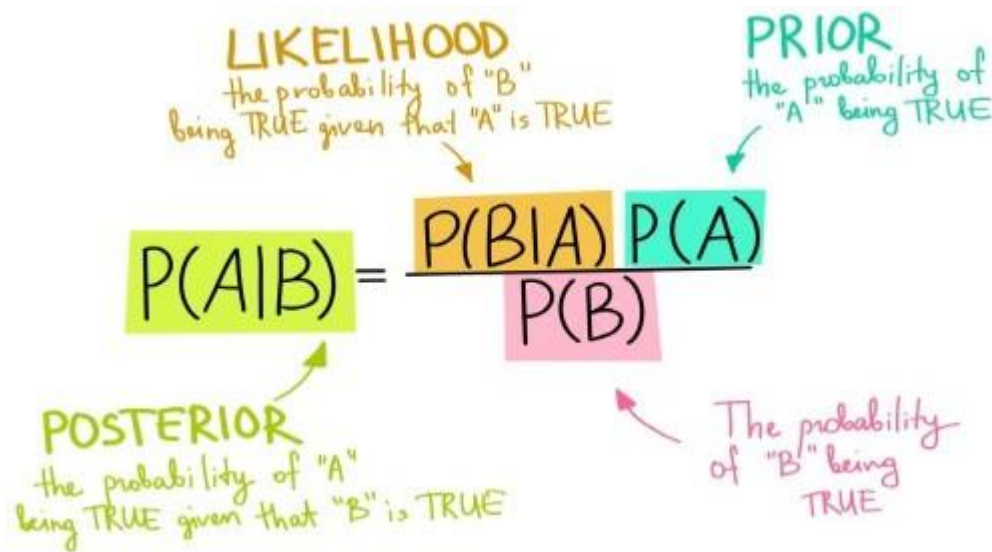
Output: Successful implementation of Naïve Bayes Classifier.

Relevant Theory / Literature Survey:

Naive Bayes

We are going to use Naive Bayes algorithm to classify our text data.

It works on the famous **Bayes theorem** which helps us to find the conditional probabilities of occurrence of two events based on the probabilities of occurrence of each individual event.



Consider we have data of student's effort level(Poor, Average and Good) and their results(Pass and Fail).

EFFORT	RESULT
Poor	Fail
Average	Pass
Average	Pass
Good	Pass
Good	Pass
Poor	Fail
Poor	Fail
Poor	Pass
Poor	Fail
Average	Pass
Average	Fail
Good	Pass
Average	Fail

Student will fail if his efforts are poor. We want to check if this statement is correct. Then we can use Bayes Theorem as,

$$P(\text{Fail} | \text{Poor}) = P(\text{Poor} | \text{Fail}) * P(\text{Fail}) / P(\text{Poor})$$

Processing

$P(\text{Fail} | \text{Poor})$ is read as Probability of getting failed given that the effort was poor.

$$P(\text{Poor} | \text{Fail}) = \text{Number of students who failed with poor efforts} / \text{Number of students failed} \\ = 4 / 6 = 0.66$$

$$P(\text{Fail}) = \text{Number of students failed} / \text{Total students} = 6 / 13$$

$$P(\text{Poor}) = \text{Number of students with poor efforts} / \text{Total students} = 5 / 13$$

$$P(\text{Fail} | \text{Poor}) = 0.66 * (6/13) / (5/13) = 0.66 * 6 / 5 = 0.792$$

This is a higher probability. We use a similar method in Naive Bayes to give the probability of different class and then label it with the class having maximum probability.

Let's take an example, where we want to tell if a fruit is tomato or not. We can tell it's a tomato from its shape, color and diameter(size). Tomato is red, it's round and has about 9-10 cm diameter. These 3 features contribute independently of each other to the probability for the fruit to be tomato. That's why these features are treated as 'Naive'.

Classifying these Naive features using Bayes theorem is known as **Naive Bayes**.

Counting how many times each attribute co-occurs with each class is the main learning idea for Naive Bayes classifier.

How to use Naive Bayes for Text?

In our case, we can't feed in text directly to our classifier. Texts are huge, they have lots of words and various combinations so we use occurrences of words in a single text data (i.e. term frequency-tf) and how many times that word comes across all text data (i.e. inverse document frequency-idf) and then we can generate feature vectors.

Feature vectors representing text contains the probabilities of appearance of the words of the text within the texts of a given category so that the algorithm can compute the likelihood of that text's belonging to the category.

Feature Vector

Document ID	Word 1	Word 2
1	2	0
2	0	1
3	3	1
4	1	3

ID	Class
1	1
2	0
3	0
4	1

1- belongs
0- Does not belong

Once we are done finding likelihoods of various categories we find the category having maximum likelihood and there we have categorized our text!

For example, let's say you have a data set as :

DOCUMENT ID	CONTENT	CLASS
1	Nvidia GPU is the best in the world.	computer graphics
2	Nvidia is giving tough competition to AMD.	computer graphics
3	We were running our application with GTX 1050(High end GPU) still it didn't work then we realized the problem was with the OS.	computer graphics
4	GPU, Ganpat Pandey University, is located in Maharashtra.	not computer graphics
5	Please buy GPU from our store.	?

We can predict the class of last data by using Naive Bayes by considering the probability of important words,

Important words for com.graphics from the data can be considered as,
{Nvidia, GPU, AMD, GTX 1050}

All words will have feature vectors of theirs representing their number of occurrences in each record,
Nvidia -> [1; 1; 0; 0]
and so on.

$P(\text{computer graphics} | \text{GPU}) = P(\text{GPU} | \text{computer graphics}) * (P(\text{computer graphics}) / P(\text{GPU}))$

$P(\text{computer graphics}) =$

Number of records having class as computer graphics / Number of total records = $3/4 = 0.75$

$P(\text{GPU}) = \text{Number of records having GPU} / \text{Total number of records} = 3/4 = 0.75$

$P(\text{GPU} | \text{computer graphics}) =$

Number of records having computer graphics with GPU in it / Total number of computer graphics' records
 $= 2/3 = 0.66$

$P(\text{computer graphics} | \text{GPU}) = 0.66 * (0.75/0.75) = 0.66$

This is greater than 0.5 so we can predict that this text data will be belonging to computer graphics.

Here, we have considered only one word GPU that is important in our test data but in real data we'll consider feature vectors of all words and then compute the probability for the class based on the occurrences of all the non primitive words.

The main advantage is that you can get really good results when data available is not much and computational resources are scarce.

Making it work

For this task, we'll need:

- Python: To run our script
- Pip: Necessary to install Python packages

Once you are done with this, let's install some packages using pip, open your terminal and type in this.

```
pip install numpy
pip install sklearn
Python
Copy
```

Numpy: Useful mathematical functions

Sklearn: Machine learning tools for python

Now, we are ready to get our hands dirty on the script

First, we'll begin by **importing the libraries** necessary.

```
import numpy as np
from sklearn.datasets import fetch_20newsgroups
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.naive_bayes import MultinomialNB
from sklearn.pipeline import Pipeline
Python
Copy
```

Here, we are importing the famous **20 News groups dataset** from the datasets available in sklearn.

```
from sklearn.datasets import fetch_20newsgroups
Python
Copy
```

Now, we define the categories we want to classify our text into and define the training data set using sklearn.

```
# We defined the categories which we want to classify
categories = ['rec.motorcycles', 'sci.electronics',
             'comp.graphics', 'sci.med']

# sklearn provides us with subset data for training and testing
train_data = fetch_20newsgroups(subset='train',
                                categories=categories, shuffle=True, random_state=42)

print(train_data.target_names)

print("\n".join(train_data.data[0].split("\n")[:3]))
print(train_data.target_names[train_data.target[0]])

# Let's look at categories of our first ten training data
```

```
for t in train_data.target[:10]:
    print(train_data.target_names[t])
Python
Copy
```

Output:

```
['comp.graphics', 'rec.motorcycles', 'sci.electronics', 'sci.med']
From: kreyling@lds.loral.com (Ed Kreyling 6966)
Subject: Sun-os and 8bit ASCII graphics
Organization: Loral Data Systems
comp.graphics
comp.graphics
comp.graphics
rec.motorcycles
comp.graphics
sci.med
sci.electronics
sci.electronics
comp.graphics
rec.motorcycles
sci.electronics
Processing
Copy
```

We defined our task into several stages in the beginning. Here we are pre-processing on text and generating feature vectors of token counts and then transform into tf-idf representation.

