

Part B Assignment 4

K-means clustering on the Iris dataset

Name : Ayan Gadpal

Roll No. 43308

K-Means clustering is one of the simplest and most commonly used clustering algorithms. It tries to find cluster centers that are representative of certain regions of the data. The algorithm alternates between two steps: assigning each data point to the closest cluster center, and then setting each cluster center as the mean of the data points that are assigned to it. The algorithm is finished when the assignment of instances to clusters no longer changes.

K-Means has the advantage that it's pretty fast, as all we're really doing is computing the distances between points and group centers. It thus has a linear complexity $O(n)$

▼ Step 1 : Load Dataset

We will first load the dataset and split it

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from sklearn import datasets

1 iris=datasets.load_iris()
2 dataset = pd.DataFrame(data= np.c_[iris['data'], iris['target']],
3                           columns= iris['feature_names'] + ['target'])
4 dataset.head()
```

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	target
0	5.1	3.5	1.4	0.2	0.0
1	4.9	3.0	1.4	0.2	0.0
2	4.7	3.2	1.3	0.2	0.0
3	4.6	3.1	1.5	0.2	0.0
4	5.0	3.6	1.4	0.2	0.0

```
1 df=pd.DataFrame({'x':iris.data[:,0],
2                  'y':iris.data[:,1],
3                  'cluster':iris.target})
4
```

▼ Step 2 : Model Building

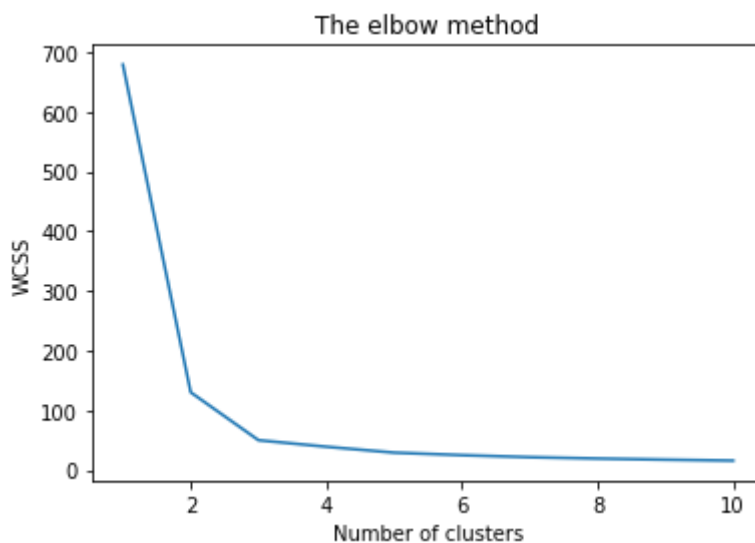
The task is to decide how many cluster do we need to get best modeling of the data

We Will use **Elbow method** to select appropriate number of clusters This method looks at the **percentage of variance** explained as a function of the number of clusters

```

1 #Finding the optimum number of clusters for k-means classification
2 x = dataset.iloc[:, [1, 2, 3, 4]].values
3 from sklearn.cluster import KMeans
4 wcss = []
5
6 for i in range(1, 11):
7     kmeans = KMeans(n_clusters = i, init = 'k-means++', max_iter = 300, n_init = 10, random_state = 0)
8     kmeans.fit(x)
9     wcss.append(kmeans.inertia_)
10
11 #Plotting the results onto a line graph, allowing us to observe 'The elbow'
12 plt.plot(range(1, 11), wcss)
13 plt.title('The elbow method')
14 plt.xlabel('Number of clusters')
15 plt.ylabel('WCSS') #within cluster sum of squares
16 plt.show()

```



As we can see the "Elbow" is formed on 3 cluster (X-axis), We will use 3 as our total number of clusters

▼ Step 3 : Apply K-mean

```

1 centroids={}
2 for i in range(3):
3     res=[]
4     res.append(df.loc[df["cluster"]==i]["x"].mean())
5     res.append(df.loc[df["cluster"]==i]["y"].mean())
6     centroids[i]=res

