

Name: Prakhar Nagpal

Register Number: RA1911033010050

UNIVERSITY PRACTICAL EXAMINATION

Subject Code: 18CSE398J

Subject Name: Machine Learning - Core Concepts with Applications 📄

Question Number 10:

Write a program to implement the k means clustering for a sample training data set from Kaggle. Compute the compute a silhouette score.

AIM:

TO IMPLEMENET THE K MEANS CLUSTERING FOR A SAMPLE TRAINING DATASET AND THEN COMPUTE IT'S SILHOUETTE SCORE

Importing all the necessary libraries

In [29]:

```
import numpy as np
import pandas as pd

from sklearn.cluster import KMeans
from sklearn.metrics import silhouette_samples, silhouette_score

import matplotlib.pyplot as plt
import matplotlib.cm as cm
import seaborn as sns
```

Importing the DataSet and here we are using Mall Customer Segmentation dataset from Kaggle

In [21]:

```
input_data = pd.read_csv("https://raw.githubusercontent.com/PrakharNagpal/minuter-weather-k
```

In [22]:



```
input_data.head()
```

Out[22]:

	CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
0	1	Male	19	15	39
1	2	Male	21	15	81
2	3	Female	20	16	6
3	4	Female	23	16	77
4	5	Female	31	17	40

In [23]:



```
input_data.shape
```

Out[23]:

```
(200, 5)
```

In [24]:



```
input_data.isnull().sum()
```

Out[24]:

```
CustomerID      0
Gender          0
Age             0
Annual Income (k$)  0
Spending Score (1-100)  0
dtype: int64
```

In [25]:



```
# selecting features
x = input_data.iloc[:,[3,4]].values
```

In [26]:



Elbow Method

```
score = []  
  
for cluster in range(1,11):  
    kmeans = KMeans(n_clusters = cluster, init="k-means++", random_state=10)  
    kmeans.fit(x)  
    score.append(kmeans.inertia_)
```

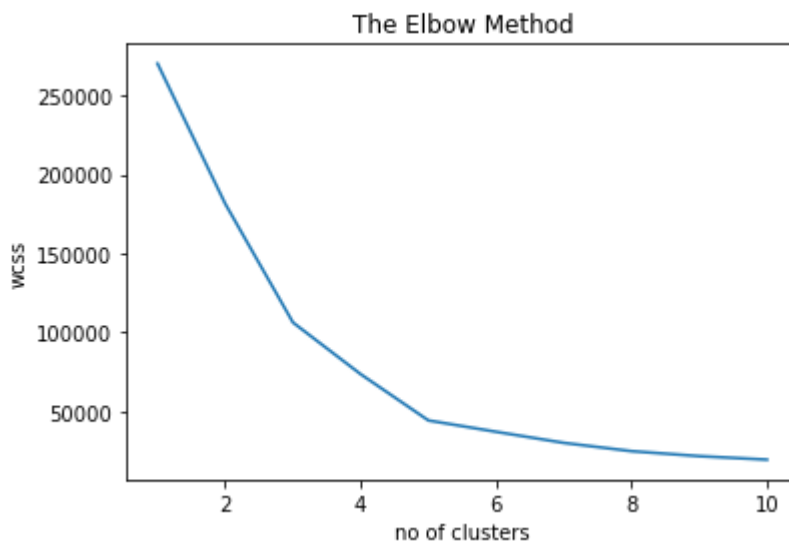
C:\Users\Prakhar\anaconda3new\lib\site-packages\sklearn\cluster_kmeans.py:81: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
warnings.warn(

In [27]:



plotting the score

```
plt.plot(range(1,11), score)  
plt.title('The Elbow Method')  
plt.xlabel('no of clusters')  
plt.ylabel('wcss')  
plt.show()
```



