**PROJECT 3 REPORT**
**SHREYAS MOHAN  1001669806**

## Project Description

- We use the Iris Dataset. This data consists of 50 samples from each of three species of Iris (Iris setosa, Iris virginica and Iris versicolor). It has four features from each sample: length and width of sepals and petals.
- We use K Means algorithm to cluster the data

    **K means works through the following iterative process:**

- Pick a value for k (the number of clusters to create).

    We use **Elbow Method** to determine the value of k and choose k as 3 as it is optimum. Also k=3, as we have **3 classes**.

- Initialize k 'centroids' (starting points) in your data

    We initialize k=3

- Creating clusters. Assign each point to the nearest centroid.

- Making clusters better. Move each centroid to the center of its cluster.

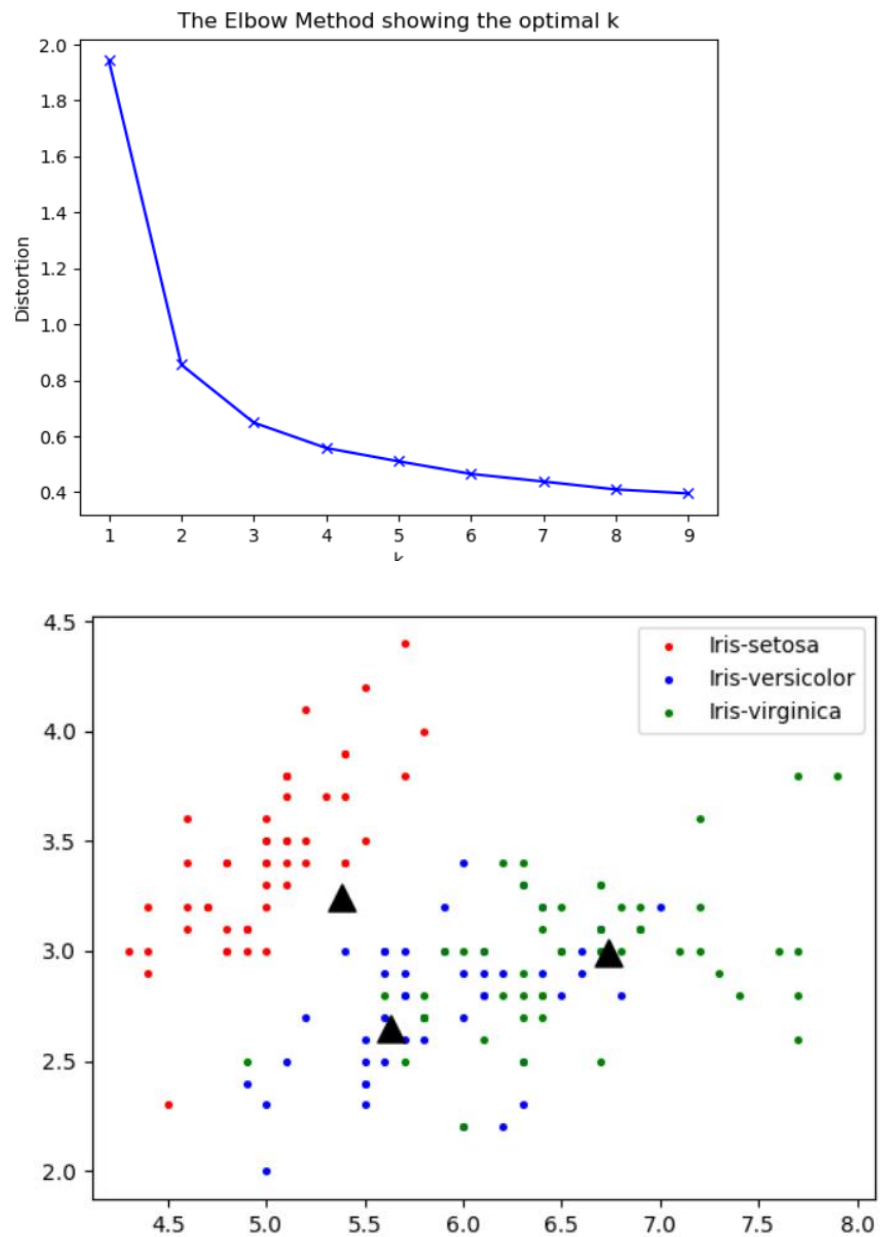- Repeat the above steps until your centroids converge.