Consider a document containing 100 words wherein the word ‘car’ appears 7 times. The term frequency (tf) for phone is then $(7 / 100) = 0.07$. Now, assume we have 1 million documents and the word car appears in one thousand of these. Then, the inverse document frequency (i.e., IDF) is calculated as $\log(10,00,000 / 100) = 4$. Thus, the Tf-IDF weight is the product of these quantities: $0.07 * 4 = 0.28$.

i.e. Apply Vectorizer=> Transformer.

```
# Builds a dictionary of features and transforms documents to feature vectors and convert our text documents to
a
# matrix of token counts (CountVectorizer)
count_vect = CountVectorizer()
X_train_counts = count_vect.fit_transform(train_data.data)

# transform a count matrix to a normalized tf-idf representation (tf-idf transformer)
tfidf_transformer = TfidfTransformer()
X_train_tfidf = tfidf_transformer.fit_transform(X_train_counts)
Python
Copy
```

We fit our Multinomial Naive Bayes classifier on train data to train it.

```
# training our classifier ; train_data.target will be having numbers assigned for each category in train data
clf = MultinomialNB().fit(X_train_tfidf, train_data.target)

# Input Data to predict their classes of the given categories
docs_new = ['I have a Harley Davidson and Yamaha.', 'I have a GTX 1050 GPU']
# building up feature vector of our input
X_new_counts = count_vect.transform(docs_new)
# We call transform instead of fit_transform because it's already been fit
X_new_tfidf = tfidf_transformer.transform(X_new_counts)
Python
Copy
```

Predict the output of our input text by using the classifier we just trained.

```
# predicting the category of our input text: Will give out number for category
predicted = clf.predict(X_new_tfidf)

for doc, category in zip(docs_new, predicted):
    print('%r => %s' % (doc, train_data.target_names[category]))
Python
Copy
```

I have a Harley Davidson and Yamaha.' => rec.motorcycles

I have a GTX 1050 GPU' => sci.med

We now finally evaluate our model by predicting the test data. Also, you'll see how to do all of the tasks of vectorizing, transforming and classifier into a single compound classifier using Pipeline.

```
# We can use Pipeline to add vectorizer -> transformer -> classifier all in a one compound classifier
text_clf = Pipeline([
    ('vect', CountVectorizer()),
    ('tfidf', TfidfTransformer()),
    ('clf', MultinomialNB()),
])
# Fitting our train data to the pipeline
text_clf.fit(train_data.data, train_data.target)

# Test data
test_data = fetch_20newsgroups(subset='test',
                               categories=categories, shuffle=True, random_state=42)
docs_test = test_data.data
# Predicting our test data
predicted = text_clf.predict(docs_test)
print('We got an accuracy of', np.mean(predicted == test_data.target)*100, '% over the test data.')
Python
Copy
```

*We got an accuracy of **91.49746192893402** % over the test data.*

With this, you have the complete knowledge of using Naive Bayes classifier from Text Classification tasks.

FAQ's :

1. What is meant by Regression?
2. What is meant by Naïve Bias classification?
3. What is classification?
4. What is dependence & independent variable?
5. What is difference between classification and regression?

Outcome:

With the completion of this assignment the students will be able to implement Naïve bias on given dataset.

Conclusion:

Assignment No:- 7

Title of Assignment:

Unsupervised Learning: Implement K-Means Clustering and Hierarchical clustering on proper data set of your choice. Compare their Convergence

During assignment students will be able to:

- Learn Unsupervised Learning & Clustering.
- Implementation of k-means clustering.

Prerequisites: Knowledge of python programming and machine learning libraries.

Concepts to be used: Unsupervised Learning & Clustering.

Input: given dataset

Output: Successful implementation of k-means clustering.

Relevant Theory / Literature Survey:

K-Means Clustering Algorithm

K-Means Clustering is an unsupervised learning algorithm that is used to solve the clustering problems in machine learning or data science. In this topic, we will learn what is K-means clustering algorithm, how the algorithm works, along with the Python implementation of k-means clustering.

What is K-Means Algorithm?

K-Means Clustering is an [Unsupervised Learning algorithm](#)

, which groups the unlabeled dataset into different clusters. Here K defines the number of pre-defined clusters that need to be created in the process, as if $K=2$, there will be two clusters, and for $K=3$, there will be three clusters, and so on.

It is an iterative algorithm that divides the unlabeled dataset into k different clusters in such a way that each dataset belongs only one group that has similar properties.

It allows us to cluster the data into different groups and a convenient way to discover the categories of groups in the unlabeled dataset on its own without the need for any training.

It is a centroid-based algorithm, where each cluster is associated with a centroid. The main aim of this algorithm is to minimize the sum of distances between the data point and their corresponding clusters.

The algorithm takes the unlabeled dataset as input, divides the dataset into k -number of clusters, and repeats the process until it does not find the best clusters. The value of k should be predetermined in this algorithm.

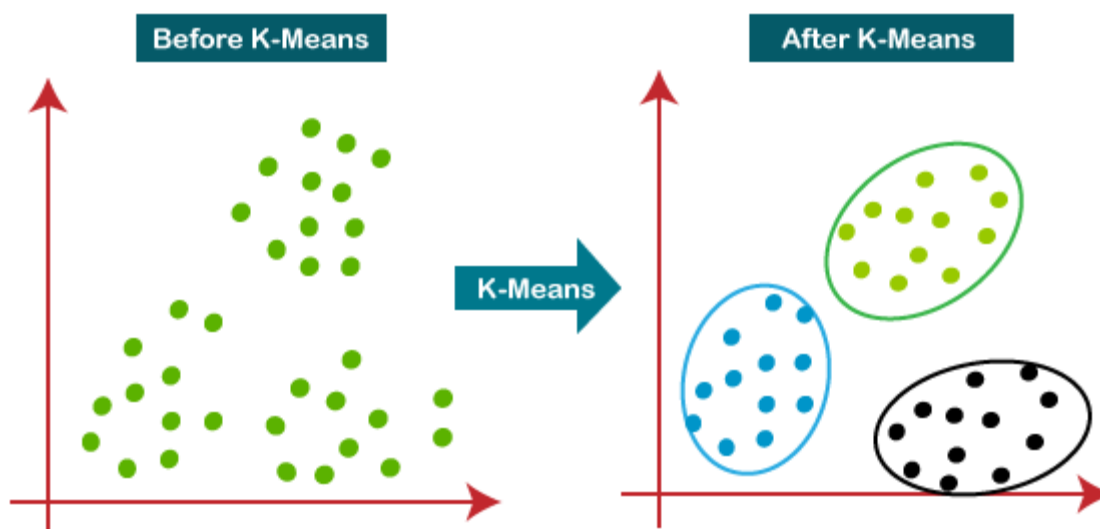
The k -means [clustering](#)

algorithm mainly performs two tasks:

- Determines the best value for K center points or centroids by an iterative process.
- Assigns each data point to its closest k -center. Those data points which are near to the particular k -center, create a cluster.

Hence each cluster has datapoints with some commonalities, and it is away from other clusters.

The below diagram explains the working of the K -means Clustering Algorithm:



How does the K -Means Algorithm Work?

The working of the K-Means algorithm is explained in the below steps:

Step-1: Select the number K to decide the number of clusters.

Step-2: Select random K points or centroids. (It can be other from the input dataset).

Step-3: Assign each data point to their closest centroid, which will form the predefined K clusters.

Step-4: Calculate the variance and place a new centroid of each cluster.

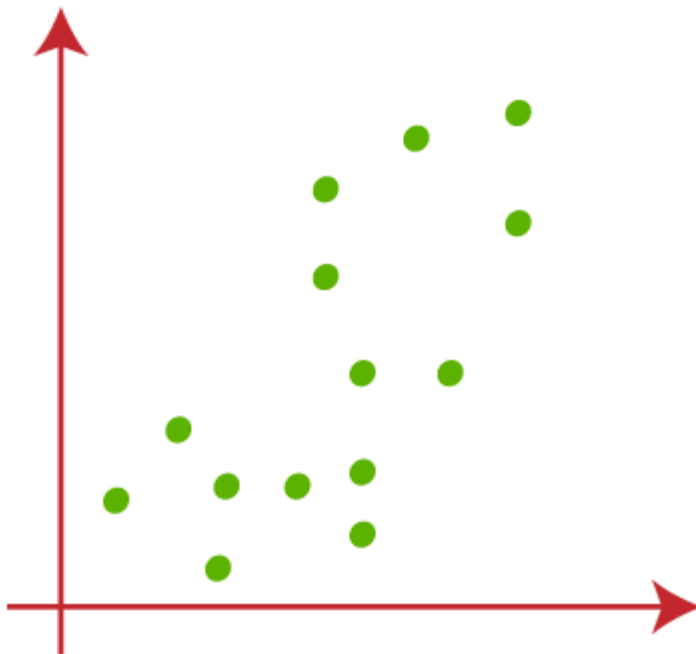
Step-5: Repeat the third steps, which means reassign each datapoint to the new closest centroid of each cluster.

Step-6: If any reassignment occurs, then go to step-4 else go to FINISH.

Step-7: The model is ready.