Hence the ideal number of clusters is 5



In [28]:

```
# Silhouette score

for n_clusters in range(2,11):
    # Create a subplot with 1 row and 2 columns
    fig, (ax1, ax2) = plt.subplots(1, 2)
    fig.set_size_inches(18, 7)

    # The 1st subplot is the silhouette plot
    # The silhouette coefficient can range from -1, 1 but in this example all
    # lie within [-0.1, 1]
    ax1.set_xlim([-0.1, 1])
    # The (n_clusters+1)*10 is for inserting blank space between silhouette
    # plots of individual clusters, to demarcate them clearly.
    ax1.set_ylim([0, len(x) + (n_clusters + 1) * 10])

    # Initialize the clusterer with n_clusters value and a random generator
    # seed of 10 for reproducibility.
    clusterer = KMeans(n_clusters=n_clusters, random_state=10)
    cluster_labels = clusterer.fit_predict(x)

    # The silhouette_score gives the average value for all the samples.
    # This gives a perspective into the density and separation of the formed
    # clusters
    silhouette_avg = silhouette_score(x, cluster_labels)
    print("For n_clusters =", n_clusters,
          "The average silhouette_score is :", silhouette_avg)

    # Compute the silhouette scores for each sample
    sample_silhouette_values = silhouette_samples(x, cluster_labels)

    y_lower = 10
    for i in range(n_clusters):
        # Aggregate the silhouette scores for samples belonging to
        # cluster i, and sort them
        ith_cluster_silhouette_values = \
            sample_silhouette_values[cluster_labels == i]

        ith_cluster_silhouette_values.sort()

        size_cluster_i = ith_cluster_silhouette_values.shape[0]
        y_upper = y_lower + size_cluster_i

        color = cm.nipy_spectral(float(i) / n_clusters)
        ax1.fill_betweenx(np.arange(y_lower, y_upper),
                          0, ith_cluster_silhouette_values,
                          facecolor=color, edgecolor=color, alpha=0.7)

        # Label the silhouette plots with their cluster numbers at the middle
        ax1.text(-0.05, y_lower + 0.5 * size_cluster_i, str(i))

        # Compute the new y_lower for next plot
        y_lower = y_upper + 10 # 10 for the 0 samples

    ax1.set_title("The silhouette plot for the various clusters.")
    ax1.set_xlabel("The silhouette coefficient values")
    ax1.set_ylabel("Cluster label")

    # The vertical line for average silhouette score of all the values
    ax1.axvline(x=silhouette_avg, color="red", linestyle="--")
```

```

ax1.set_yticks([]) # Clear the yaxis labels / ticks
ax1.set_xticks([-0.1, 0, 0.2, 0.4, 0.6, 0.8, 1])

# 2nd Plot showing the actual clusters formed
colors = cm.nipy_spectral(cluster_labels.astype(float) / n_clusters)
ax2.scatter(x[:, 0], x[:, 1], marker='.', s=30, lw=0, alpha=0.7,
            c=colors, edgecolor='k')

# Labeling the clusters
centers = clusterer.cluster_centers_
# Draw white circles at cluster centers
ax2.scatter(centers[:, 0], centers[:, 1], marker='o',
            c="white", alpha=1, s=200, edgecolor='k')

for i, c in enumerate(centers):
    ax2.scatter(c[0], c[1], marker='$%d$' % i, alpha=1,
                s=50, edgecolor='k')

ax2.set_title("The visualization of the clustered data.")
ax2.set_xlabel("Feature space for the 1st feature")
ax2.set_ylabel("Feature space for the 2nd feature")

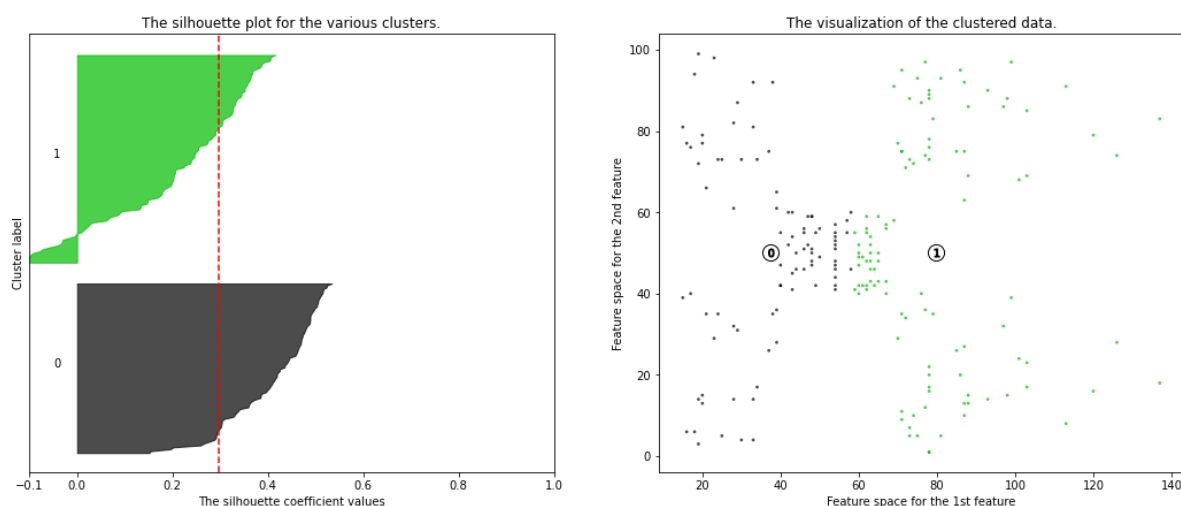
plt.suptitle(("Silhouette analysis for KMeans clustering on sample data "
             "with n_clusters = %d" % n_clusters),
             fontsize=14, fontweight='bold')

plt.show()

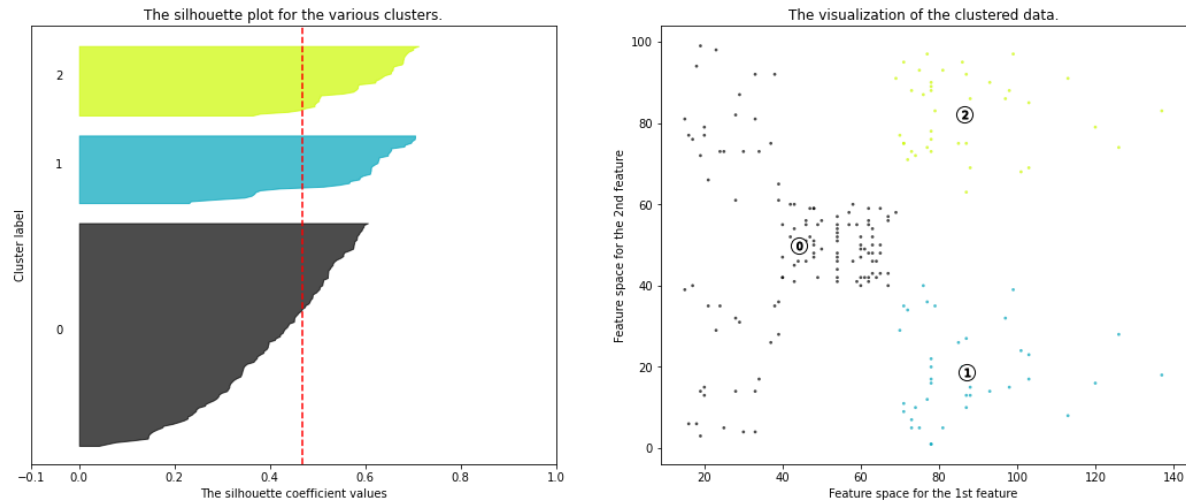
```

For n_clusters = 2 The average silhouette_score is : 0.2968969162503008
 For n_clusters = 3 The average silhouette_score is : 0.46761358158775435
 For n_clusters = 4 The average silhouette_score is : 0.4931963109249047
 For n_clusters = 5 The average silhouette_score is : 0.553931997444648
 For n_clusters = 6 The average silhouette_score is : 0.5376203956398481
 For n_clusters = 7 The average silhouette_score is : 0.5270287298101395
 For n_clusters = 8 The average silhouette_score is : 0.4572211842776841
 For n_clusters = 9 The average silhouette_score is : 0.45872989167156364
 For n_clusters = 10 The average silhouette_score is : 0.4467356774401869

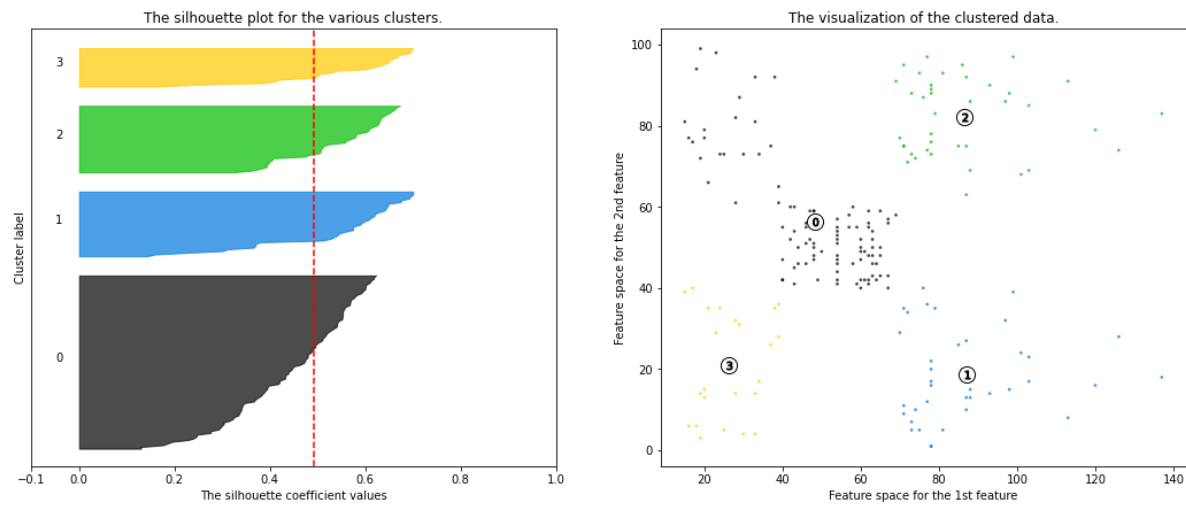
Silhouette analysis for KMeans clustering on sample data with n_clusters = 2



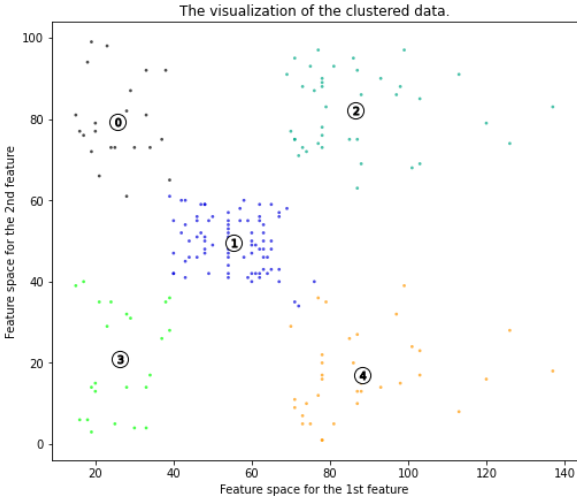
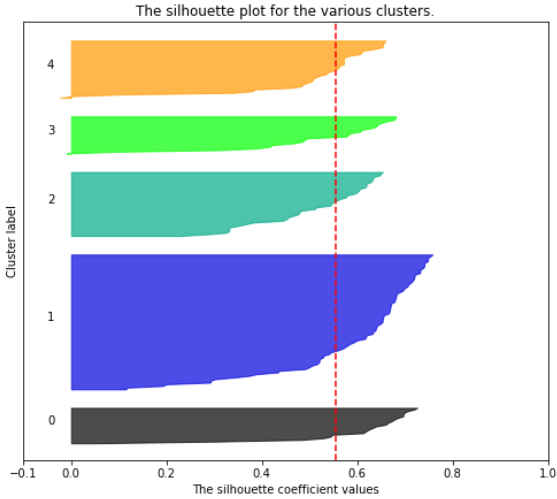
Silhouette analysis for KMeans clustering on sample data with n_clusters = 3



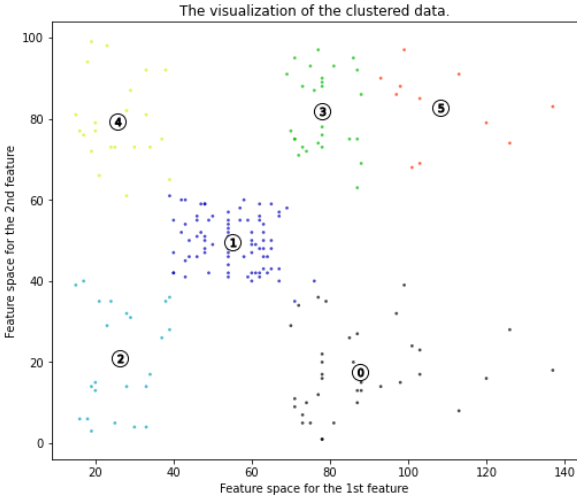
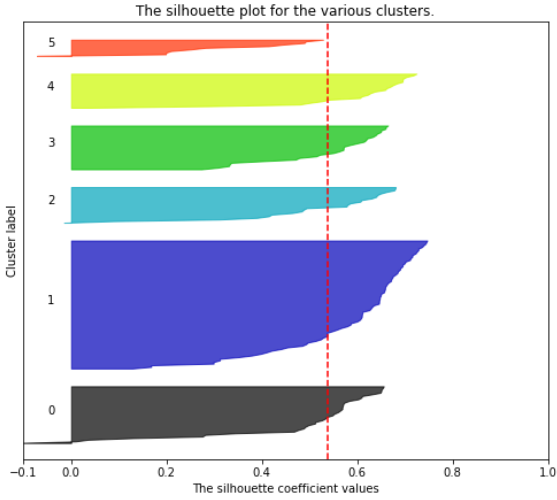
Silhouette analysis for KMeans clustering on sample data with n_clusters = 4



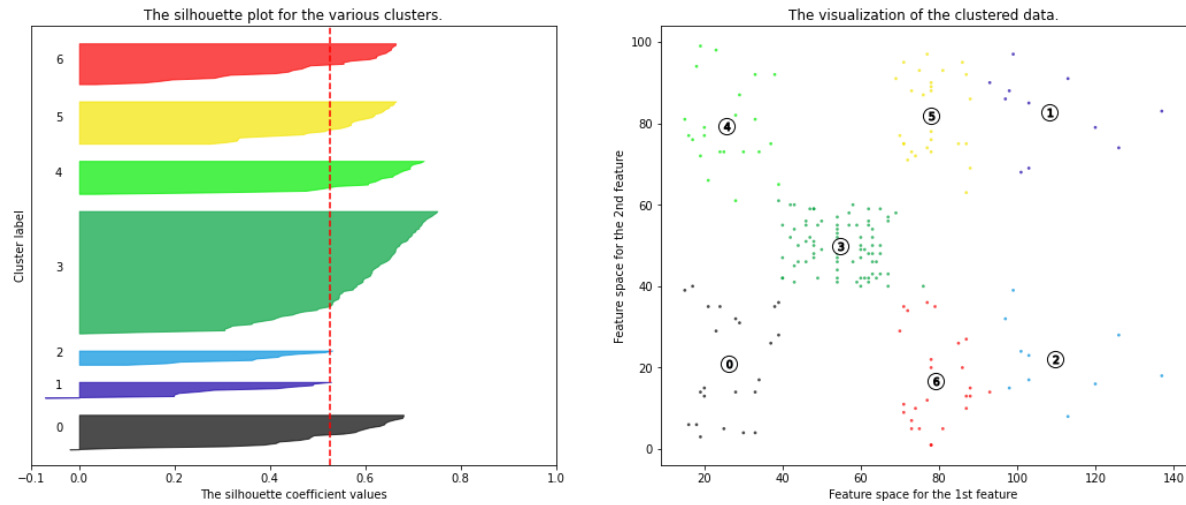
Silhouette analysis for KMeans clustering on sample data with n_clusters = 5



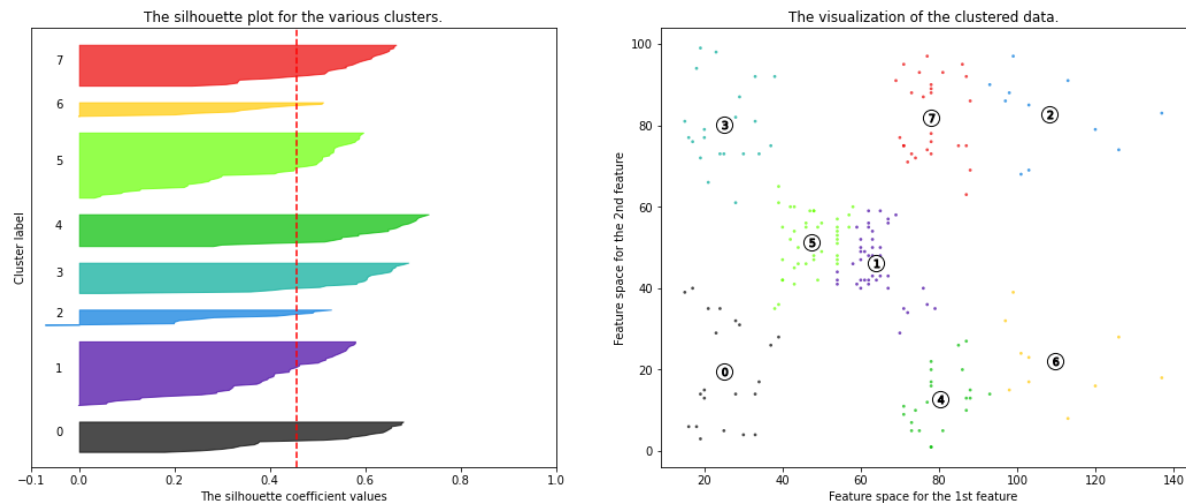
Silhouette analysis for KMeans clustering on sample data with n_clusters = 6



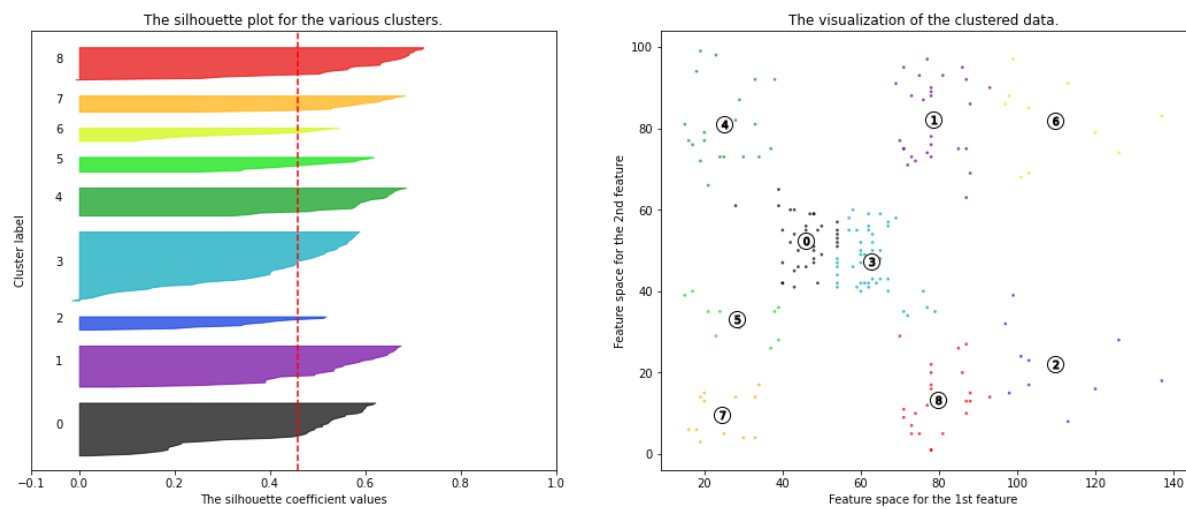
Silhouette analysis for KMeans clustering on sample data with n_clusters = 7



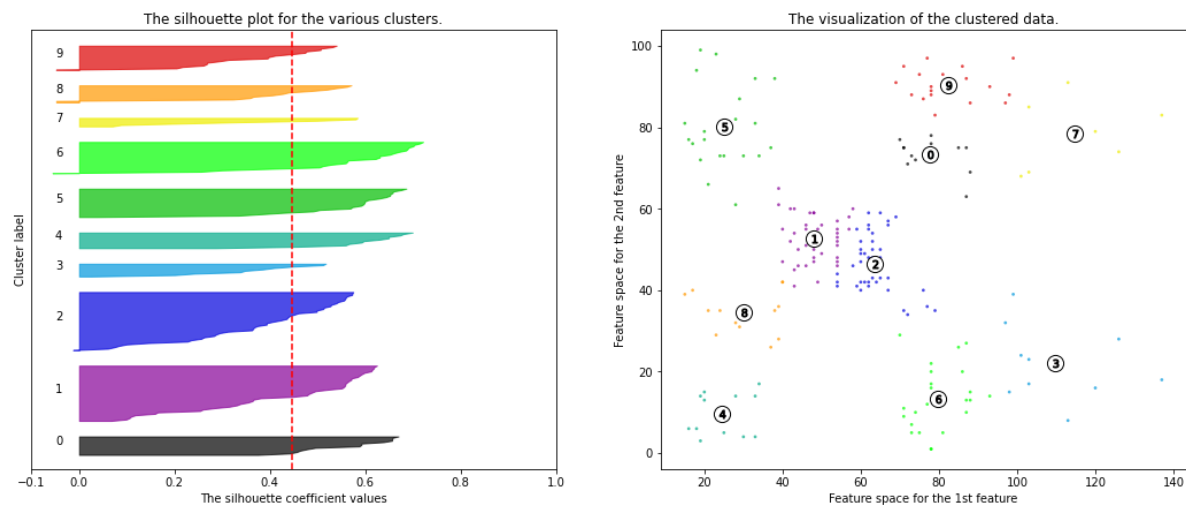
Silhouette analysis for KMeans clustering on sample data with n_clusters = 8



Silhouette analysis for KMeans clustering on sample data with n_clusters = 9



Silhouette analysis for KMeans clustering on sample data with n_clusters = 10



Inference:

Hence here we have determined silhouette scores for n_clusters ranging from 2 to 10 and the score were :

For n_clusters = 2 The average silhouette_score is : 0.2968969162503008

For n_clusters = 3 The average silhouette_score is : 0.46761358158775435

For n_clusters = 4 The average silhouette_score is : 0.4931963109249047

For n_clusters = 5 The average silhouette_score is : 0.553931997444648

For n_clusters = 6 The average silhouette_score is : 0.5376203956398481

For n_clusters = 7 The average silhouette_score is : 0.5270287298101395

For n_clusters = 8 The average silhouette_score is : 0.4572211842776841

For n_clusters = 9 The average silhouette_score is : 0.45872989167156364

For n_clusters = 10 The average silhouette_score is : 0.4467356774401869

Result:

Hence the model was successfully built for Mall Segmentation Dataset and silhouette scores were observed for various n_clusters ranging from 2 to 10