

CSBB 311 : MACHINE LEARNING

LAB ASSIGNMENT 7 : Implementation of K-Means Clustering From Scratch and using Inbuilt Library

Submitted By:

Name: KARTIK MITTAL

Roll No: 221210056

Branch: CSE

Semester: 5th Sem

Group: 2

Submitted To: Dr. Preeti Verma

Department of Computer Science and Engineering

NATIONAL INSTITUTE OF TECHNOLOGY DELHI



2024

Code (Using Inbuilt Library) -

```
1  import numpy as np
2  import pandas as pd
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5  from sklearn.cluster import KMeans
6  from sklearn.metrics import confusion_matrix
7  from sklearn.preprocessing import StandardScaler
8
9  # Load the data
10 df = pd.read_csv('Mall_Customers.csv')
11
12 # Delete rows with missing values
13 df.dropna(inplace=True)
14
15 # Selecting features for clustering
16 X = df[['Annual Income (k$)', 'Spending Score (1-100)']].values
17
18 # Step 1: Initial scatter plot of the data points
19 plt.figure(figsize=(10, 6))
20 plt.scatter(X[:, 0], X[:, 1], color='blue', s=100)
21 plt.title('Initial Scatter Plot of Data Points')
22 plt.xlabel('Annual Income (k$)')
23 plt.ylabel('Spending Score (1-100)')
24 plt.show()
25
26 # Step 2: Standardizing the features
27 scaler = StandardScaler()
28 X_scaled = scaler.fit_transform(X)
```

```

30 # Step 3: Elbow curve to find the optimal number of clusters
31 inertia = []
32 K = range(1, 11)
33 for k in K:
34     kmeans = KMeans(n_clusters=k, random_state=42)
35     kmeans.fit(X_scaled)
36     inertia.append(kmeans.inertia_)
37
38 # Plot the elbow curve
39 plt.figure(figsize=(8, 6))
40 plt.plot(K, inertia, marker='o', linestyle='--')
41 plt.title('Elbow Curve')
42 plt.xlabel('Number of clusters')
43 plt.ylabel('Inertia (Sum of Squared Distances)')
44 plt.show()
45
46 # Step 4: Applying K-Means with the optimal number of clusters (e.g., 5 from elbow curve)
47 optimal_clusters = 5
48 kmeans = KMeans(n_clusters=optimal_clusters, random_state=42)
49 kmeans.fit(X_scaled)
50
51 # Step 5: Visualizing the clusters with centroids
52 plt.figure(figsize=(10, 6))
53 plt.scatter(X_scaled[:, 0], X_scaled[:, 1], c=kmeans.labels_, cmap='viridis', s=100)
54 plt.scatter(kmeans.cluster_centers[:, 0], kmeans.cluster_centers[:, 1], color='red',
55             marker='*', s=200, label='Centroids')
56 plt.title('K-Means Clustering with Centroids (Scikit-Learn)')
57 plt.xlabel('Annual Income (Standardized)')
58 plt.ylabel('Spending Score (Standardized)')

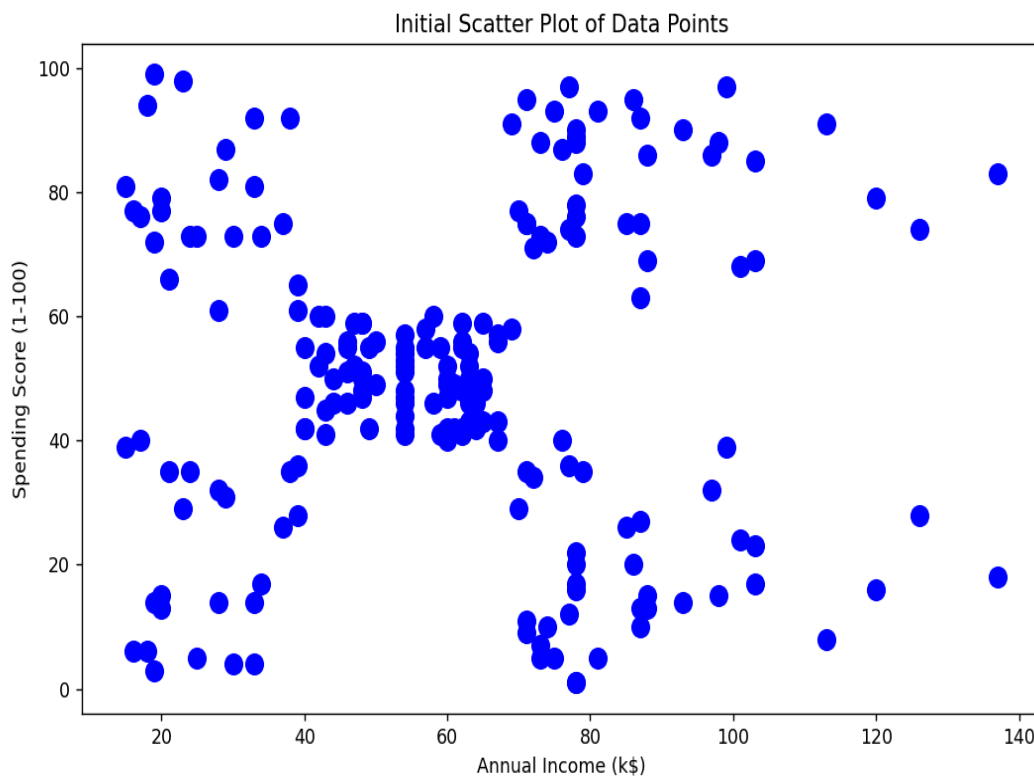
```