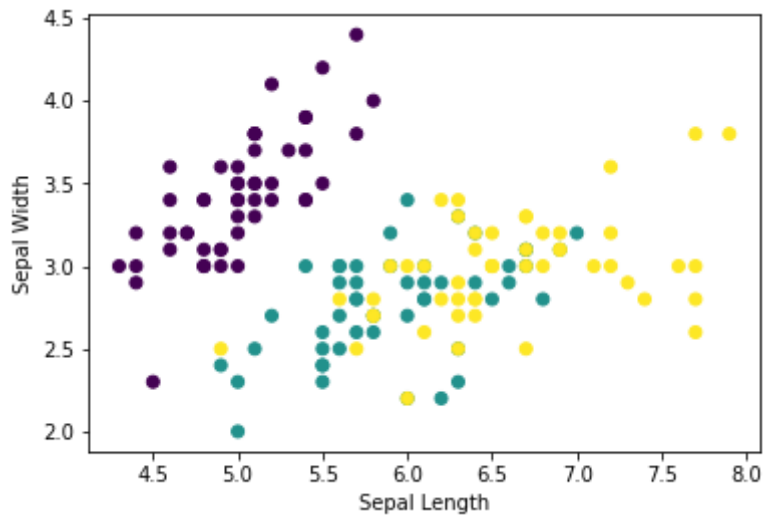
```

1 centroids

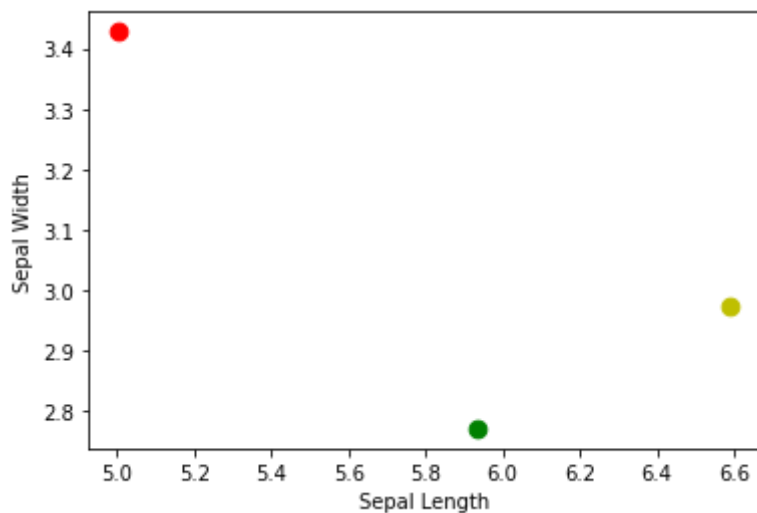
```
{0: [5.005999999999999, 3.428000000000001],
1: [5.936, 2.7700000000000005],
2: [6.587999999999998, 2.973999999999998]}
```

▼ Visualize

```
1 # fig=plt.figure(figsize=(5,5))
2 plt.scatter(df["x"],df["y"],c=iris.target)
3 plt.xlabel("Sepal Length")
4 plt.ylabel("Sepal Width")
5 plt.show()
```



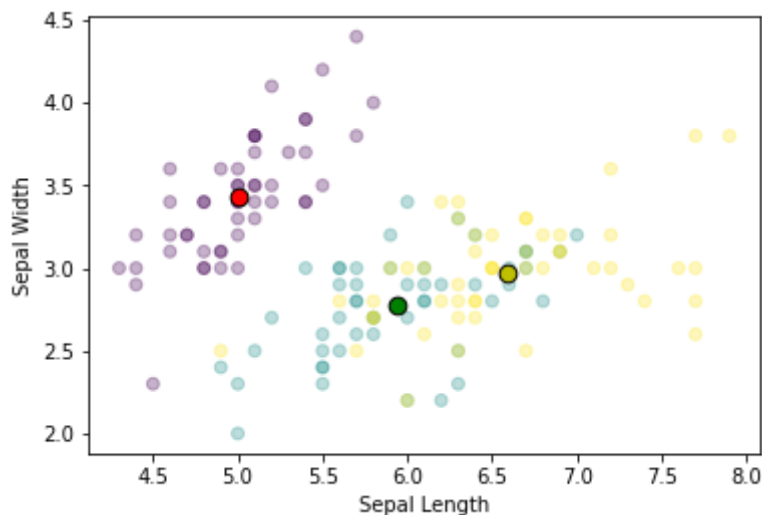
```
1 centers={0:'r',1:'g',2:'y'}
2 for i in range(3):
3     plt.scatter(centroids[i][0],centroids[i][1],s=75,color=centers[i])
4 plt.xlabel("Sepal Length")
5 plt.ylabel("Sepal Width")
6 plt.show()
```