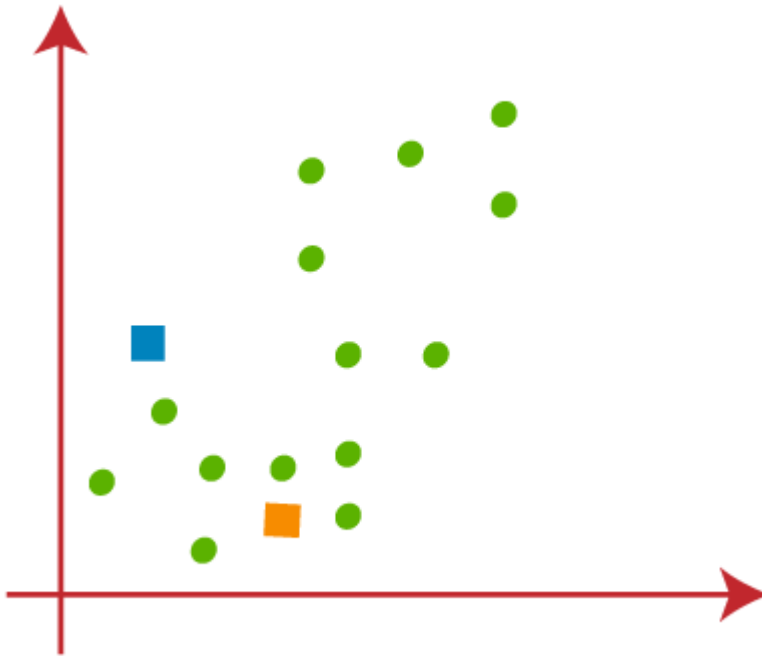
Let's understand the above steps by considering the visual plots:

Suppose we have two variables M1 and M2. The x-y axis scatter plot of these two variables is given below:



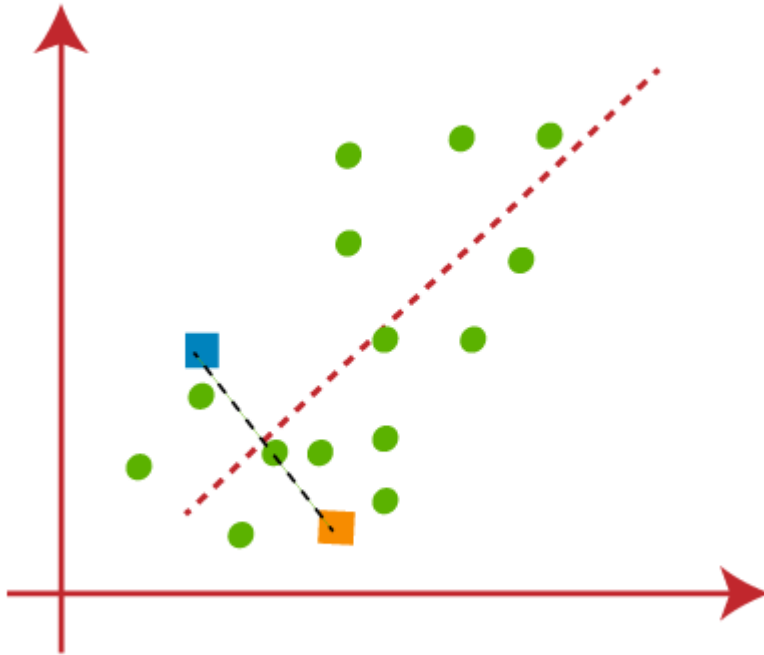
- Let's take number k of clusters, i.e., $K=2$, to identify the dataset and to put them into different clusters. It means here we will try to group these datasets into two different clusters.

- We need to choose some random k points or centroid to form the cluster. These points can be either the points from the dataset or any other point. So, here we are selecting the below two points as k points, which are not the part of our dataset. Consider the below image:

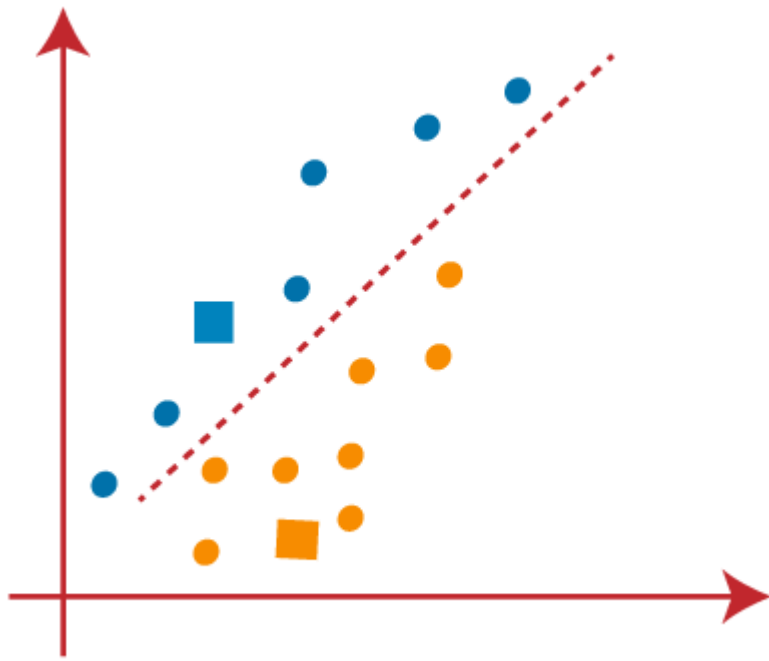


- Now we will assign each data point of the scatter plot to its closest K-point or centroid. We will compute it by applying some mathematics that we have studied to calculate the distance between two points. So, we will draw a median between both the centroids. Consider the below

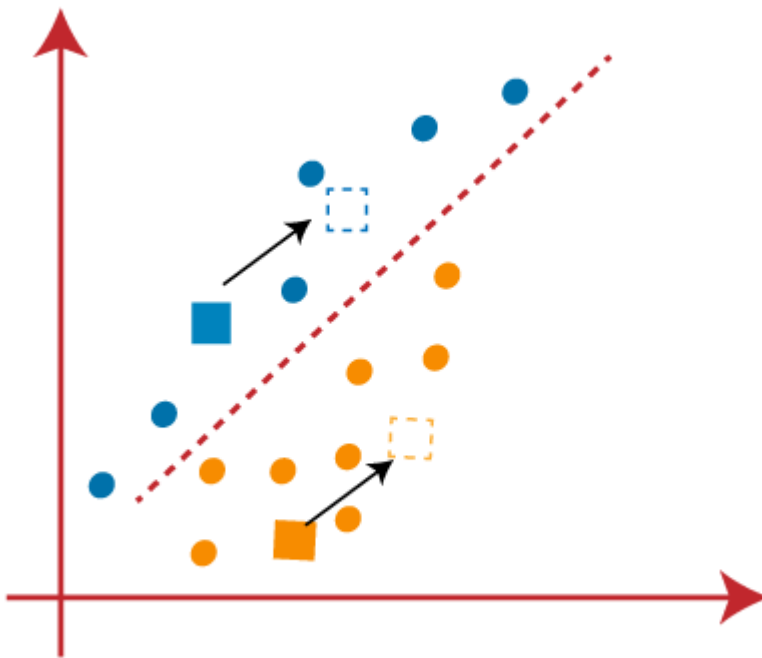
image:



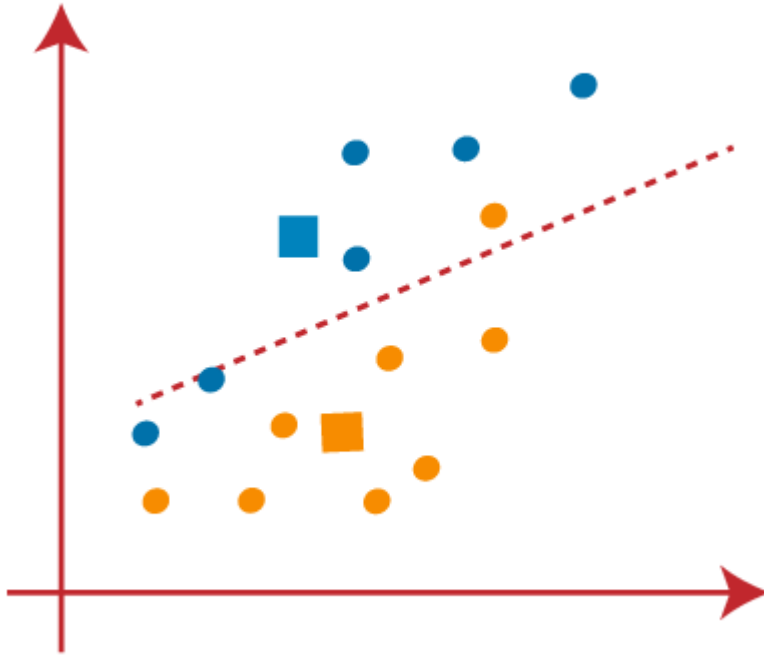
From the above image, it is clear that points left side of the line is near to the K1 or blue centroid, and points to the right of the line are close to the yellow centroid. Let's color them as blue and yellow for clear visualization.



