

Performance Metrics Clustering-Silhouetter Coefficient

- Aryan Shriva

In [1]:

```
1 from sklearn.datasets import make_blobs
2 from sklearn.cluster import KMeans
3 from sklearn.metrics import silhouette_samples, silhouette_score
4
5 import matplotlib.pyplot as plt
6 import matplotlib.cm as cm
7 import numpy as np
8
```

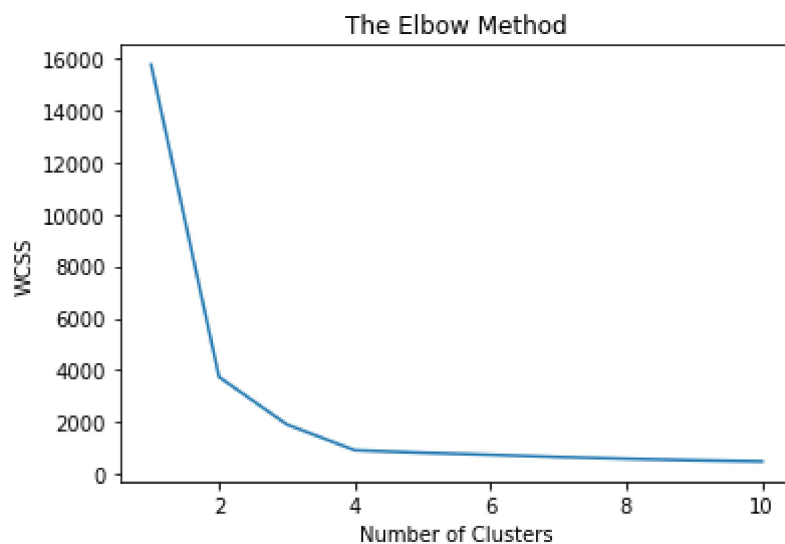
In [2]:

```
1 # Generating the sample data from make_blobs
2 # This particular setting has one distinct cluster and 3 clusters placed clo
3 # together.
4 X, y = make_blobs(n_samples=500,
5                   n_features=2,
6                   centers=4,
7                   cluster_std=1,
8                   center_box=(-10.0, 10.0),
9                   shuffle=True,
10                  random_state=1) # For reproducibility
11
12 range_n_clusters = [2, 3, 4, 5, 6]
```

```
In [3]: 1 from sklearn.cluster import KMeans
2
3 wcss=[]
4 for i in range(1,11):
5     kmeans=KMeans(n_clusters=i, init='k-means++',random_state=0)
6     kmeans.fit(X)
7     wcss.append(kmeans.inertia_)
8
9 plt.plot(range(1,11),wcss)
10 plt.title('The Elbow Method')
11 plt.xlabel('Number of Clusters')
12 plt.ylabel('WCSS')
13 plt.show()
```

C:\Users\arya shriva\anaconda3\lib\site-packages\sklearn\cluster_kmeans.py:88
 1: 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 e
 nvironment variable OMP_NUM_THREADS=2.