```

1 plt.scatter(df["x"],df["y"],c=iris.target,alpha=0.3)
2 for i in range(3):
3     plt.scatter(centroids[i][0],centroids[i][1],color=centers[i],s=75,edgecolors="black"
4 plt.xlabel("Sepal Length")
5 plt.ylabel("Sepal Width")
6 plt.show()
7

```



▼ Calculate the closest points to the centroids and assigning the new clusters

```

1 def assignment(df, centroids):
2     for i in range(3):
3         # sqrt((x1 - x2)^2 + (y1 - y2)^2)
4         df['distance_from_{}'.format(i)] = (np.sqrt((df['x'] - centroids[i][0]) ** 2 +
5             centroid_distance_cols = ['distance_from_{}'.format(i) for i in centroids.keys()])
6         df['closest'] = df.loc[:, centroid_distance_cols].idxmin(axis=1)
7         df['closest'] = df['closest'].map(lambda x: int(x.lstrip('distance_from_')))
8         df['color'] = df['closest'].map(lambda x: centers[x])
9     return df

```

```

1 df=assignment(df,centroids)
2 df.head()

```

	x	y	cluster	distance_from_0	distance_from_1	distance_from_2	closest	col
0	5.1	3.5	0	0.118406	1.051594	1.764836	0	
1	4.9	3.0	0	0.440931	0.926140	1.914215	0	
2	4.7	3.2	0	0.381602	1.187514	2.116492	0	
3	4.6	3.1	0	0.521939	1.242335	2.212913	0	
4	5.0	3.6	0	0.172105	1.192508	1.887407	0	

▼ Calculate new Centroids

```

1 def update():
2     for i in range(3):
3         centroids[i][0] = df.loc[df["closest"]==i]["x"].mean()
4         centroids[i][1] = df.loc[df["closest"]==i]["y"].mean()

```

```

1 update()
2 centroids

```

```

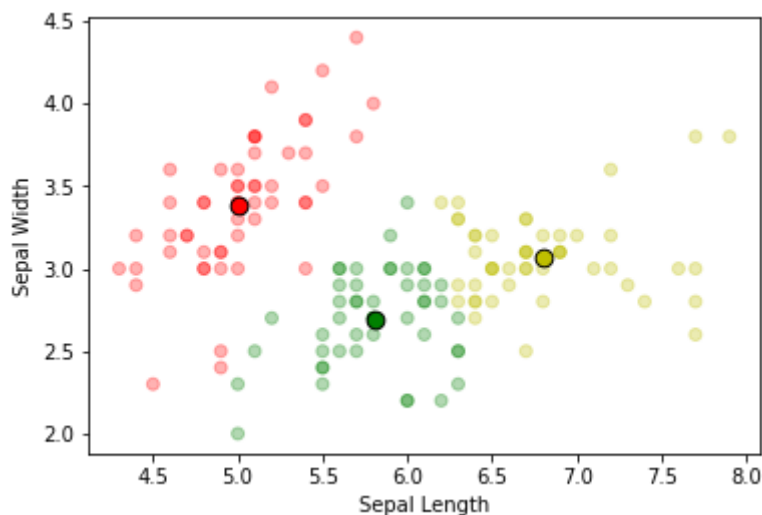
{0: [5.00943396226415, 3.383018867924529],
1: [5.806122448979591, 2.693877551020408],
2: [6.802083333333331, 3.0687499999999996]}

```

```

1 plt.scatter(df['x'],df['y'],color=df["color"],alpha=0.3)
2 for i in centroids.keys():
3     plt.scatter(centroids[i][0],centroids[i][1],color=centers[i],s=75,edgecolors="k")
4 plt.xlabel("Sepal Length")
5 plt.ylabel("Sepal Width")
6 plt.show()

```



▼ Continue until all assigned clusters do not change

```

1 while True:
2     closest_centroids = df['closest'].copy(deep=True)
3     update()
4     df = assignment(df, centroids)
5     if closest_centroids.equals(df['closest']):
6         break

```

▼ Step 4 : Visualize the result

```

1 plt.scatter(df['x'],df['y'],color=df["color"],alpha=0.3)
2 for i in centroids.keys():
3     plt.scatter(centroids[i][0],centroids[i][1],color=centers[i],s=75,edgecolors="k")

```

```
4 plt.xlabel("Sepal Length")  
5 plt.ylabel("Sepal Width")  
6 plt.show()
```