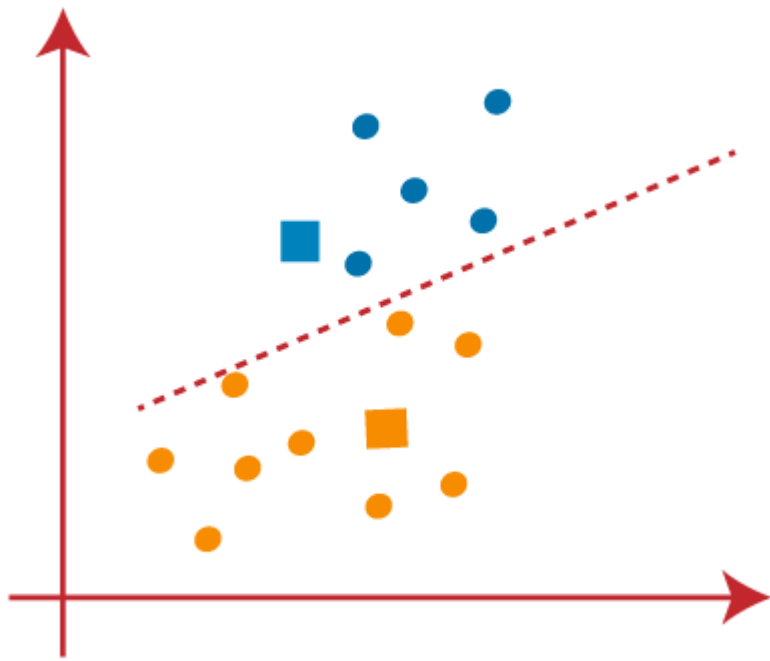
- As we need to find the closest cluster, so we will repeat the process by choosing a **new centroid**. To choose the new centroids, we will compute the center of gravity of these centroids, and will find new centroids as below:



- Next, we will reassign each datapoint to the new centroid. For this, we will repeat the same process of finding a median line. The median will be like below image:

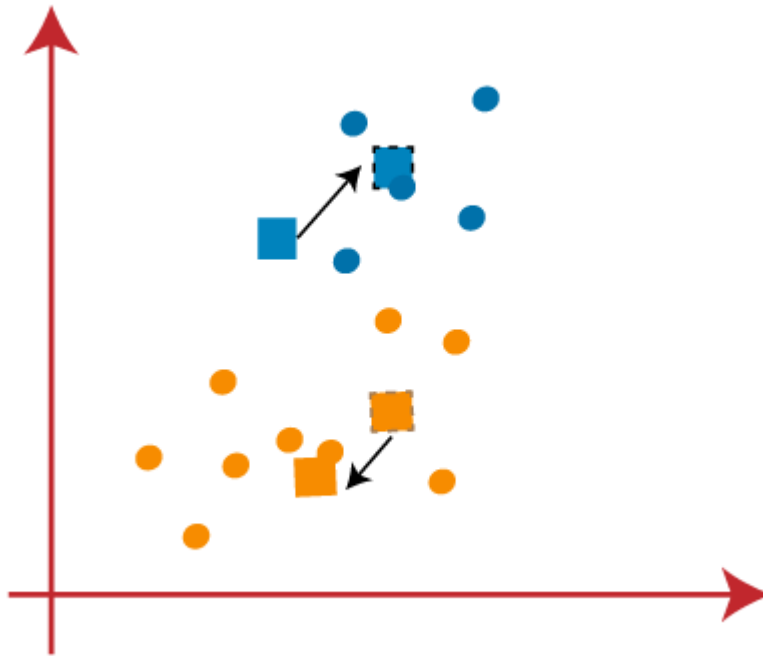


From the above image, we can see, one yellow point is on the left side of the line, and two blue points are right to the line. So, these three points will be assigned to new centroids.

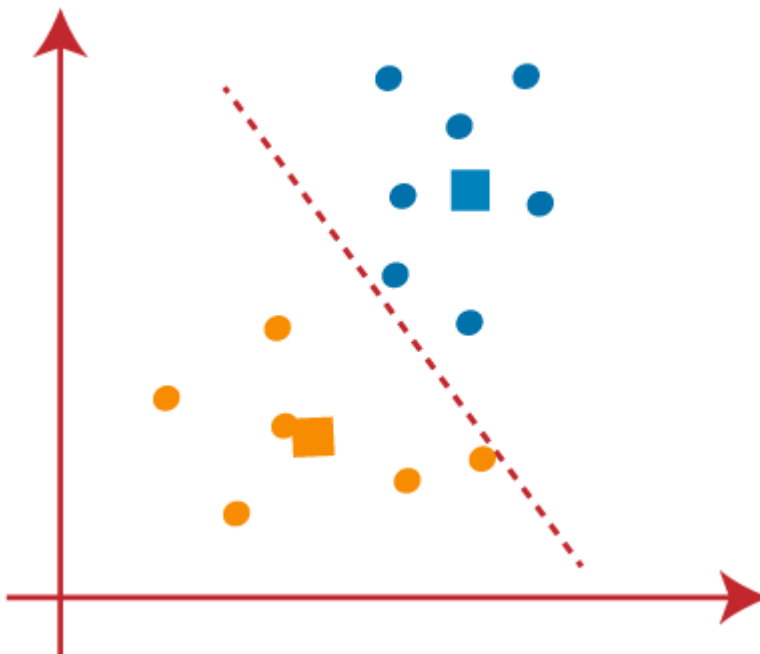


As reassignment has taken place, so we will again go to the step-4, which is finding new centroids or K-points.

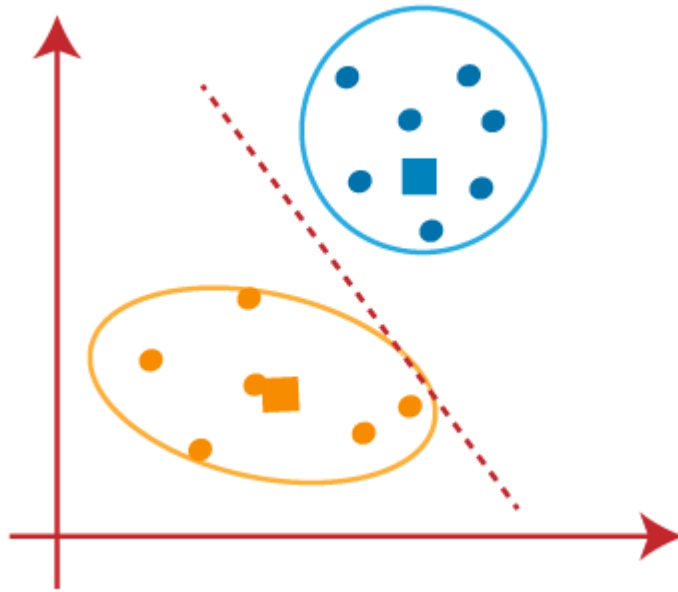
- We will repeat the process by finding the center of gravity of centroids, so the new centroids will be as shown in the below image:



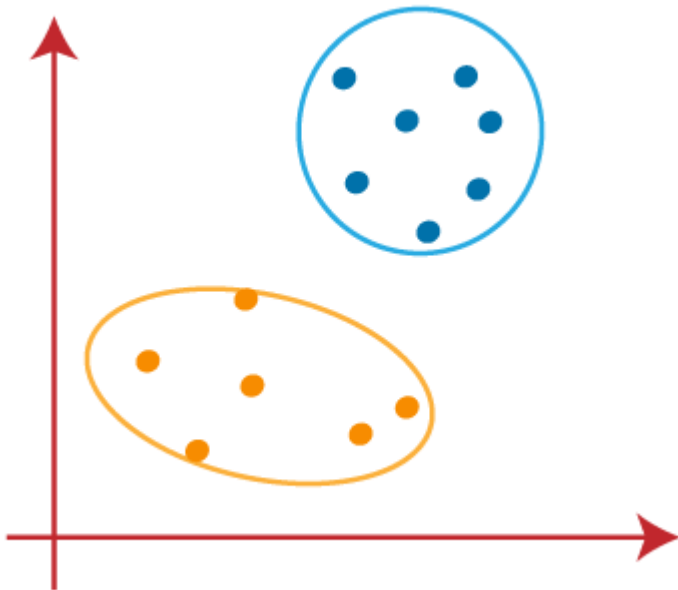
- As we got the new centroids so again will draw the median line and reassign the data points. So, the image will be:



- We can see in the above image; there are no dissimilar data points on either side of the line, which means our model is formed. Consider the below image:



As our model is ready, so we can now remove the assumed centroids, and the two final clusters will be as shown in the below image:



How to choose the value of "K number of clusters" in K-means Clustering?

The performance of the K-means clustering algorithm depends upon highly efficient clusters that it forms. But choosing the optimal number of clusters is a big task. There are some different ways to find the optimal number of clusters, but here we are discussing the most appropriate method to find the number of clusters or value of K. The method is given below:

Elbow Method

The Elbow method is one of the most popular ways to find the optimal number of clusters. This method uses the concept of WCSS value. **WCSS** stands for **Within Cluster Sum of Squares**, which defines the total variations within a cluster. The formula to calculate the value of WCSS (for 3 clusters) is given below:

$$WCSS = \sum_{P_i \text{ in Cluster1}} \text{distance}(P_i C_1)^2 + \sum_{P_i \text{ in Cluster2}} \text{distance}(P_i C_2)^2 + \sum_{P_i \text{ in Cluster3}} \text{distance}(P_i C_3)^2$$

In the above formula of WCSS,

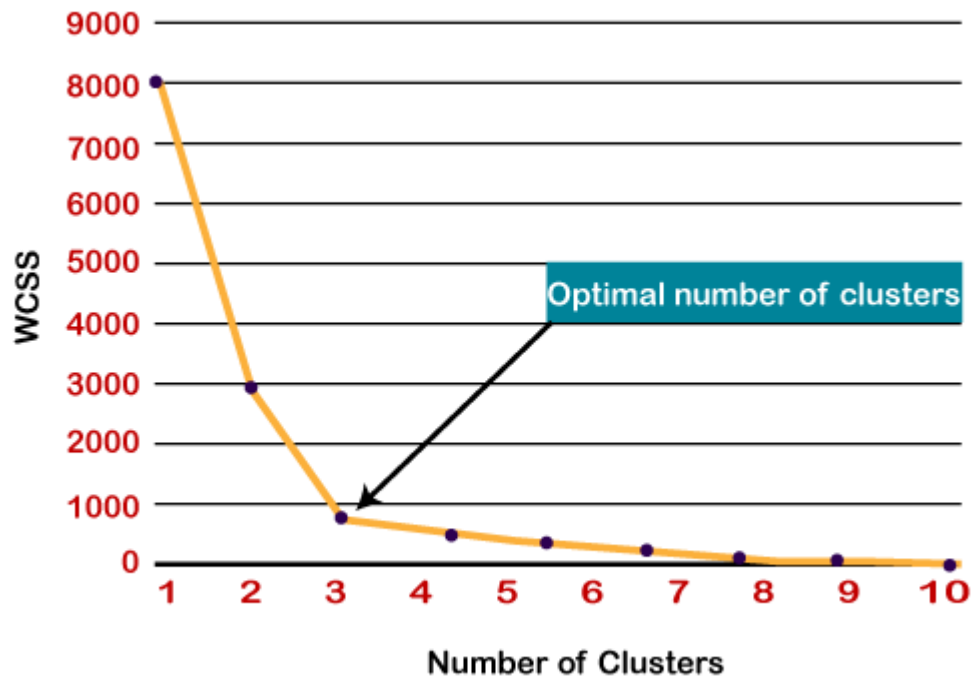
$\sum_{P_i \text{ in Cluster1}} \text{distance}(P_i C_1)^2$: It is the sum of the square of the distances between each data point and its centroid within a cluster1 and the same for the other two terms.

