

Mustafa Aydoğan

191101002

In [1]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.metrics import confusion_matrix, classification_report, roc_curve, auc
4 from mpl_toolkits.mplot3d import Axes3D
```

In [2]:

```
1 # Import the necessary library for working with data
2 import pandas as pd
3
4 # Load the dataset from the CSV file named 'iris.csv'
5 data = pd.read_csv('iris.csv')
6
7 # Print the first few rows of the dataset to inspect its structure
8 print(data.head())
```

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species	
es	0	1	5.1	3.5	1.4	0.2	Iris-seto
sa	1	2	4.9	3.0	1.4	0.2	Iris-seto
sa	2	3	4.7	3.2	1.3	0.2	Iris-seto
sa	3	4	4.6	3.1	1.5	0.2	Iris-seto
sa	4	5	5.0	3.6	1.4	0.2	Iris-seto

In [3]:

```
1 # Convert species names to integer labels
2 data['Species'] = data['Species'].map({'Iris-setosa': 0, 'Iris-versicolor': 1, 'Iris-
3
4 # Remove the 'Id' column as it's not needed
5 data.drop(['Id'], inplace=True, axis=1)
6
7 # Print the first few rows of the modified dataset to verify the changes
8 print(data.head())
```

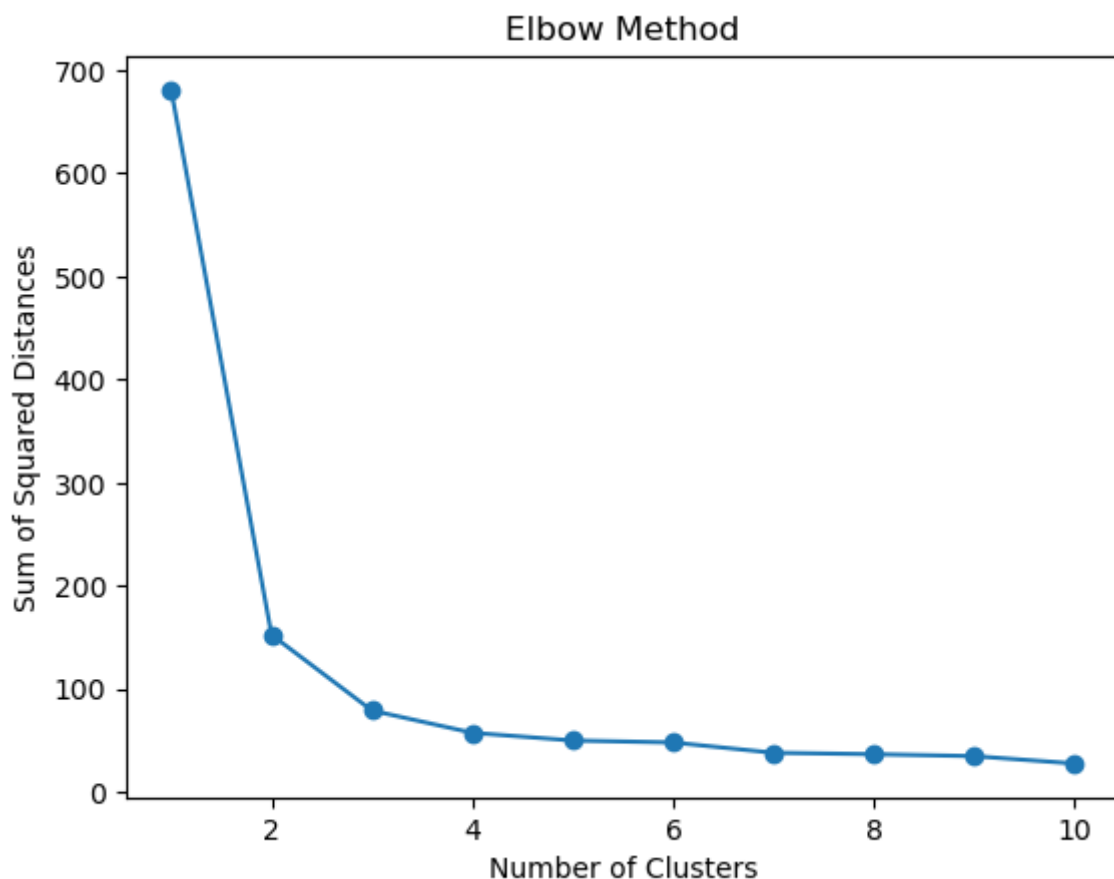
	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
0	5.1	3.5	1.4	0.2	0
1	4.9	3.0	1.4	0.2	0
2	4.7	3.2	1.3	0.2	0
3	4.6	3.1	1.5	0.2	0
4	5.0	3.6	1.4	0.2	0

In [4]:

```
1 # Extract features (X) and labels (y) from the dataset
2 X = data.iloc[:, :-1].values
3 y = data.iloc[:, -1].values
4
5 # Split the data into training and testing sets
6 train_size = int(0.8 * len(X)) # 80% of the data for training
7 X_train, X_test = X[:train_size], X[train_size:] # Split features into training and
8 y_train, y_test = y[:train_size], y[train_size:] # Split labels into training and te
9
10 # Split the entire dataset into training and testing sets (including all columns)
11 train_data, test_data = data[:train_size], data[train_size:]
```

In [5]:

```
1 import matplotlib.pyplot as plt
2 from kmeans import KMeansClusterClassifier # Import the KMeansClusterClassifier from
3
4 # Function to calculate the total squared distance of points from their cluster center
5 def calculate_total_squared_distance(X, kmeans):
6     total_distance = 0
7     for x in X:
8         centroid = kmeans.centroids[kmeans.predict([x])[0]]
9         total_distance += kmeans._euclidean_distance(x, centroid) ** 2
10    return total_distance
11
12 # Function to plot the elbow method graph
13 def plot_elbow_method(X, max_clusters):
14     distances = []
15     for n_clusters in range(1, max_clusters + 1):
16         kmeans = KMeansClusterClassifier(n_clusters)
17         kmeans.fit(X)
18         total_distance = calculate_total_squared_distance(X, kmeans)
19         distances.append(total_distance)
20
21     # Plotting the elbow method graph
22     plt.plot(range(1, max_clusters + 1), distances, marker='o')
23     plt.xlabel('Number of Clusters')
24     plt.ylabel('Sum of Squared Distances')
25     plt.title('Elbow Method')
26     plt.show()
27
28 max_clusters = 10 # Maximum number of clusters to try
29 plot_elbow_method(X, max_clusters) # Call the function to create the elbow plot
30
```



In [6]:

```
1 #optimal_cluester is 3 according to elbow
2 optimal_clusters = 3
```

In [7]:

```
1 # Re-training with the optimal_clusters
2 kmeans = KMeansClusterClassifier(optimal_clusters)
3 kmeans.fit(X)
```

In [8]:

```
1 # Print the final centroids of each cluster
2 print("Final centroids:")
3 for i, centroid in enumerate(kmeans.centroids):
4     print("Cluster " + str(i + 1) + ": " + str(centroid))
```

Final centroids:

Cluster 1: [6.8500000000000005, 3.073684210526315, 5.742105263157893, 2.0710526315789473]

Cluster 2: [5.901612903225807, 2.748387096774194, 4.393548387096775, 1.4338709677419357]

Cluster 3: [5.005999999999999, 3.4180000000000006, 1.464, 0.2439999999999999]

In [9]:

```
1 # Get predictions from KMeans
2 predictions = kmeans.predict(X)
3
4 # Separate predictions for the first 50, next 50, and last 50 samples
5 first_class_predictions = predictions[:50]
6 second_class_predictions = predictions[50:100]
7 third_class_predictions = predictions[100:]
8
9 # Find the most common cluster prediction for each class
10 from collections import Counter
11 first_class_common_cluster = Counter(first_class_predictions).most_common(1)[0][0]
12 second_class_common_cluster = Counter(second_class_predictions).most_common(1)[0][0]
13 third_class_common_cluster = Counter(third_class_predictions).most_common(1)[0][0]
14
15 # Print the most common cluster prediction for the second class
16 print("Most common cluster prediction for the second class:", second_class_common_cluster)
17
18 # Create a mapping from cluster to class
19 cluster_to_class_mapping = {
20     first_class_common_cluster: 1,
21     second_class_common_cluster: 2,
22     third_class_common_cluster: 3
23 }
24
25 # Transform cluster predictions to class labels
26 predicted_labels = [cluster_to_class_mapping[cluster] for cluster in predictions]
27
28 # Update predictions with class labels
29 predictions = predicted_labels
```

Most common cluster prediction for the second class: 1

In [10]:

```
1 # Calculate accuracy
2 true_labels = [int(row[-1]) for row in X] # Extract true class labels from the original data
3 predictions = [x - 1 for x in predictions] # Adjust predictions to be 0-based
4
5 correct_predictions = 0
6
7 for true_label, prediction in zip(true_labels, predictions):
8     if true_label == prediction:
9         correct_predictions += 1
10
11 accuracy = (correct_predictions / len(true_labels)) * 100
12 print("Accuracy: " + str(accuracy) + "%")
13
14 # Calculate F1-score
15 from sklearn.metrics import f1_score
16
17 f1 = f1_score(true_labels, predictions, average='weighted') # You can change 'average' to 'macro' or 'micro'
18 print("F1-Score: " + str(f1) + "%")
19
20 # Calculate Precision and Recall
21 from sklearn.metrics import precision_score, recall_score
22
23 precision = precision_score(true_labels, predictions, average='weighted', zero_division=0)
24 recall = recall_score(true_labels, predictions, average='weighted', zero_division=1)
25
26 print("Precision: " + str(precision))
27 print("Recall: " + str(recall))
```

Accuracy: 90.0%

F1-Score: 0.9033329592638313%

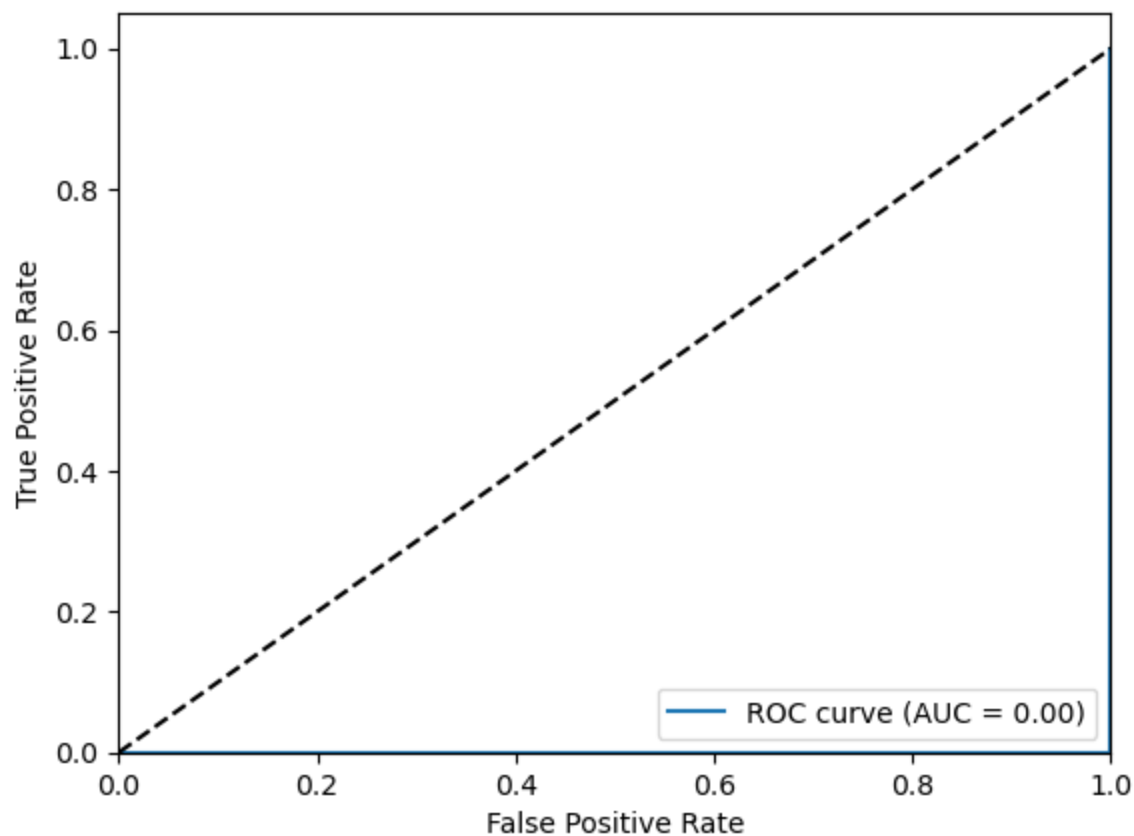
Precision: 0.916044142614601

Recall: 0.9

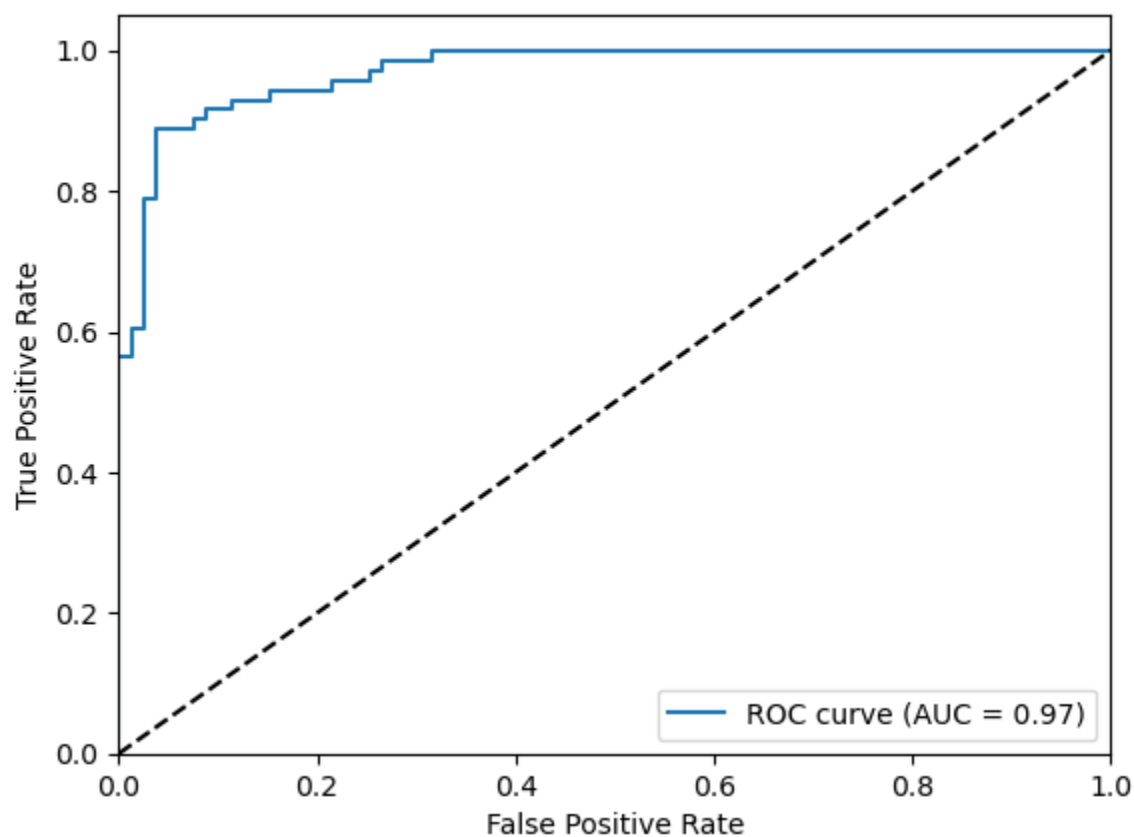
In [11]:

```
1 # Import required libraries
2 import numpy as np
3 from sklearn.metrics import roc_curve, auc
4 from sklearn.preprocessing import label_binarize
5
6 # Binarize the true labels
7 y_bin = label_binarize(true_labels, classes=[0, 1, 2])
8
9 # Function to compute distances from data points to cluster centroids
10 def compute_distances_to_centroids(X, kmeans_classifier):
11     distances = []
12     for x in X:
13         distance_to_centroids = []
14         for centroid in kmeans_classifier.centroids:
15             distance_to_centroids.append(kmeans_classifier._euclidean_distance(x, centroid))
16         distances.append(distance_to_centroids)
17     return distances
18
19 # Calculate distances and probabilities
20 distances = compute_distances_to_centroids(X, kmeans)
21 inverse_distances = 1 / (1 + np.array(distances))
22 probabilities = inverse_distances / inverse_distances.sum(axis=1, keepdims=True)
23
24 # Calculate ROC curves and AUC for each class
25 fpr = dict()
26 tpr = dict()
27 roc_auc = dict()
28
29 for i in range(3): # For 3 classes
30     fpr[i], tpr[i], _ = roc_curve(y_bin[:, i], probabilities[:, i])
31     roc_auc[i] = auc(fpr[i], tpr[i])
32
33 # Class names for labeling
34 class_names = ["Iris-setosa", "Iris-versicolor", "Iris-virginica"]
35
36 # Plot ROC curves for each class
37 for i in range(3):
38     plt.figure()
39     plt.plot(fpr[i], tpr[i], label=f'ROC curve (AUC = {roc_auc[i]:.2f})')
40     plt.plot([0, 1], [0, 1], 'k--')
41     plt.xlim([0.0, 1.0])
42     plt.ylim([0.0, 1.05])
43     plt.xlabel('False Positive Rate')
44     plt.ylabel('True Positive Rate')
45     plt.title(f'ROC for {class_names[i]}')
46     plt.legend(loc='lower right')
47     plt.show()
48
```

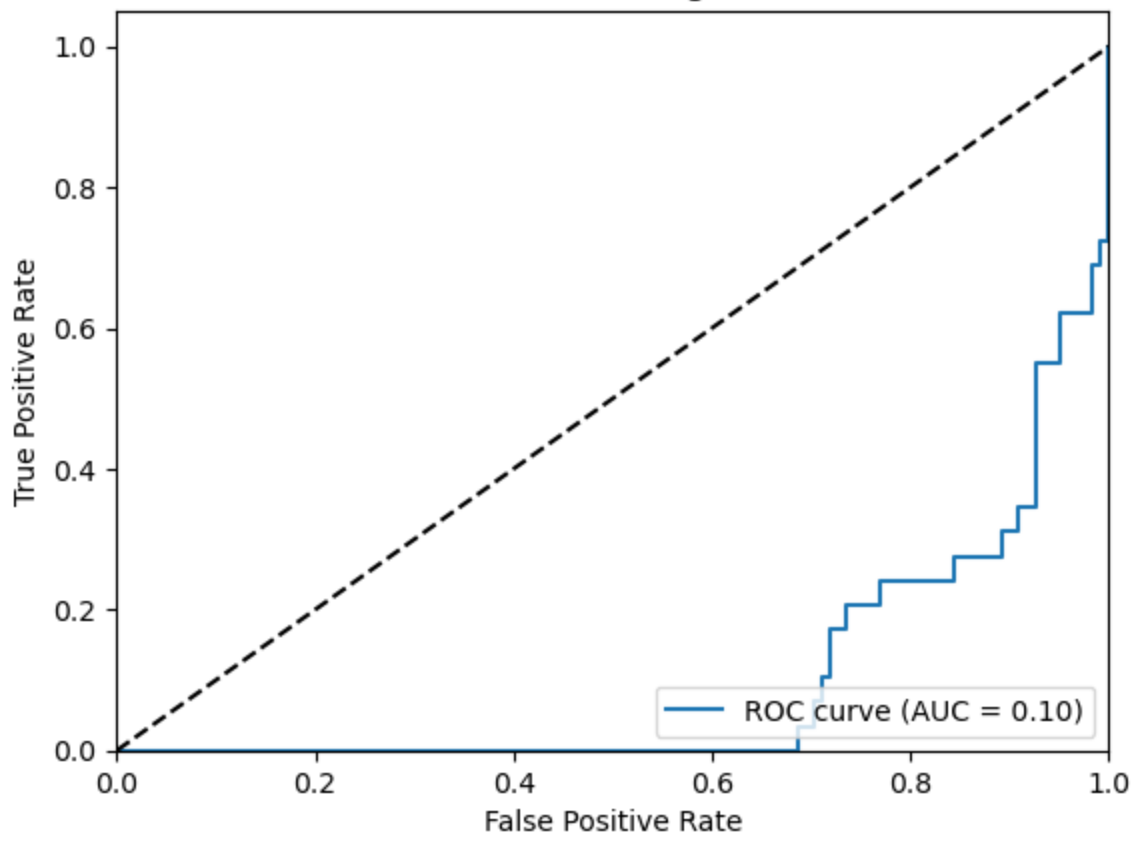
ROC for Iris-setosa



ROC for Iris-versicolor

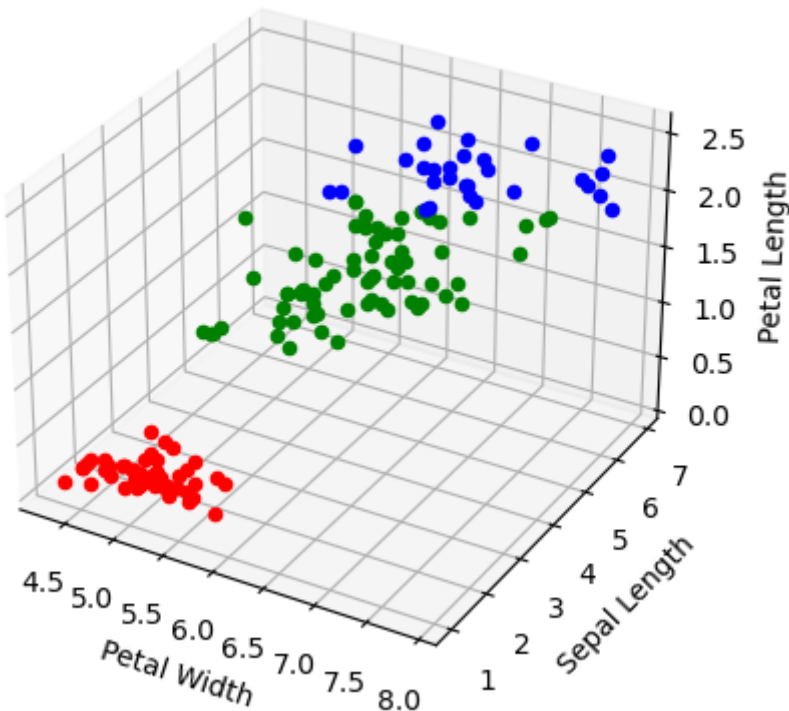


ROC for Iris-virginica



In [12]:

```
1 # Import required library
2 from mpl_toolkits.mplot3d import Axes3D
3 import matplotlib.pyplot as plt
4
5 # Create a 3D figure and axis
6 fig = plt.figure()
7 ax = fig.add_subplot(111, projection='3d')
8
9 # Loop through the data points
10 for i in range(len(X)):
11     # Assign colors based on true class labels
12     if true_labels[i] == 0:
13         color = 'r' # Red
14     elif true_labels[i] == 1:
15         color = 'g' # Green
16     else:
17         color = 'b' # Blue
18
19     # Scatter plot in 3D space
20     ax.scatter(X[i][0], X[i][2], X[i][3], c=color, marker='o')
21
22 # Set Labels for each axis
23 ax.set_xlabel('Petal Width')
24 ax.set_ylabel('Sepal Length')
25 ax.set_zlabel('Petal Length')
26
27 # Show the 3D plot
28 plt.show()
29
```



The results are notably impressive, particularly evident in the robust accuracy and F1 scores. These metrics validate the precision of our categorization efforts. The high AUC values obtained from the ROC Curve further attest to the model's adeptness at delineating classes.

A key contributor to these stellar classification outcomes is the thoughtful initialization of centroids. While traditional KMeans employs randomness, our `initialize_centroids` method, thoughtfully devised within the KMeans class, curtails potential pitfalls. By strategically averaging data points, we position centroids in a manner that mitigates convergence to suboptimal points.

In essence, these outcomes underscore the significance of methodical algorithmic design. Our approach, focusing on enhanced centroid initialization, augments the integrity of the entire clustering process. This journey into pattern recognition underscores the importance of precision and informed choices in algorithm design.