## Structure of Code:

- First, the iris data file is copied to an iris csv file and headers 'SepalLength', 'SepalWidth', 'PetalLength', 'PetalWidth', and 'Species' are added to the data points for accessing the data points.
- Features of the iris are stored in a list while categories are stored in a separate list.
- We use elbow_method() to determine the value of k and as there are only three classes, we initialize k to 3.
- We find the centroids of the cluster using the mean and standard deviation formula which includes error.
- We remove the error by normalizing the distance between each centroid and data points and assign it to the closest centroid and repeating this until the error comes to zero.
- Output displays the number of data points that are assigned to an incorrect cluster and also displays the centroids.
- Plotting the data points and centroids in the graph and saving it in an external file 'KMeans.png'

**Screenshot of the output:**

**You can run on command line or any Python IDE.**

**Refer readme file for execution instructions.**



The Elbow Method showing the optimal k

File   Edit   Search   Source   Run   Debug   Consoles   Projects   Tools   View   Help

C:\Users\Shreyas Mohan\Documents\Fall 19\Machine Learning\KMeans_algorithm-master

Editor - C:\Users\Shreyas Mohan\Documents\Fall 19\Machine Learning\KMeans_algorithm-master\Project3.py

Project3.py

```python
70
71 # When, after an update, the estimate of that center stays the same, exit loop
72 while error != 0:
73     # Measure the distance to every center
74     for i in range(k):
75         distances[:,i] = np.linalg.norm(data - centers[i], axis=1)
76     # Assign all training data to closest center
77     clusters = np.argmin(distances, axis = 1)
78
79     centers_old = deepcopy(centers_new)
80     # Calculate mean for every cluster and update the center
81     for i in range(k):
82         centers_new[i] = np.mean(data[clusters == i], axis=0)
83     error = np.linalg.norm(centers_new - centers_old)
84 incorrect = 0
85 for i, j in zip(category,clusters):
86     if i != j:
87         incorrect +=1
88 print("Number of data points incorrectly clustered")
89 print(incorrect)
90 # Plot the data and the centers generated as random
91 colors=['red', 'blue', 'green']
92 for i in range(n):
93     if colors[int(category[i])] == 'red':
94         l1 = plt.scatter(data[i, 0], data[i,1], s=7, color = colors[int(category[i])], label = 'Iris-setosa')
95     elif colors[int(category[i])] == 'blue':
96         l2 = plt.scatter(data[i, 0], data[i,1], s=7, color = colors[int(category[i])], label = 'Iris-versicolor'
97     elif colors[int(category[i])] == 'green':
98         l3 = plt.scatter(data[i, 0], data[i,1], s=7, color = colors[int(category[i])], label = 'Iris-virginica'
99 handles, labels = plt.gca().get_legend_handles_labels()
100 handle_list, label_list = [], []
101 for handle, label in zip(handles, labels):
102     if label not in label_list:
103         handle_list.append(handle)
104         label_list.append(label)
105 plt.legend(handle_list, label_list)
106 plt.scatter(centers_new[:,0], centers_new[:,1], marker='^', c='black', s=150, label = 'Centroids')
107 plt.savefig('KMeans')
108 plt.show()
109 #To print the centroids
110 print(centers_new)
```
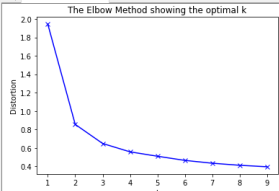
Variable explorer

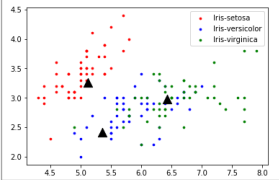| Name | Type | Size | Value |
|---|---|---|---|
| c | int | 1 | 4 |
| category | float64 | (150,) | [0. 0. 0. ... 2. 2. 2.] |
| centers | float64 | (3, 4) | [[4.290248   2.36012406 4.36605347 1.08625329]<br>[5.63222992 3.46556345 ... |
| centers_new | float64 | (3, 4) | [[5.36       2.41       4.15       1.28     ]<br>[6.4308642  2.98024691 ... |

Variable explorer   Help

IPython console

Console 1/A



Number of data points incorrectly clustered
116



IPython console   History log

Permissions: RW   End-of-lines: CRLF   Encoding: ASCII   Line: 110   Column: 1   Memory: 69 %