# MachineLearningAssignment4

Name: Sanjana Arkatala

Student ID: 700687517

#Question1:

Read_csv() is used to read the dataset and train_test_split() is used to split the data into training and testing sets.

Mean_squared_error is used to find the mean square error.

```
In [15]:    1  import pandas as pd
            2  import numpy as np
            3  import matplotlib.pyplot as plt
            4  from sklearn.model_selection import train_test_split
            5  from sklearn.linear_model import LinearRegression
            6  from sklearn import metrics
            7  from sklearn import preprocessing
            8  from sklearn.metrics import mean_squared_error
            9  from sklearn.cluster import KMeans
           10  from sklearn.impute import SimpleImputer
           11  from sklearn.decomposition import PCA
           12  from sklearn.preprocessing import LabelEncoder, StandardScaler
           13  import seaborn as sns
           14  sns.set(style="white", color_codes=True)
           15  import warnings
           16  warnings.filterwarnings("ignore")
```

```
In [14]:    1  df=pd.read_csv("./Salary_Data.csv")
            2  df.head()
```
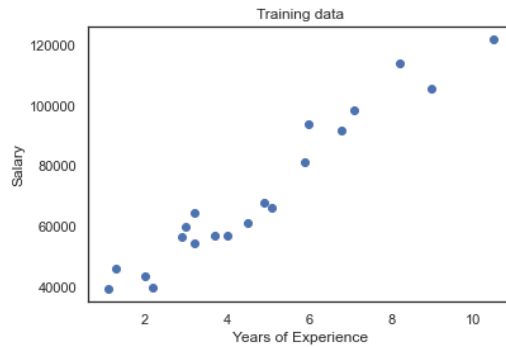
Out[14]:

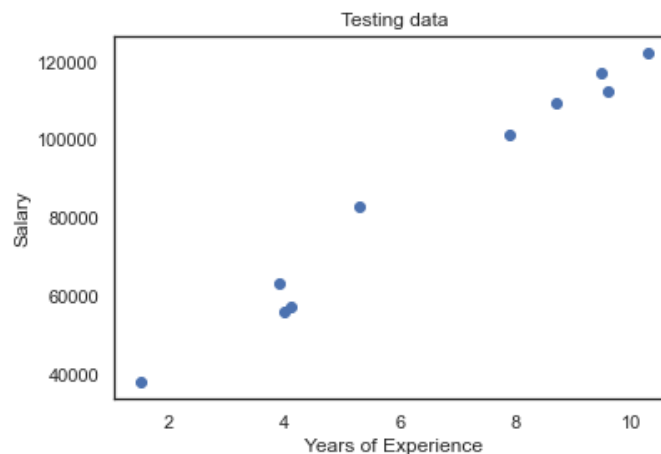|   | YearsExperience | Salary |
|---|---|---|
| 0 | 1.1 | 39343.0 |
| 1 | 1.3 | 46205.0 |
| 2 | 1.5 | 37731.0 |
| 3 | 2.0 | 43525.0 |
| 4 | 2.2 | 39891.0 |

```
In [16]:  1  X = df.iloc[:, :-1].values
          2  Y = df.iloc[:, 1].values
          3  X_Training, X_Testing, Y_Training, Y_Testing = train_test_split(X,Y, test_size=1/3,random_state = 0)
          4
          5  regressor = LinearRegression()
          6  regressor.fit(X_Training, Y_Training)
          7
          8  Y_Predict = regressor.predict(X_Testing)
          9
         10  mean_squared_error(Y_Testing,Y_Predict)
         11
```

Out[16]: 21026037.329511296

```
In [17]:  1  plt.title('Training data')
          2  plt.xlabel('Years of Experience')
          3  plt.ylabel('Salary')
          4  plt.scatter(X_Training, Y_Training)
          5  plt.show()
```



```
In [18]:  1  plt.title('Testing data')
          2  plt.xlabel('Years of Experience')
          3  plt.ylabel('Salary')
          4  plt.scatter(X_Testing, Y_Testing)
          5  plt.show()
```

#Question-2

All the missing values are replaced with mean by using a simple imputer and the model is fitted with the resulting data.