warnings.warn(



```
In [4]: 1 clusterer = KMeans(n_clusters=4, random_state=10)
2 cluster_labels = clusterer.fit_predict(X)
3 print(cluster_labels)
```

```
[2 2 3 1 0 1 0 0 0 0 2 2 0 1 0 2 0 2 1 0 3 3 0 1 0 0 1 1 3 0 2 1 0 2 0 2 3
 3 2 3 0 3 1 0 0 2 3 0 1 1 1 3 3 0 2 3 3 3 3 0 1 1 3 0 1 0 2 0 3 3 2 3 0 2
 0 0 2 0 0 3 1 1 3 1 1 3 3 1 3 3 1 2 3 0 1 2 2 0 2 1 1 2 1 3 1 0 0 1 1 3 0
 2 1 3 1 3 1 0 1 0 3 2 2 3 0 3 1 2 2 0 1 3 3 3 3 2 1 0 1 1 0 2 0 1 1 1 0 0
 2 2 3 3 1 2 1 3 3 3 3 3 3 3 3 1 2 2 2 0 1 2 3 0 2 1 3 3 3 3 2 0 3 1 2 2
 3 0 2 2 0 1 1 2 2 0 1 0 2 2 1 2 3 1 0 0 2 0 3 2 0 3 0 3 2 0 0 0 1 3 1 0 2
 3 0 3 3 3 1 3 1 2 3 2 3 1 1 3 2 1 2 0 3 2 2 2 2 0 3 2 3 0 1 1 0 0 1 3 0 3
 1 0 1 3 3 1 0 2 2 3 3 3 0 1 1 0 1 3 2 1 2 1 2 2 1 2 1 1 0 3 3 3 0 0 3 2 1
 2 2 2 0 3 0 2 3 2 2 3 2 2 3 1 2 0 0 1 1 3 2 1 1 0 2 1 1 0 3 1 3 0 2 2 1 3
 2 0 1 1 0 0 0 2 0 1 1 3 1 1 1 1 2 2 0 1 3 0 2 1 3 1 0 1 3 0 3 1 0 0 2 1 2
 2 2 2 2 3 2 1 2 1 1 3 1 0 3 3 2 1 3 1 0 2 3 3 2 3 3 1 1 2 3 0 1 0 0 2 2
 0 2 3 3 2 3 2 3 1 2 1 3 0 1 3 0 1 2 0 1 1 3 0 3 0 2 1 2 0 1 2 2 2 3 1 0 2
 0 0 3 3 2 0 0 0 0 0 2 0 3 2 0 1 0 1 0 3 3 1 1 1 3 0 3 2 3 1 0 2 1 2 1 2
 0 1 1 2 3 0 2 3 3 3 2 0 1 3 0 2 2 2 0]
```

```

In [5]: 1 for n_clusters in range_n_clusters:
2         # Create a subplot with 1 row and 2 columns
3         fig, (ax1, ax2) = plt.subplots(1, 2)
4         fig.set_size_inches(18, 7)
5
6         # The 1st subplot is the silhouette plot
7         # The silhouette coefficient can range from -1, 1 but in this example al
8         # lie within [-0.1, 1]
9         ax1.set_xlim([-0.1, 1])
10        # The (n_clusters+1)*10 is for inserting blank space between silhouette
11        # plots of individual clusters, to demarcate them clearly.
12        ax1.set_ylim([0, len(X) + (n_clusters + 1) * 10])
13
14        # Initialize the clusterer with n_clusters value and a random generator
15        # seed of 10 for reproducibility.
16        clusterer = KMeans(n_clusters=n_clusters, random_state=10)
17        cluster_labels = clusterer.fit_predict(X)
18
19        # The silhouette_score gives the average value for all the samples.
20        # This gives a perspective into the density and separation of the formed
21        # clusters
22        silhouette_avg = silhouette_score(X, cluster_labels)
23        print("For n_clusters =", n_clusters,
24              "The average silhouette_score is :", silhouette_avg)
25
26        # Compute the silhouette scores for each sample
27        sample_silhouette_values = silhouette_samples(X, cluster_labels)
28
29        y_lower = 10
30        for i in range(n_clusters):
31            # Aggregate the silhouette scores for samples belonging to
32            # cluster i, and sort them
33            ith_cluster_silhouette_values = \
34                sample_silhouette_values[cluster_labels == i]
35
36            ith_cluster_silhouette_values.sort()
37
38            size_cluster_i = ith_cluster_silhouette_values.shape[0]
39            y_upper = y_lower + size_cluster_i
40
41            color = cm.nipy_spectral(float(i) / n_clusters)
42            ax1.fill_betweenx(np.arange(y_lower, y_upper),
43                             0, ith_cluster_silhouette_values,
44                             facecolor=color, edgecolor=color, alpha=0.7)
45
46            # Label the silhouette plots with their cluster numbers at the middl
47            ax1.text(-0.05, y_lower + 0.5 * size_cluster_i, str(i))
48
49            # Compute the new y_lower for next plot
50            y_lower = y_upper + 10 # 10 for the 0 samples
51
52        ax1.set_title("The silhouette plot for the various clusters.")
53        ax1.set_xlabel("The silhouette coefficient values")
54        ax1.set_ylabel("Cluster label")
55
56        # The vertical line for average silhouette score of all the values