To measure the distance between data points and centroid, we can use any method such as Euclidean distance or Manhattan distance.

To find the optimal value of clusters, the elbow method follows the below steps:

- It executes the K-means clustering on a given dataset for different K values (ranges from 1-10).
- For each value of K, calculates the WCSS value.
- Plots a curve between calculated WCSS values and the number of clusters K.
- The sharp point of bend or a point of the plot looks like an arm, then that point is considered as the best value of K.

Since the graph shows the sharp bend, which looks like an elbow, hence it is known as the elbow method. The graph for the elbow method looks like the below image:



Note: We can choose the number of clusters equal to the given data points. If we choose the number of clusters equal to the data points, then the value of WCSS becomes zero, and that will be the endpoint of the plot.

Python Implementation of K-means Clustering Algorithm

In the above section, we have discussed the K-means algorithm, now let's see how it can be implemented using [Python](#)

.

Before implementation, let's understand what type of problem we will solve here. So, we have a dataset of **Mall_Customers**, which is the data of customers who visit the mall and spend there.

In the given dataset, we have **Customer_Id**, **Gender**, **Age**, **Annual Income (\$)**, and **Spending Score** (which is the calculated value of how much a customer has spent in the mall, the more the value, the more he has spent). From this dataset, we need to calculate some patterns, as it is an unsupervised method, so we don't know what to calculate exactly.

The steps to be followed for the implementation are given below:

- **Data Pre-processing**
- **Finding the optimal number of clusters using the elbow method**
- **Training the K-means algorithm on the training dataset**
- **Visualizing the clusters**

Step-1: Data pre-processing Step

The first step will be the data pre-processing, as we did in our earlier topics of Regression and Classification. But for the clustering problem, it will be different from other models. Let's discuss it:

- **Importing Libraries**

As we did in previous topics, firstly, we will import the libraries for our model, which is part of data pre-processing. The code is given below:

5. # importing libraries
6. **import** numpy as nm
7. **import** matplotlib.pyplot as mtp
8. **import** pandas as pd

In the above code, the **numpy** we have imported for the performing mathematics calculation, **matplotlib** is for plotting the graph, and **pandas** are for managing the dataset.

- **Importing the Dataset:**

Next, we will import the dataset that we need to use. So here, we are using the Mall_Customer_data.csv dataset. It can be imported using the below code:

3. # Importing the dataset
4. dataset = pd.read_csv('Mall_Customers_data.csv')

By executing the above lines of code, we will get our dataset in the Spyder IDE. The dataset looks like the below image:

Index	CustomerID	Genre	Age	Annual Income (k\$)	Spending Score (1-100)
0	1	Male	19	15	39
1	2	Male	21	15	81
2	3	Female	20	16	6
3	4	Female	23	16	77
4	5	Female	31	17	40
5	6	Female	22	17	76
6	7	Female	35	18	6
7	8	Female	23	18	94
8	9	Male	64	19	3
9	10	Female	30	19	72
10	11	Male	67	19	14
11	12	Female	35	19	99
12	13	Female	58	20	15
13	14	Female	24	20	77
14	15	Male	37	20	13
15	16	Male	22	20	79

From the above dataset, we need to find some patterns in it.

- **Extracting Independent Variables**

Here we don't need any dependent variable for data pre-processing step as it is a clustering problem, and we have no idea about what to determine. So we will just add a line of code for the matrix of features.

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

As we can see, we are extracting only 3rd and 4th feature. It is because we need a 2d plot to visualize the model, and some features are not required, such as customer_id.

Step-2: Finding the optimal number of clusters using the elbow method

In the second step, we will try to find the optimal number of clusters for our clustering problem. So, as discussed above, here we are going to use the elbow method for this purpose.

As we know, the elbow method uses the WCSS concept to draw the plot by plotting WCSS values on the Y-axis and the number of clusters on the X-axis. So we are going to calculate the value for WCSS for different k values ranging from 1 to 10. Below is the code for it:

```
15. #finding optimal number of clusters using the elbow method
16. from sklearn.cluster import KMeans
17. wcss_list= [] #Initializing the list for the values of WCSS
18.
19. #Using for loop for iterations from 1 to 10.
20. for i in range(1, 11):
21.     kmeans = KMeans(n_clusters=i, init='k-means++', random_state= 42)
22.     kmeans.fit(x)
23.     wcss_list.append(kmeans.inertia_)
24. mtp.plot(range(1, 11), wcss_list)
25. mtp.title('The Elbow Method Graph')
26. mtp.xlabel('Number of clusters(k)')
27. mtp.ylabel('wcss_list')
28. mtp.show()
```

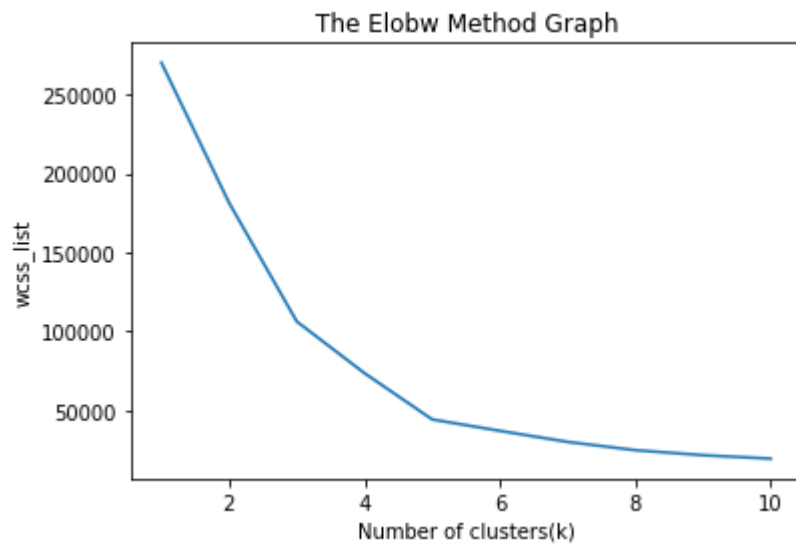
As we can see in the above code, we have used the **KMeans** class of sklearn. cluster library to form the clusters.

Next, we have created the **wcss_list** variable to initialize an empty list, which is used to contain the value of wcss computed for different values of k ranging from 1 to 10.

After that, we have initialized the for loop for the iteration on a different value of k ranging from 1 to 10; since for loop in Python, exclude the outbound limit, so it is taken as 11 to include 10th value.

The rest part of the code is similar as we did in earlier topics, as we have fitted the model on a matrix of features and then plotted the graph between the number of clusters and WCSS.

Output: After executing the above code, we will get the below output:



From the above plot, we can see the elbow point is at **5**. So the number of clusters here will be **5**.

wcss_list - List (10 elements)			
Index	Type	Size	Value
0	float64	1	269981.28
1	float64	1	181363.59595959596
2	float64	1	106348.37306211118
3	float64	1	73679.78903948834
4	float64	1	44448.45544793371
5	float64	1	37233.81451071001
6	float64	1	30259.65720728547
7	float64	1	25011.83934915659
8	float64	1	21850.165282585633
9	float64	1	19672.07284901432

Save and Close
Close

Step- 3: Training the K-means algorithm on the training dataset

As we have got the number of clusters, so we can now train the model on the dataset.