The elbow graph bends at 2 so we consider this as the number of clusters.

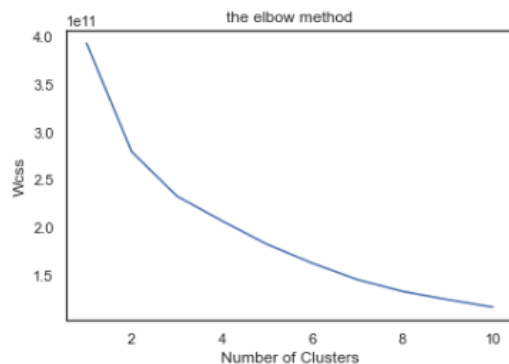The silhouette score is reduced after scaling

```
In [19]:   1  #question_2
```

```
In [22]:   1  df2=pd.read_csv("./K-Mean_Dataset.csv")
           2  df2.head()
```

Out[22]:

| | CUST_ID | BALANCE | BALANCE_FREQUENCY | PURCHASES | ONEOFF_PURCHASES | INSTALLMENTS_PURCHASES | CASH_ADVANCE | PURCHASES_FREQUEN |
|---|---|---|---|---|---|---|---|---|
| 0 | C10001 | 40.900749 | 0.818182 | 95.40 | 0.00 | 95.4 | 0.000000 | 0.166 |
| 1 | C10002 | 3202.467416 | 0.909091 | 0.00 | 0.00 | 0.0 | 6442.945483 | 0.000 |
| 2 | C10003 | 2495.148862 | 1.000000 | 773.17 | 773.17 | 0.0 | 0.000000 | 1.000 |
| 3 | C10004 | 1666.670542 | 0.636364 | 1499.00 | 1499.00 | 0.0 | 205.788017 | 0.083 |
| 4 | C10005 | 817.714335 | 1.000000 | 16.00 | 16.00 | 0.0 | 0.000000 | 0.083 |

```
In [34]:   1  X = df2.iloc[:,1:].values
           2
           3  imputer = SimpleImputer(missing_values=np.nan, strategy='mean')
           4  imputer = imputer.fit(X)
           5  x = imputer.transform(X)
```

```
In [35]:   1  wcss = []
           2  for i in range(1,11):
           3      kmeans = KMeans(n_clusters=i,init='k-means++',max_iter=300,n_init=10,random_state=0)
           4      kmeans.fit(x)
           5      wcss.append(kmeans.inertia_)
           6
           7  plt.plot(range(1,11),wcss)
           8  plt.title('the elbow method')
           9  plt.xlabel('Number of Clusters')
          10  plt.ylabel('Wcss')
          11  plt.show()
```



Type *Markdown* and LaTeX: $\alpha^2$

#Question3

Silhouette score is reduced after scaling

```
In [24]:   1  from sklearn.cluster import KMeans
           2  nclusters = 2
           3  km = KMeans(n_clusters=nclusters)
           4  km.fit(x)
```

Out[24]:  KMeans(n_clusters=2)

```
In [25]:   1  y_cluster_kmeans = km.predict(x)
           2  from sklearn import metrics
           3  score = metrics.silhouette_score(x, y_cluster_kmeans)
           4  print('Silhouette score:',score)
```

Silhouette score: 0.511639269641848

```
In [8]:    1  #Question_3
```

```
In [26]:   1  scaler = preprocessing.StandardScaler()
           2  scaler.fit(x)
           3  X_scaled_array = scaler.transform(x)
           4  X_scaled = pd.DataFrame(X_scaled_array)
```

```
In [27]:   1  from sklearn.cluster import KMeans
           2  nclusters = 2
           3  km = KMeans(n_clusters=nclusters)
           4  km.fit(X_scaled)
```

Out[27]:  KMeans(n_clusters=2)

```
In [28]:   1  y_scaled_cluster_kmeans = km.predict(X_scaled)
           2  from sklearn import metrics
           3  score = metrics.silhouette_score(X_scaled, y_scaled_cluster_kmeans)
           4  print('Silhouette score after applying scaling:',score)
```

Silhouette score after applying scaling: 0.20946870767221923

Github link: https://github.com/Sanjana9791/MachineLearningAssignment4.git