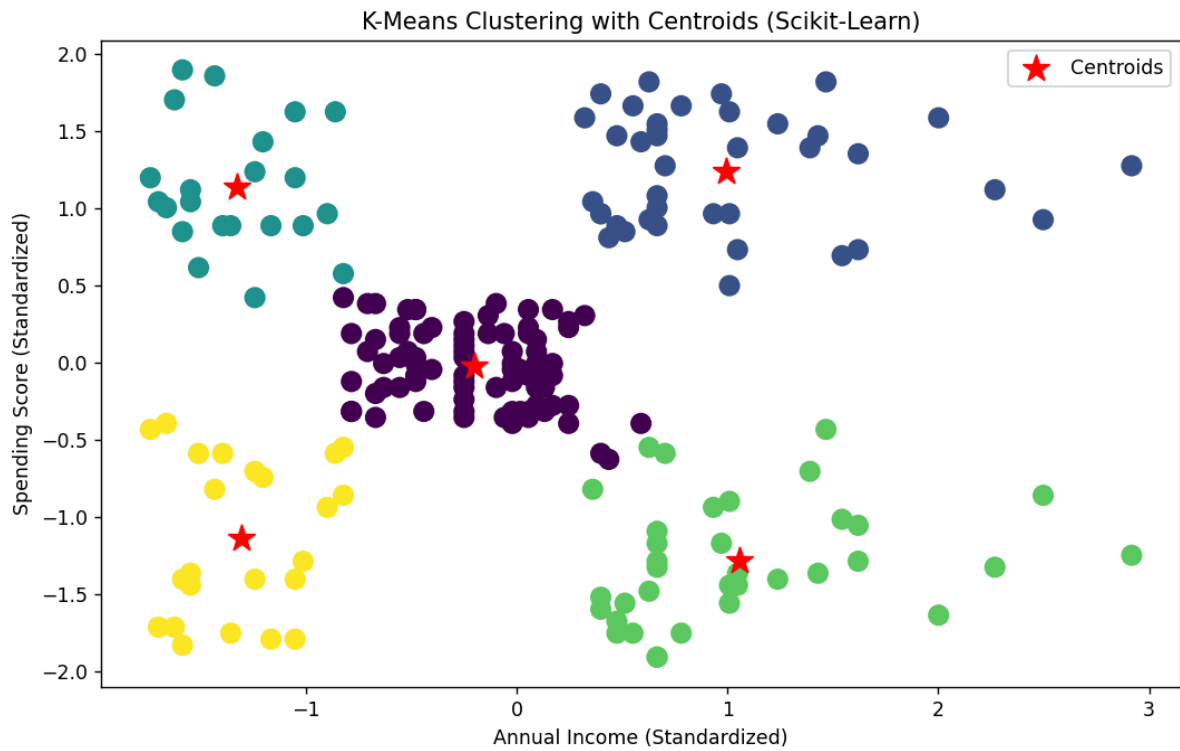
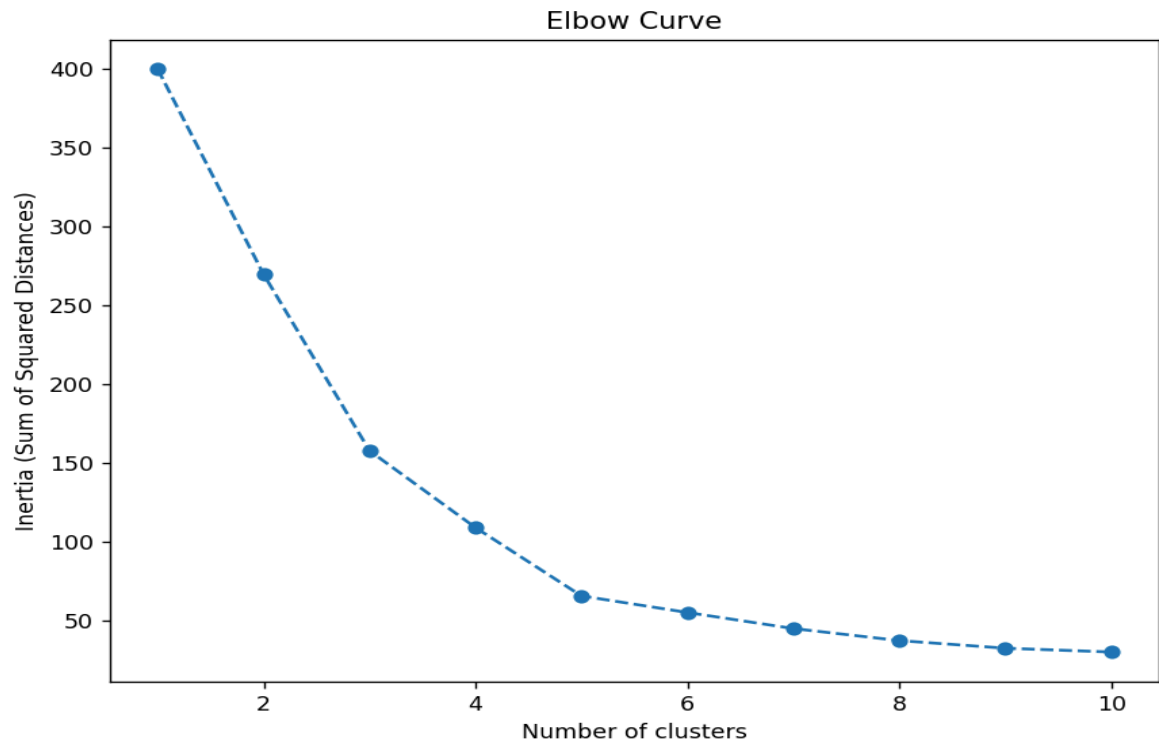
```

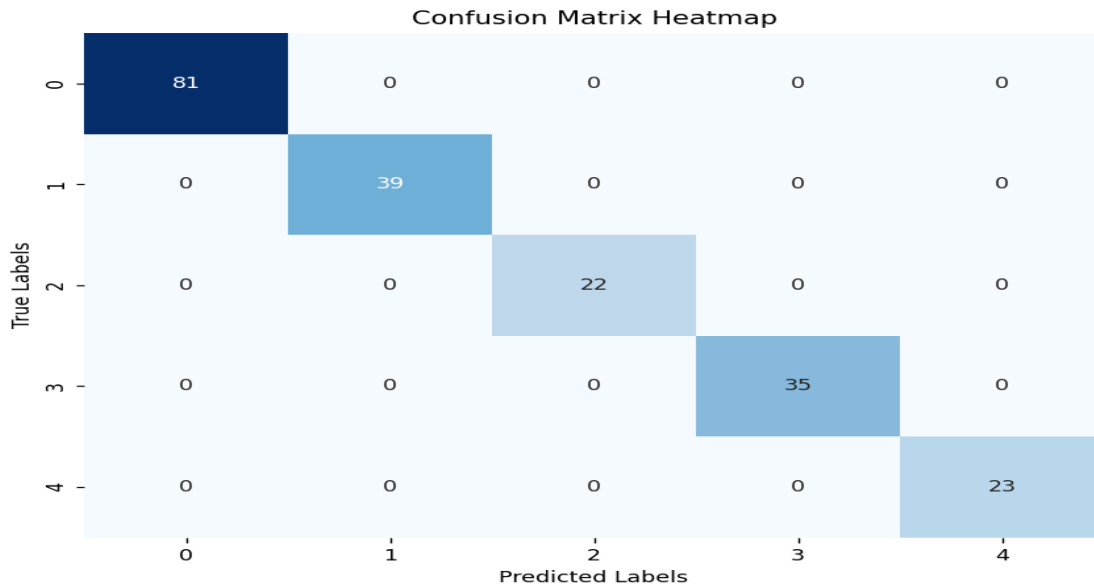
59 plt.legend()
60 plt.show()
61
62 # Step 6: Confusion matrix heatmap (assuming cluster predictions vs true labels)
63 # In this case, we don't have true labels, but we can use the predicted labels for simulation
64 y_true = kmeans.labels_ # In practice, this should be true labels if available
65 y_pred = kmeans.predict(X_scaled)
66
67 # Create confusion matrix
68 conf_matrix = confusion_matrix(y_true, y_pred)
69
70 # Plotting the confusion matrix heatmap
71 plt.figure(figsize=(8, 6))
72 sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues', cbar=False)
73 plt.title('Confusion Matrix Heatmap')
74 plt.xlabel('Predicted Labels')
75 plt.ylabel('True Labels')
76 plt.show()
77
78 # Output cluster centers
79 print("Cluster centers (standardized):")
80 print(kmeans.cluster_centers_)

```

Output -







Code (From Scratch) -

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.metrics import confusion_matrix
6 from sklearn.preprocessing import StandardScaler
7
8 # Load the data
9 df = pd.read_csv('Mall_Customers.csv')
10
11 # Delete rows with missing values
12 df.dropna(inplace=True)
13
14 # Selecting features for clustering
15 X = df[['Annual Income (k$)', 'Spending Score (1-100)']].values
16
17 # Step 1: Initial scatter plot of the data points
18 plt.figure(figsize=(10, 6))
19 plt.scatter(X[:, 0], X[:, 1], color='blue', s=100)
20 plt.title('Initial Scatter Plot of Data Points')
21 plt.xlabel('Annual Income (k$)')
22 plt.ylabel('Spending Score (1-100)')
23 plt.show()
24
25 # Standardizing the features
26 scaler = StandardScaler()
27 X_scaled = scaler.fit_transform(X)
```

```

29 # Function to calculate the Euclidean distance between points
30 def euclidean_distance(a, b):
31     return np.sqrt(np.sum((a - b) ** 2))
32
33 # K-Means Clustering from scratch
34 class KMeansScratch:
35     def __init__(self, n_clusters, max_iter=100):
36         self.n_clusters = n_clusters
37         self.max_iter = max_iter
38
39     def fit(self, X):
40         # Randomly initialize cluster centers
41         np.random.seed(42)
42         random_idx = np.random.permutation(X.shape[0])
43         self.centroids = X[random_idx[:self.n_clusters]]
44
45         for i in range(self.max_iter):
46             # Assign clusters
47             self.labels = self.assign_clusters(X)
48
49             # Compute new centroids
50             new_centroids = self.calculate_centroids(X)
51
52             # If centroids do not change, break
53             if np.all(self.centroids == new_centroids):
54                 break
55             self.centroids = new_centroids
56

```

```

def assign_clusters(self, X):
    # Assign each point to the nearest centroid
    labels = []
    for point in X:
        distances = [euclidean_distance(point, centroid) for centroid in self.centroids]
        # np.argmin(distances) returns the index of the minimum value in the array
        labels.append(np.argmin(distances))
    return np.array(labels)

def calculate_centroids(self, X):
    # Compute the centroids as the mean of the points assigned to each cluster
    centroids = np.zeros((self.n_clusters, X.shape[1]))
    for idx in range(self.n_clusters):
        # Get all points that belong to cluster idx
        points = X[self.labels == idx]
        centroids[idx] = np.mean(points, axis=0) if len(points) > 0 else self.centroids[idx]
    return centroids

def predict(self, X):
    return self.assign_clusters(X)

```

Step 2: Plot the elbow curve to find optimal number of clusters

```

inertia = []
for k in range(1, 11):
    kmeans = KMeansScratch(n_clusters=k)
    kmeans.fit(X_scaled)

```



```

83     # Inertia is the sum of squared distances to the nearest centroid
84     inertia_val = np.sum([euclidean_distance(X_scaled[i],
85         kmeans.centroids[kmeans.labels[i]])**2 for i in range(X_scaled.shape[0])])
86     inertia.append(inertia_val)
87
88 # Plot the elbow curve
89 plt.figure(figsize=(8, 6))
90 plt.plot(range(1, 11), inertia, marker='o', linestyle='--')
91 plt.title('Elbow Curve')
92 plt.xlabel('Number of clusters')
93 plt.ylabel('Inertia (Sum of Squared Distances)')
94 plt.show()
95
96 # Step 3: Applying K-Means with the optimal number of clusters (e.g., 5 from elbow curve)
97 kmeans_scratch = KMeansScratch(n_clusters=5)
98 kmeans_scratch.fit(X_scaled)
99
100 # Step 4: Visualizing the clusters with centroids
101 plt.figure(figsize=(10, 6))
102 plt.scatter(X_scaled[:, 0], X_scaled[:, 1], c=kmeans_scratch.labels, cmap='viridis', s=100)
103 plt.scatter(kmeans_scratch.centroids[:, 0], kmeans_scratch.centroids[:, 1],
104     color='red', marker='*', s=200, label='Centroids')
105 plt.title('K-Means Clustering with Centroids (Scratch)')
106 plt.xlabel('Annual Income (Standardized)')
107 plt.ylabel('Spending Score (Standardized)')
108 plt.legend()
109 plt.show()

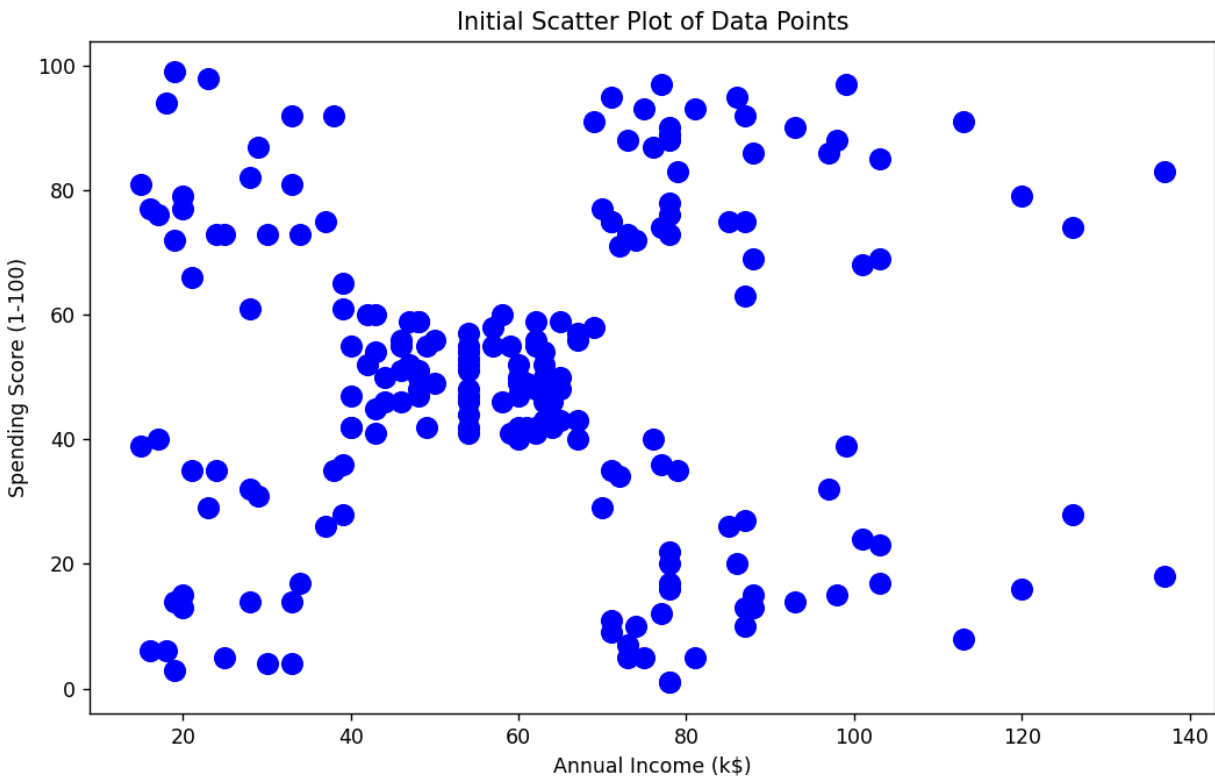
```

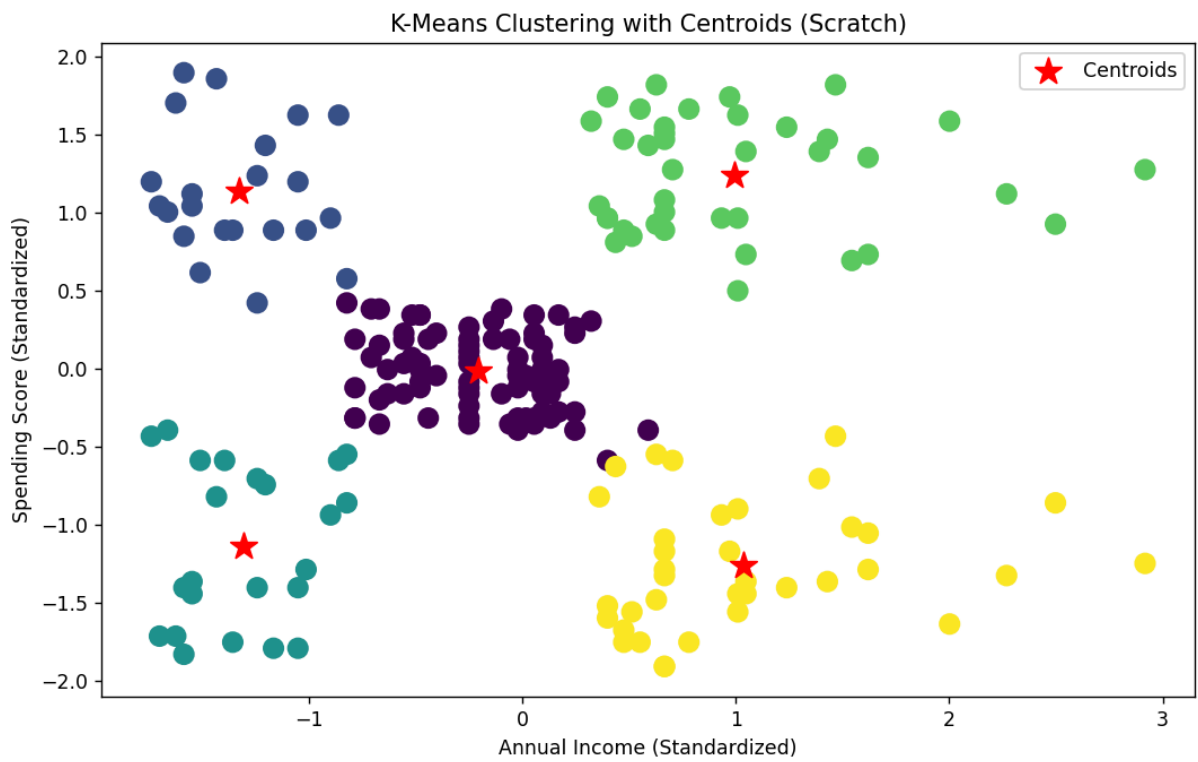
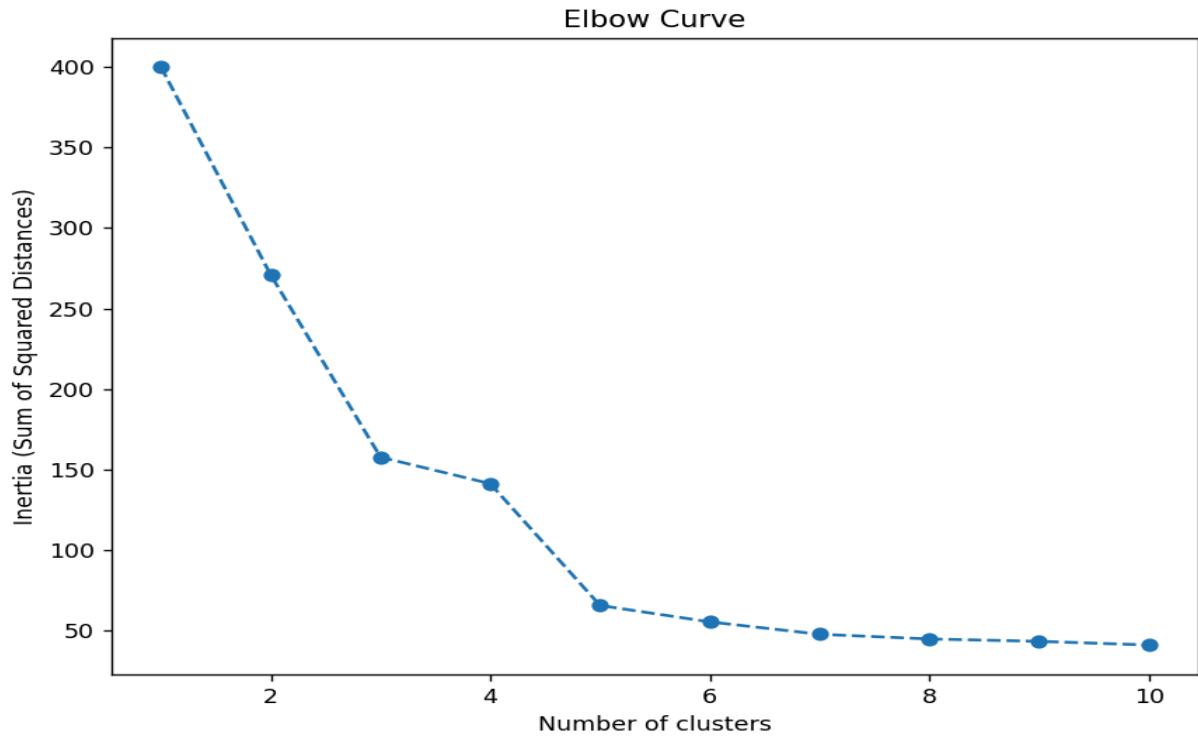
```