To train the model, we will use the same two lines of code as we have used in the above section, but here instead of using `i`, we will use `5`, as we know there are 5 clusters that need to be formed. The code is given below:

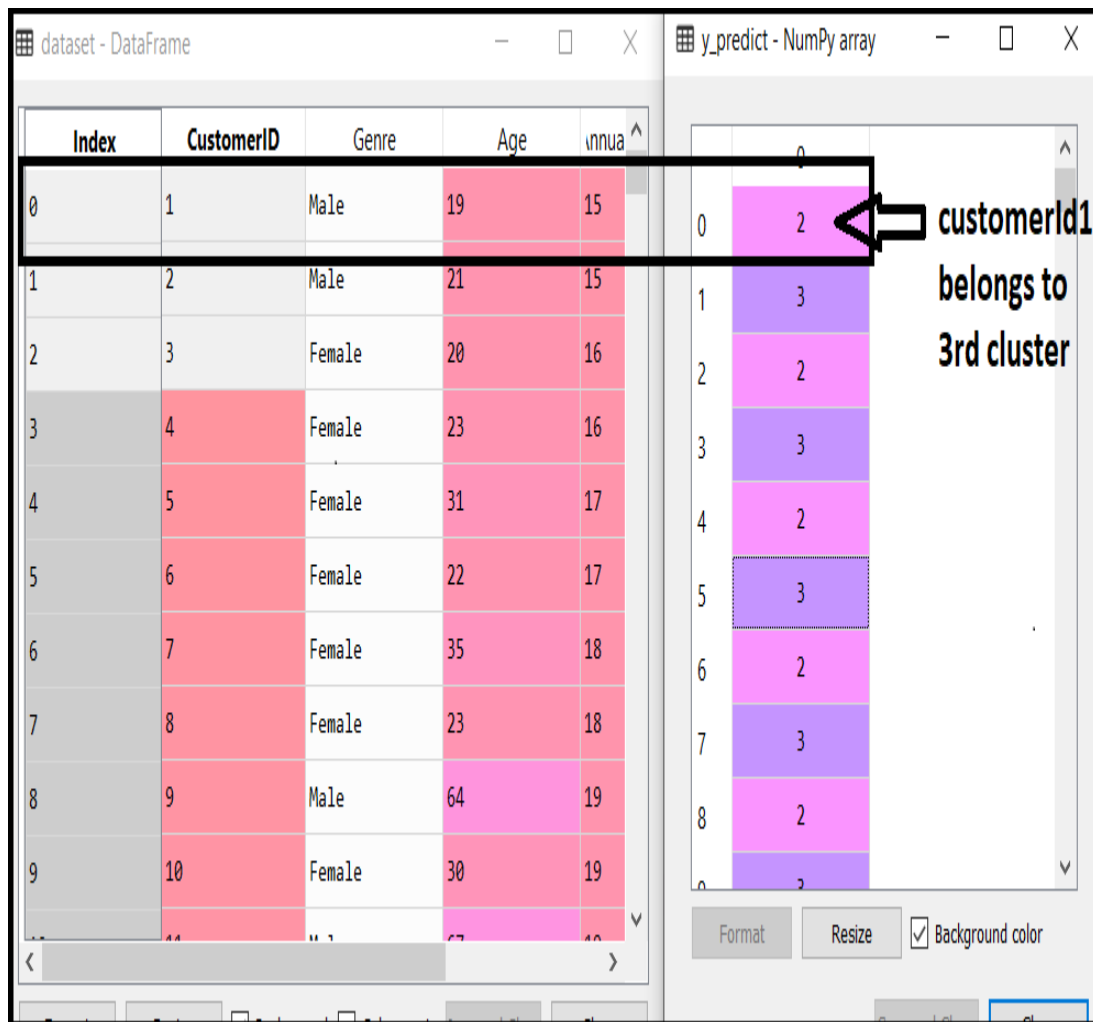
- #training the K-means model on a dataset
- `kmeans = KMeans(n_clusters=5, init='k-means++', random_state= 42)`

6. `y_predict= kmeans.fit_predict(x)`

The first line is the same as above for creating the object of KMeans class.

In the second line of code, we have created the dependent variable **y_predict** to train the model.

By executing the above lines of code, we will get the `y_predict` variable. We can check it under **the variable explorer** option in the Spyder IDE. We can now compare the values of `y_predict` with our original dataset. Consider the below image:



From the above image, we can now relate that the CustomerID 1 belongs to a cluster

3(as index starts from 0, hence 2 will be considered as 3), and 2 belongs to cluster 4, and so on.

Step-4: Visualizing the Clusters

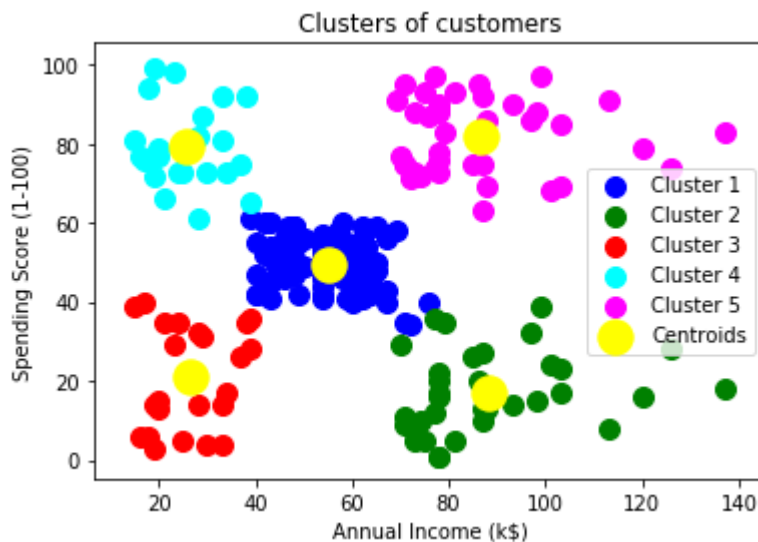
The last step is to visualize the clusters. As we have 5 clusters for our model, so we will visualize each cluster one by one.

To visualize the clusters will use scatter plot using `mtp.scatter()` function of matplotlib.

```
13. #visulaizing the clusters
14. mtp.scatter(x[y_predict == 0, 0], x[y_predict == 0, 1], s = 100, c = 'blue', label = 'Cluster 1') #for first cluster
15. mtp.scatter(x[y_predict == 1, 0], x[y_predict == 1, 1], s = 100, c = 'green', label = 'Cluster 2') #for second cluster
16. mtp.scatter(x[y_predict == 2, 0], x[y_predict == 2, 1], s = 100, c = 'red', label = 'Cluster 3') #for third cluster
17. mtp.scatter(x[y_predict == 3, 0], x[y_predict == 3, 1], s = 100, c = 'cyan', label = 'Cluster 4') #for fourth cluster
18. mtp.scatter(x[y_predict == 4, 0], x[y_predict == 4, 1], s = 100, c = 'magenta', label = 'Cluster 5') #for fifth cluster
19. mtp.scatter(kmeans.cluster_centers[:, 0], kmeans.cluster_centers[:, 1], s = 300, c = 'yellow', label = 'Centroid')
20. mtp.title('Clusters of customers')
21. mtp.xlabel('Annual Income (k$)')
22. mtp.ylabel('Spending Score (1-100)')
23. mtp.legend()
24. mtp.show()
```

In above lines of code, we have written code for each clusters, ranging from 1 to 5. The first coordinate of the `mtp.scatter`, i.e., `x[y_predict == 0, 0]` containing the x value for the showing the matrix of features values, and the `y_predict` is ranging from 0 to 1.

Output:



The output image is clearly showing the five different clusters with different colors. The clusters are formed between two parameters of the dataset; Annual income of customer and Spending. We can change the colors and labels as per the requirement or choice. We can also observe some points from the above patterns, which are given below:

- **Cluster1** shows the customers with average salary and average spending so we can categorize these customers as
- Cluster2 shows the customer has a high income but low spending, so we can categorize them as **careful**.
- Cluster3 shows the low income and also low spending so they can be categorized as sensible.
- Cluster4 shows the customers with low income with very high spending so they can be categorized as **careless**.
- Cluster5 shows the customers with high income and high spending so they can be categorized as target, and these customers can be the most profitable customers for the mall owner.

FAQ'S

4. What is meant by Unsupervised learning?
5. What is meant by clustering?
6. What is k-means clustering?

Outcome:

With the completion of this assignment the students will be able to implement k-means clustering on given dataset.

Conclusion:

Assignment No:- 8

Title of Assignment:

Implement Support Vector Machine algorithm using suitable dataset.

During assignment students will be able to:

- Learn Classification.
- Implementation of Support Vector Machine.

Prerequisites: Knowledge of python programming and machine learning libraries.

Concepts to be used: Classification.

Input: Any dataset

Output: Successful implementation of Support Vector Machine (SVM)

Relevant Theory / Literature Survey:

SVM offers very high accuracy compared to other classifiers such as logistic regression, and decision trees. It is known for its kernel trick to handle nonlinear input spaces. It is used in a variety of applications such as face detection, intrusion detection, classification of emails, news articles and web pages, classification of genes, and handwriting recognition.

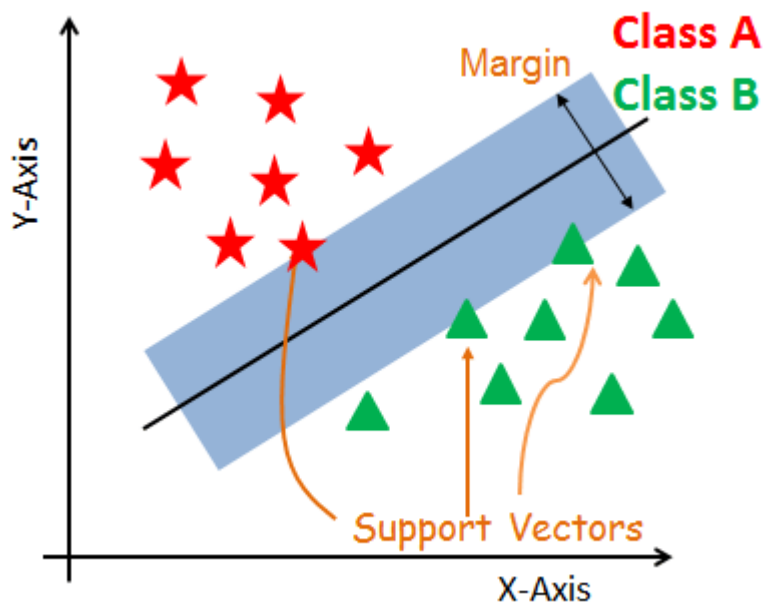
SVM is an exciting algorithm and the concepts are relatively simple. The classifier separates data points using a hyperplane with the largest amount of margin. That's why an SVM classifier is also known as a discriminative classifier. SVM finds an optimal hyperplane which helps in classifying new data points.