```

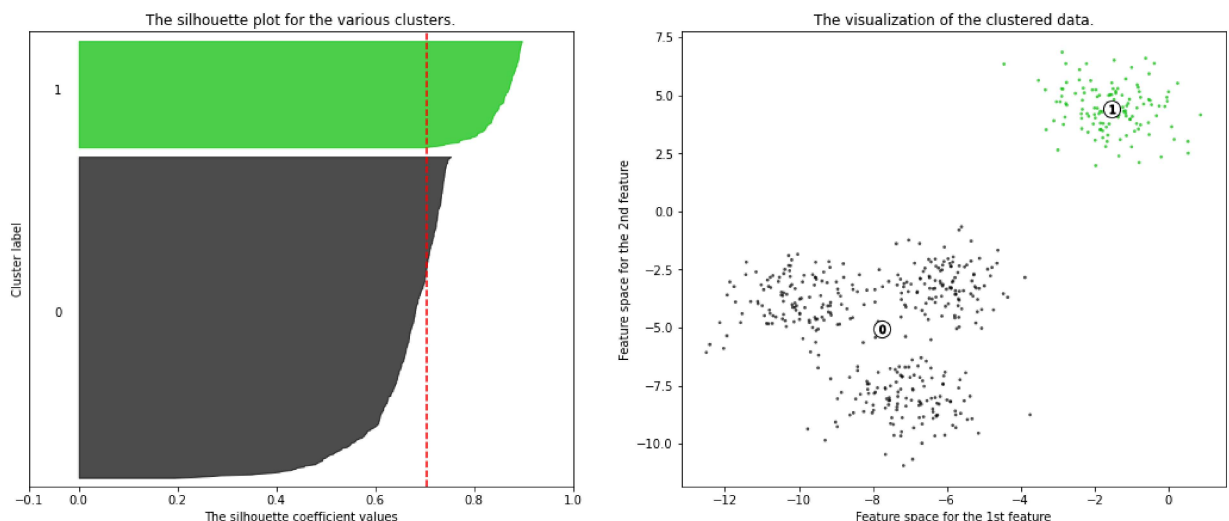
```

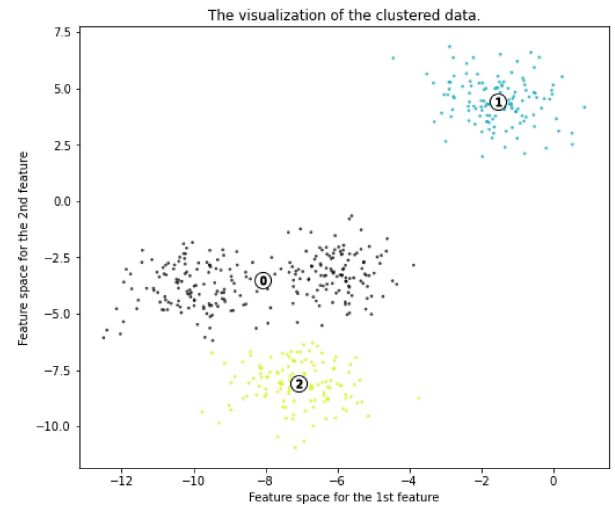
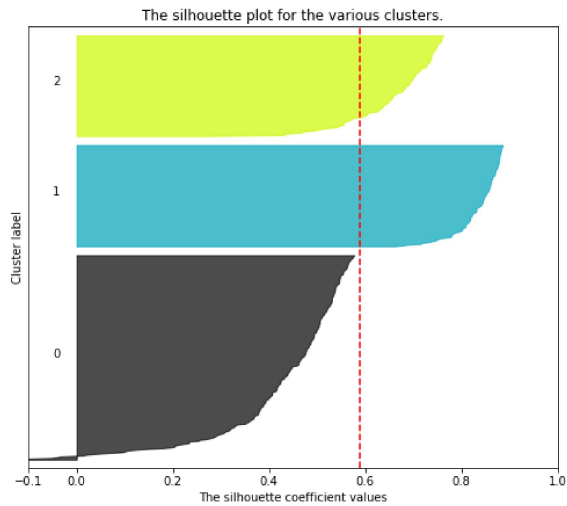
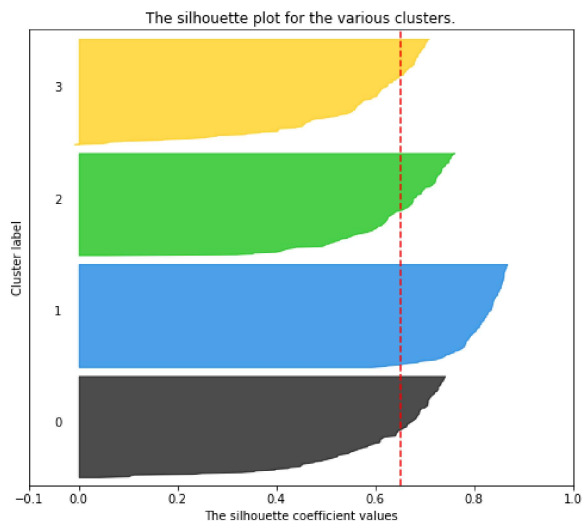
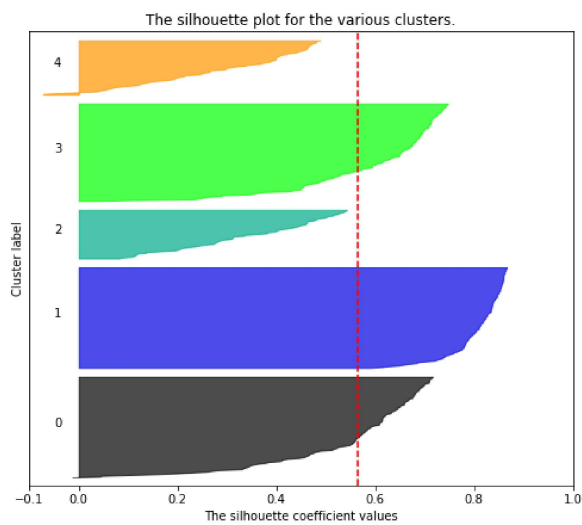
57 ax1.axvline(x=silhouette_avg, color="red", linestyle="--")
58
59 ax1.set_yticks([]) # Clear the yaxis labels / ticks
60 ax1.set_xticks([-0.1, 0, 0.2, 0.4, 0.6, 0.8, 1])
61
62 # 2nd Plot showing the actual clusters formed
63 colors = cm.nipy_spectral(cluster_labels.astype(float) / n_clusters)
64 ax2.scatter(X[:, 0], X[:, 1], marker='.', s=30, lw=0, alpha=0.7,
65             c=colors, edgecolor='k')
66
67 # Labeling the clusters
68 centers = clusterer.cluster_centers_
69 # Draw white circles at cluster centers
70 ax2.scatter(centers[:, 0], centers[:, 1], marker='o',
71             c="white", alpha=1, s=200, edgecolor='k')
72
73 for i, c in enumerate(centers):
74     ax2.scatter(c[0], c[1], marker='$%d$' % i, alpha=1,
75                 s=50, edgecolor='k')
76
77 ax2.set_title("The visualization of the clustered data.")
78 ax2.set_xlabel("Feature space for the 1st feature")
79 ax2.set_ylabel("Feature space for the 2nd feature")
80
81 plt.suptitle(("Silhouette analysis for KMeans clustering on sample data
82              with n_clusters = %d" % n_clusters),
83              fontsize=14, fontweight='bold')
84
85 plt.show()

```

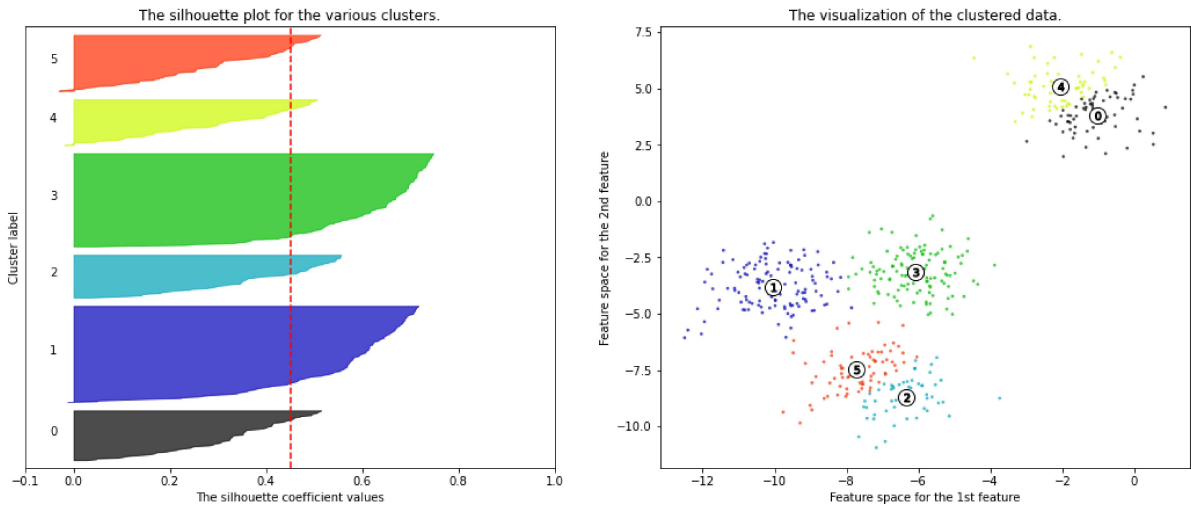
For n_clusters = 2 The average silhouette_score is : 0.7049787496083262
 For n_clusters = 3 The average silhouette_score is : 0.5882004012129721
 For n_clusters = 4 The average silhouette_score is : 0.6505186632729437
 For n_clusters = 5 The average silhouette_score is : 0.56376469026194
 For n_clusters = 6 The average silhouette_score is : 0.4504666294372765

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



In []:

1