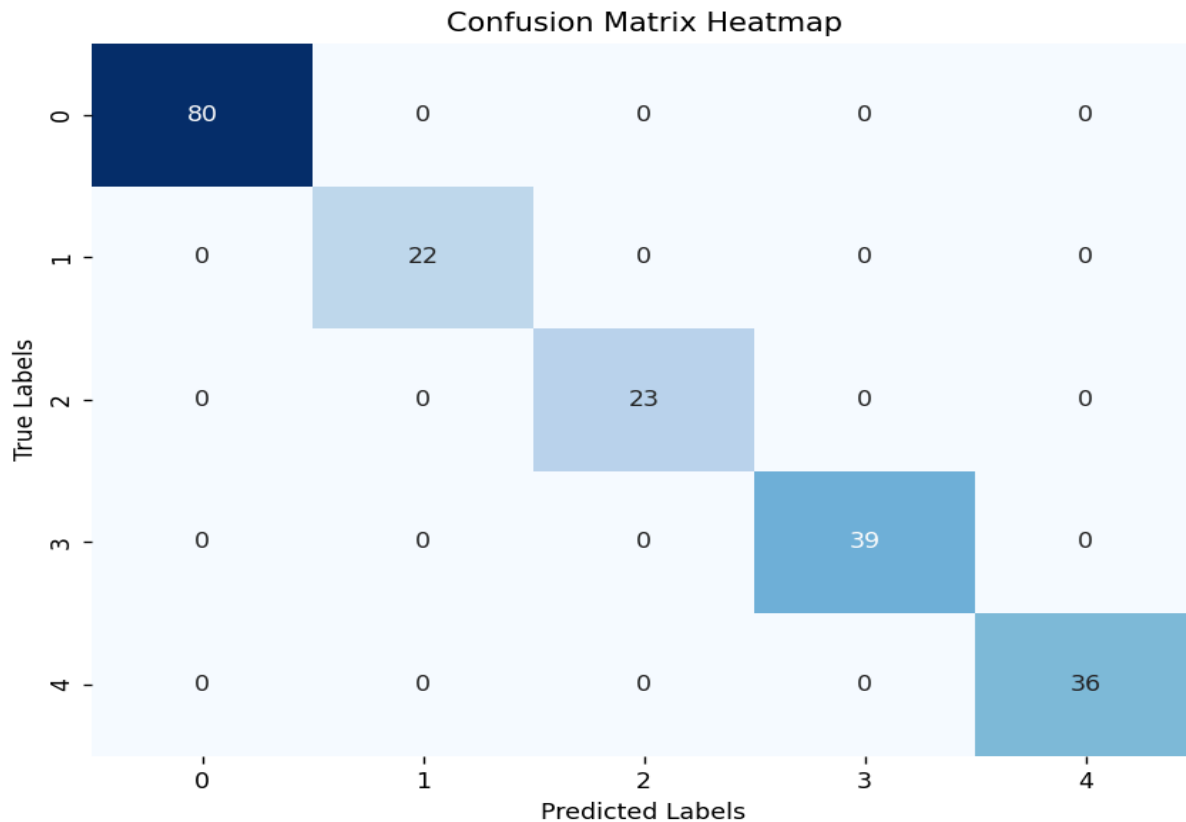
111 # Step 5: Confusion matrix heatmap (assuming cluster predictions vs true labels)
112 # Here, since we don't have true labels, we simulate a confusion matrix using the predicted labels
113 y_true = kmeans_scratch.labels # Predicted labels (in real scenario, compare with actual labels)
114 y_pred = kmeans_scratch.predict(X_scaled)
115
116 conf_matrix = confusion_matrix(y_true, y_pred)
117
118 # Plotting the confusion matrix heatmap
119 plt.figure(figsize=(8, 6))
120 sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues', cbar=False)
121 plt.title('Confusion Matrix Heatmap')
122 plt.xlabel('Predicted Labels')
123 plt.ylabel('True Labels')
124 plt.show()
125
126 # Output cluster centers
127 print("Cluster centers (standardized):")
128 print(kmeans_scratch.centroids)

```

Output -







Conclusion -

- The inbuilt K-Means (e.g., Scikit-learn) is more reliable and accurate because it has optimized algorithms and edge-case handling, ensuring better predictions, especially with large dataset.
- So inbuilt has better prediction results.