In this practical, you are going to cover following topics:

- Support Vector Machines
- How does it work?
- Kernels
- Classifier building in Scikit-learn
- Tuning Hyperparameters
- Advantages and Disadvantages

Support Vector Machines

Generally, Support Vector Machines is considered to be a classification approach, it but can be employed in both types of classification and regression problems. It can easily handle multiple continuous and categorical variables. SVM constructs a hyperplane in multidimensional space to separate different classes. SVM generates optimal hyperplane in an iterative manner, which is used to minimize an error. The core idea of SVM is to find a maximum marginal hyperplane(MMH) that best divides the dataset into classes.



Support Vectors

Support vectors are the data points, which are closest to the hyperplane. These points will define the separating line better by calculating margins. These points are more relevant to the construction of the classifier.

Hyperplane

A hyperplane is a decision plane which separates between a set of objects having different class memberships.

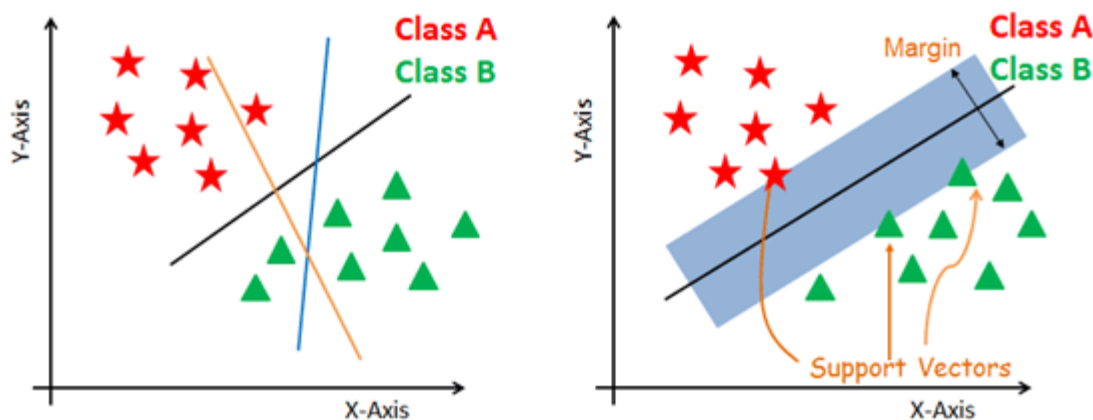
Margin

A margin is a gap between the two lines on the closest class points. This is calculated as the perpendicular distance from the line to support vectors or closest points. If the margin is larger in between the classes, then it is considered a good margin, a smaller margin is a bad margin.

How does SVM work?

The main objective is to segregate the given dataset in the best possible way. The distance between the either nearest points is known as the margin. The objective is to select a hyperplane with the maximum possible margin between support vectors in the given dataset. SVM searches for the maximum marginal hyperplane in the following steps:

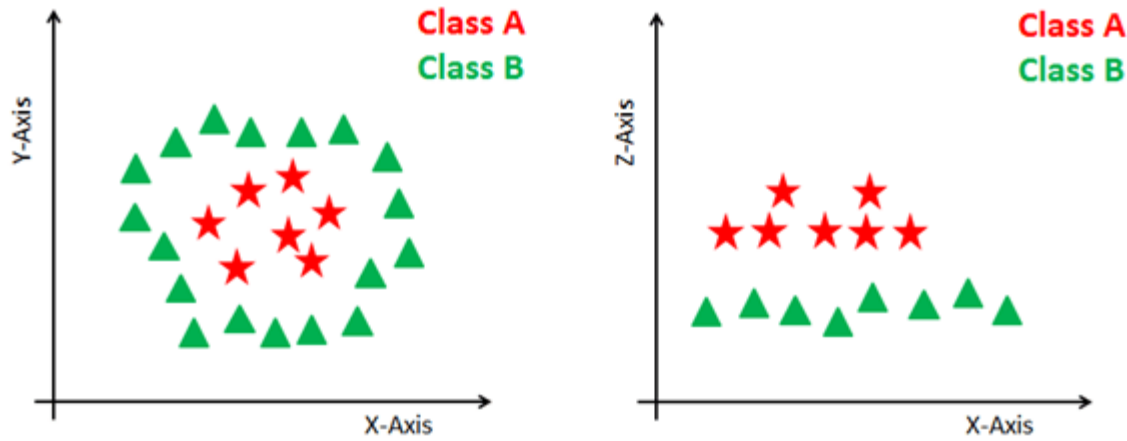
1. Generate hyperplanes which segregates the classes in the best way. Left-hand side figure showing three hyperplanes black, blue and orange. Here, the blue and orange have higher classification error, but the black is separating the two classes correctly.
2. Select the right hyperplane with the maximum segregation from the either nearest data points as shown in the right-hand side figure.



Dealing with non-linear and inseparable planes

Some problems can't be solved using linear hyperplane, as shown in the figure below (left-hand side).

In such situation, SVM uses a kernel trick to transform the input space to a higher dimensional space as shown on the right. The data points are plotted on the x-axis and z-axis (Z is the squared sum of both x and y : $z=x^2+y^2$). Now you can easily segregate these points using linear separation.



SVM Kernels

The SVM algorithm is implemented in practice using a kernel. A kernel transforms an input data space into the required form. SVM uses a technique called the kernel trick. Here, the kernel takes a low-dimensional input space and transforms it into a higher dimensional space. In other words, you can say that it converts nonseparable problem to separable problems by adding more dimension to it. It is most useful in non-linear separation problem. Kernel trick helps you to build a more accurate classifier.

- **Linear Kernel** A linear kernel can be used as normal dot product any two given observations. The product between two vectors is the sum of the multiplication of each pair of input values.

$$K(x, x_i) = \sum (x * x_i)$$

- **Polynomial Kernel** A polynomial kernel is a more generalized form of the linear kernel. The polynomial kernel can distinguish curved or nonlinear input space.

$$K(x, x_i) = 1 + \sum (x * x_i)^d$$

Where d is the degree of the polynomial. d=1 is similar to the linear transformation. The degree needs to be manually specified in the learning algorithm.

- **Radial Basis Function Kernel** The Radial basis function kernel is a popular kernel function commonly used in support vector machine classification. RBF can map an input space in infinite dimensional space.

$$K(x, x_i) = \exp(-\gamma * \sum ((x - x_i)^2))$$

Here gamma is a parameter, which ranges from 0 to 1. A higher value of gamma will perfectly fit the training dataset, which causes over-fitting. Gamma=0.1 is considered to be a good default value. The value of gamma needs to be manually specified in the learning algorithm.

Classifier Building in Scikit-learn

Until now, you have learned about the theoretical background of SVM. Now you will learn about its implementation in Python using scikit-learn.

In the model the building part, you can use the cancer dataset, which is a very famous multi-class classification problem. This dataset is computed from a digitized image of a fine needle aspirate (FNA) of a breast mass. They describe characteristics of the cell nuclei present in the image.

The dataset comprises 30 features (mean radius, mean texture, mean perimeter, mean area, mean smoothness, mean compactness, mean concavity, mean concave points, mean symmetry, mean fractal dimension, radius error, texture error, perimeter error, area error, smoothness error, compactness error, concavity error, concave points error, symmetry error, fractal dimension error, worst radius, worst texture, worst perimeter, worst area, worst smoothness, worst compactness, worst concavity, worst concave points, worst symmetry, and worst fractal dimension) and a target (type of cancer).

This data has two types of cancer classes: malignant (harmful) and benign (not harmful). Here, you can build a model to classify the type of cancer. The dataset is available in the scikit-learn library or you can also download it from the UCI Machine Learning Library.

Loading Data

Let's first load the required dataset you will use.

```
#Import scikit-learn dataset library
from sklearn import datasets

#Load dataset
cancer = datasets.load_breast_cancer()
```

Exploring Data

After you have loaded the dataset, you might want to know a little bit more about it. You can check feature and target names.

```
# print the names of the 13 features
print("Features: ", cancer.feature_names)

# print the label type of cancer('malignant' 'benign')
print("Labels: ", cancer.target_names)

Features:  ['mean radius' 'mean texture' 'mean perimeter' 'mean area'
            'mean smoothness' 'mean compactness' 'mean concavity'
            'mean concave points' 'mean symmetry' 'mean fractal dimension'
            'radius error' 'texture error' 'perimeter error' 'area error'
            'smoothness error' 'compactness error' 'concavity error'
            'concave points error' 'symmetry error' 'fractal dimension error'
            'worst radius' 'worst texture' 'worst perimeter' 'worst area'
            'worst smoothness' 'worst compactness' 'worst concavity'
            'worst concave points' 'worst symmetry' 'worst fractal dimension']
Labels:  ['malignant' 'benign']
```

Let's explore it for a bit more. You can also check the shape of the dataset using shape.

```
# print data(feature) shape
cancer.data.shape

(569, 30)
```

Let's check top 5 records of the feature set.

```
# print the cancer data features (top 5 records)
print(cancer.data[0:5])

[[1.799e+01  1.038e+01  1.228e+02  1.001e+03  1.184e-01  2.776e-01  3.001e-01
  1.471e-01  2.419e-01  7.871e-02  1.095e+00  9.053e-01  8.589e+00  1.534e+02
  6.399e-03  4.904e-02  5.373e-02  1.587e-02  3.003e-02  6.193e-03  2.538e+01
```

```

1.733e+01 1.846e+02 2.019e+03 1.622e-01 6.656e-01 7.119e-01 2.654e-01
4.601e-01 1.189e-01]
[2.057e+01 1.777e+01 1.329e+02 1.326e+03 8.474e-02 7.864e-02 8.690e-02
7.017e-02 1.812e-01 5.667e-02 5.435e-01 7.339e-01 3.398e+00 7.408e+01
5.225e-03 1.308e-02 1.860e-02 1.340e-02 1.389e-02 3.532e-03 2.499e+01
2.341e+01 1.588e+02 1.956e+03 1.238e-01 1.866e-01 2.416e-01 1.860e-01
2.750e-01 8.902e-02]
[1.969e+01 2.125e+01 1.300e+02 1.203e+03 1.096e-01 1.599e-01 1.974e-01
1.279e-01 2.069e-01 5.999e-02 7.456e-01 7.869e-01 4.585e+00 9.403e+01
6.150e-03 4.006e-02 3.832e-02 2.058e-02 2.250e-02 4.571e-03 2.357e+01
2.553e+01 1.525e+02 1.709e+03 1.444e-01 4.245e-01 4.504e-01 2.430e-01
3.613e-01 8.758e-02]
[1.142e+01 2.038e+01 7.758e+01 3.861e+02 1.425e-01 2.839e-01 2.414e-01
1.052e-01 2.597e-01 9.744e-02 4.956e-01 1.156e+00 3.445e+00 2.723e+01
9.110e-03 7.458e-02 5.661e-02 1.867e-02 5.963e-02 9.208e-03 1.491e+01
2.650e+01 9.887e+01 5.677e+02 2.098e-01 8.663e-01 6.869e-01 2.575e-01
6.638e-01 1.730e-01]
[2.029e+01 1.434e+01 1.351e+02 1.297e+03 1.003e-01 1.328e-01 1.980e-01
1.043e-01 1.809e-01 5.883e-02 7.572e-01 7.813e-01 5.438e+00 9.444e+01
1.149e-02 2.461e-02 5.688e-02 1.885e-02 1.756e-02 5.115e-03 2.254e+01
1.667e+01 1.522e+02 1.575e+03 1.374e-01 2.050e-01 4.000e-01 1.625e-01
2.364e-01 7.678e-02]]

```

Let's take a look at the target set.

```

# print the cancer labels (0:malignant, 1:benign)
print(cancer.target)
[0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 0 0 0 0 0 0 0 0 0 0 0 0
1 0 0 0 0 0 0 0 0 1 0 1 1 1 1 1 0 0 1 0 0 1 1 1 1 0 1 0 0 1 1 1 0 1 0 0

```

```

1 0 1 0 0 1 1 1 0 0 1 0 0 0 1 1 1 0 1 1 0 0 1 1 1 1 0 1 1 0 1 1
1 1 1 1 1 1 0 0 0 1 0 0 1 1 1 0 0 1 0 1 0 0 1 0 0 1 1 0 1 1 0 1 1 1 0 1
1 1 1 1 1 1 1 1 0 1 1 1 1 0 0 1 0 1 1 0 0 1 1 0 0 1 1 1 1 0 1 1 0 0 0 1 0
1 0 1 1 1 0 1 1 0 0 1 0 0 0 0 1 0 0 0 1 0 1 0 1 1 0 1 0 0 0 0 1 1 0 0 1 1
1 0 1 1 1 1 1 0 0 1 1 0 1 1 0 0 1 0 1 1 1 1 0 1 1 1 1 1 0 1 0 0 0 0 0 0 0
0 0 0 0 0 0 0 1 1 1 1 1 1 1 0 1 0 1 1 0 1 1 0 1 0 0 1 1 1 1 1 1 1 1 1 1 1
1 0 1 1 0 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 0 1 0 1 1 1 1 0 0 0 1 1
1 1 0 1 0 1 0 1 1 1 0 1 1 1 1 1 1 1 1 0 0 0 1 1 1 1 1 1 1 1 1 1 0 0 1 0 0
0 1 0 0 1 1 1 1 1 0 1 1 1 1 1 1 0 1 1 1 0 1 1 0 0 1 1 1 1 1 1 0 1 1 1 1 1
1 0 1 1 1 1 1 0 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 0 1 0 0 1 0 1 1 1 1 1 0 1 1
0 1 0 1 1 0 1 0 1 1 1 1 1 1 1 1 1 0 0 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 0 1
1 1 1 1 1 1 0 1 0 1 1 1 0 1 1 1 1 1 1 0 0 1 0 1 0 1 1 1 1 1 0 1 1 0 1 0 0
1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 0 1 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1 1 1 1 0 0 0 0 0 0 0 1]

```

Splitting Data

To understand model performance, dividing the dataset into a training set and a test set is a good strategy.

Split the dataset by using the function `train_test_split()`. you need to pass 3 parameters features, target, and test_set size. Additionally, you can use random_state to select records randomly.

```

# Import train_test_split function

from sklearn.model_selection import train_test_split

# Split dataset into training set and test set

X_train, X_test, y_train, y_test = train_test_split(cancer.data,
cancer.target, test_size=0.3, random_state=109) # 70% training and 30% test

```

Generating Model

Let's build support vector machine model. First, import the SVM module and create support vector classifier object by passing argument kernel as the linear kernel in `SVC()` function.

Then, fit your model on train set using `fit()` and perform prediction on the test set using `predict()`.

```
#Import svm model
from sklearn import svm

#Create a svm Classifier
clf = svm.SVC(kernel='linear') # Linear Kernel

#Train the model using the training sets
clf.fit(X_train, y_train)

#Predict the response for test dataset
y_pred = clf.predict(X_test)
```

Evaluating the Model

Let's estimate how accurately the classifier or model can predict the breast cancer of patients.

Accuracy can be computed by comparing actual test set values and predicted values.

```
#Import scikit-learn metrics module for accuracy calculation
from sklearn import metrics

# Model Accuracy: how often is the classifier correct?
print("Accuracy:",metrics.accuracy_score(y_test, y_pred))

Accuracy: 0.9649122807017544
```

Well, you got a classification rate of 96.49%, considered as very good accuracy.

For further evaluation, you can also check precision and recall of model.

```
# Model Precision: what percentage of positive tuples are labeled as such?
```



```
print("Precision:",metrics.precision_score(y_test, y_pred))
```

```
# Model Recall: what percentage of positive tuples are labelled as such?
```

```
print("Recall:",metrics.recall_score(y_test, y_pred))
```

```
Precision: 0.9811320754716981
```

```
Recall: 0.9629629629629629
```

Well, you got a precision of 98% and recall of 96%, which are considered as very good values.

Tuning Hyperparameters

- **Kernel:** The main function of the kernel is to transform the given dataset input data into the required form. There are various types of functions such as linear, polynomial, and radial basis function (RBF). Polynomial and RBF are useful for non-linear hyperplane. Polynomial and RBF kernels compute the separation line in the higher dimension. In some of the applications, it is suggested to use a more complex kernel to separate the classes that are curved or nonlinear. This transformation can lead to more accurate classifiers.
- **Regularization:** Regularization parameter in python's Scikit-learn C parameter used to maintain regularization. Here C is the penalty parameter, which represents misclassification or error term. The misclassification or error term tells the SVM optimization how much error is bearable. This is how you can control the trade-off between decision boundary and misclassification term. A smaller value of C creates a small-margin hyperplane and a larger value of C creates a larger-margin hyperplane.
- **Gamma:** A lower value of Gamma will loosely fit the training dataset, whereas a higher value of gamma will exactly fit the training dataset, which causes over-fitting. In other words, you can say a low value of gamma considers only nearby points in calculating the separation line, while the a value of gamma considers all the data points in the calculation of the separation line.

Advantages

SVM Classifiers offer good accuracy and perform faster prediction compared to Naïve Bayes algorithm. They also use less memory because they use a subset of training points in the decision phase. SVM works well with a clear margin of separation and with high dimensional space.

Disadvantages

SVM is not suitable for large datasets because of its high training time and it also takes more time in training compared to Naïve Bayes. It works poorly with overlapping classes and is also sensitive to the type of kernel used.

FAQ'S

1. What is meant by Support Vector?
2. What is meant by Hyperplane?
3. What is kernal?
4. What is gamma in SVM?
5. Discuss pros and cons of SVM?

Outcome:

With the completion of this assignment the students will be able to implement Support Vector Machine on given dataset.

Conclusion: