

2. Clustering

2.1. First dataset

Scatter Plot

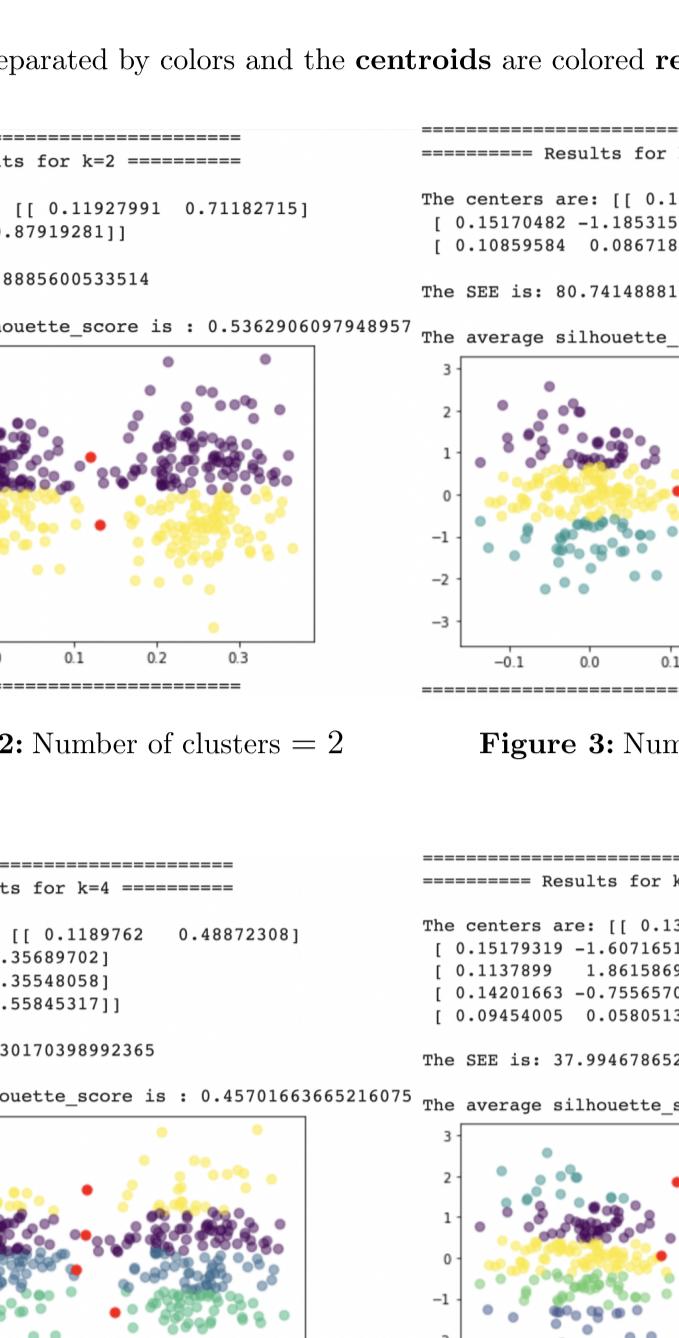


Figure 1: Scatter plot for the First dataset

Type of clusters: Center-based Clusters

- A cluster is a set of objects such that an object in a cluster is closer (more similar) to the “center” of a cluster, than to the center of any other cluster.
- The center of a cluster is often a centroid, the average of all the points in the cluster, or a medoid, the most “representative” point of a cluster.

Predicted number of clusters: $K = 2$

Lloyd's k-means Clustering

Clusters are separated by colors and the centroids are colored red.

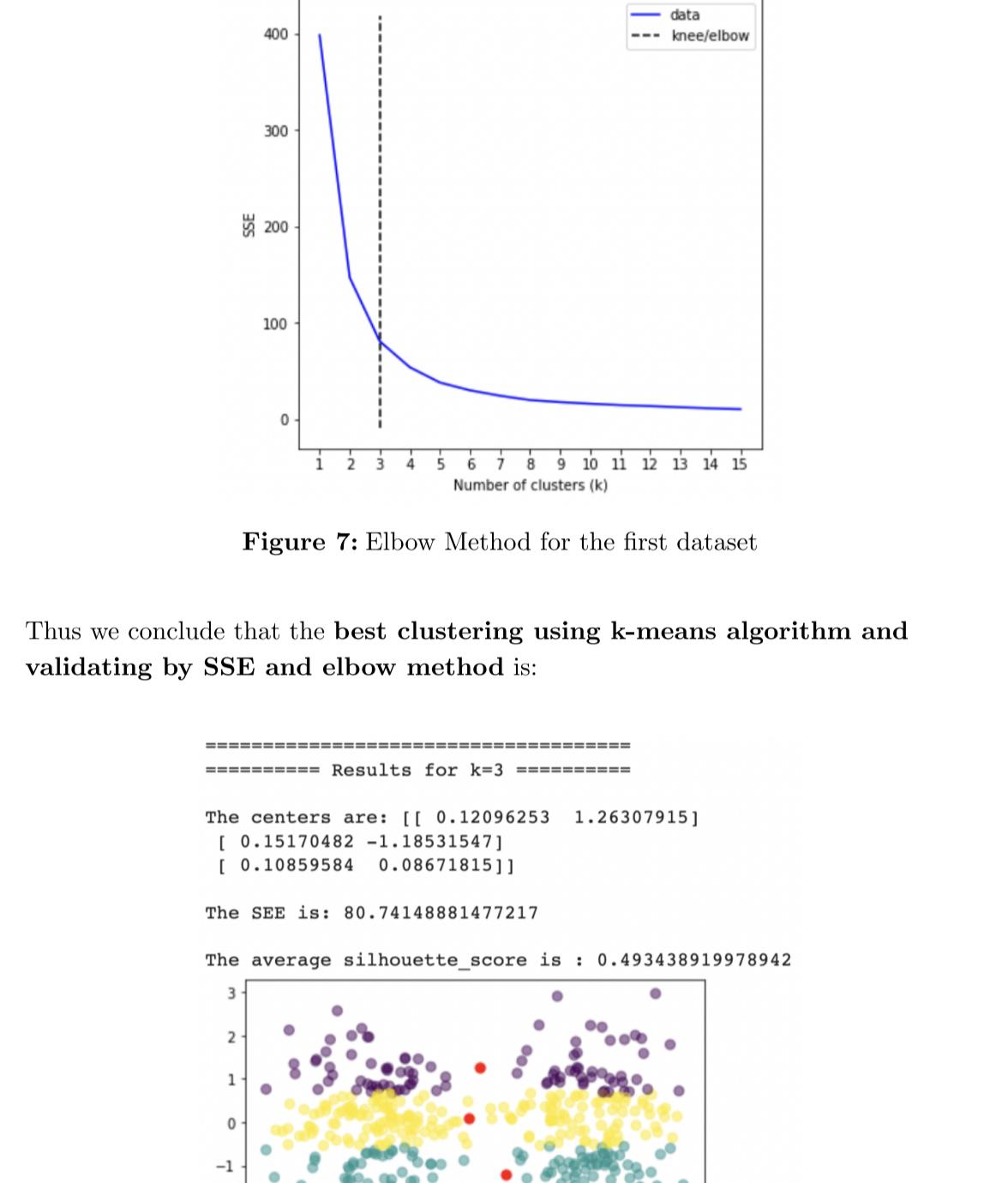


Figure 2: Number of clusters = 2

Figure 3: Number of clusters = 3

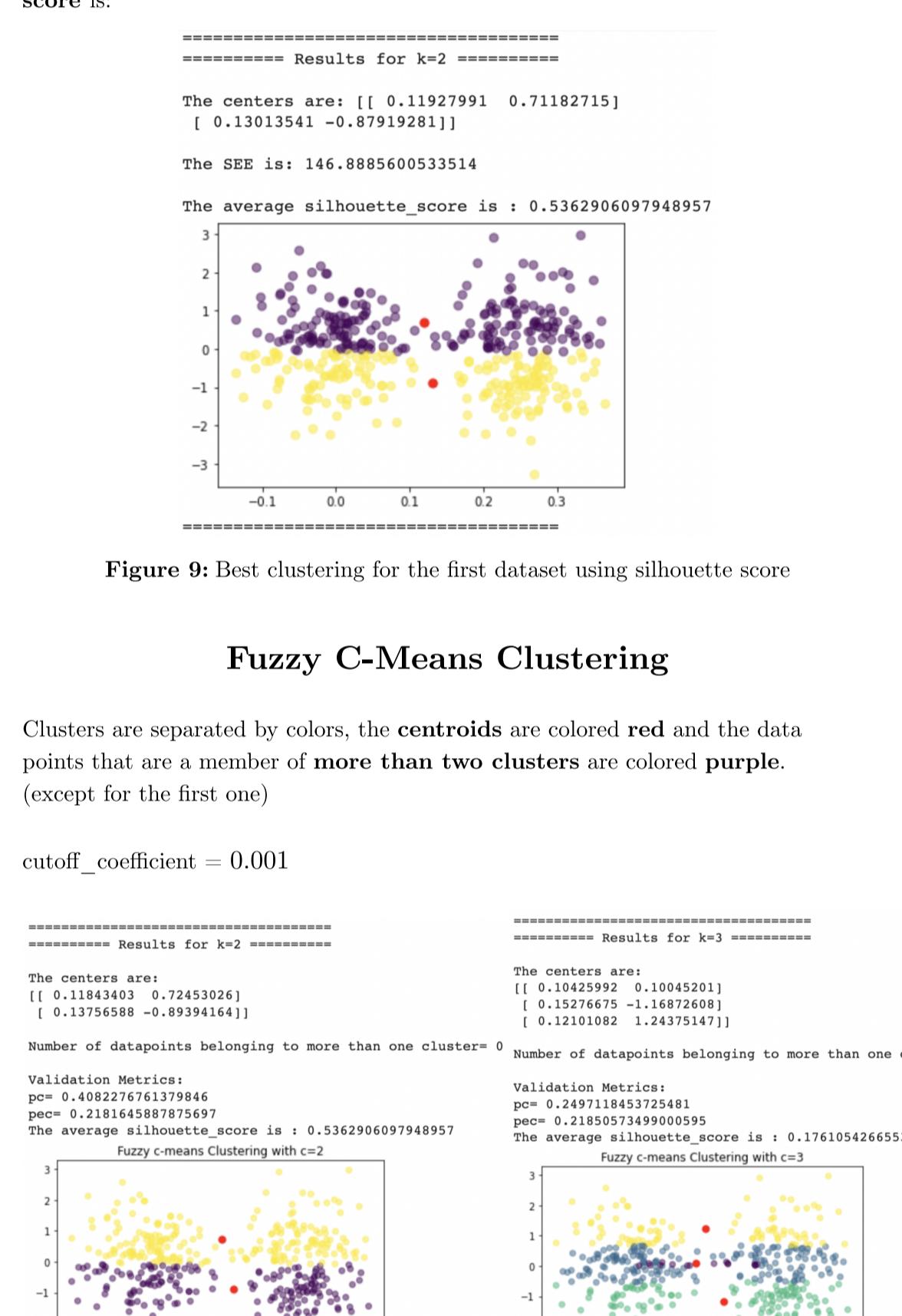


Figure 4: Number of clusters = 4

Figure 5: Number of clusters = 5

Figure 6: Number of clusters = 6

Minimum SSE: 29

Number of clusters with minimum SSE: 6

** Be careful! a small number of SSE does not guarantee a good clustering! Elbow method is used to solve this in the next part.**

Maximum silhouette score: 0.54

Number of clusters with maximum silhouette score: 2

K-Means Elbow Method

Elbow point: 3
SSE for this point is: 53.82732990087291
<Figure size 360x360 with 0 Axes>

The Elbow Method showing the optimal k

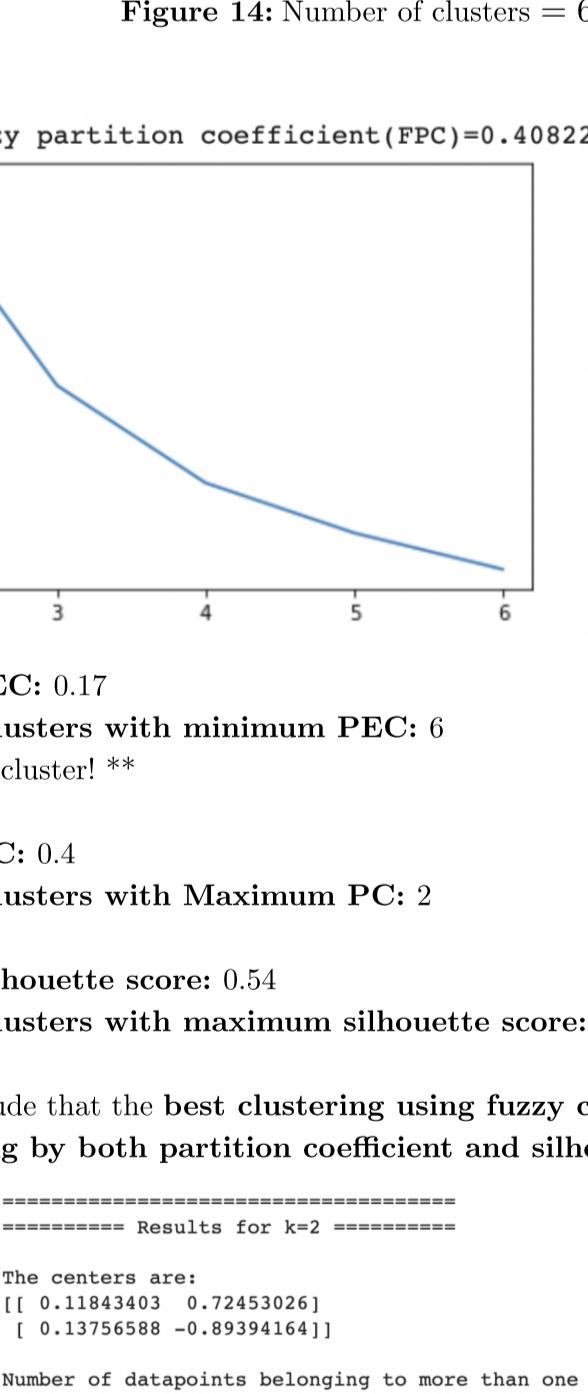


Figure 7: Elbow Method for the first dataset

Thus we conclude that the best clustering using k-means algorithm and validating by SSE and elbow method is:

===== Results for k=2 =====

The centers are: [[0.12096253 1.26307915]

[0.10859584 0.08671815]]

The SEE is: 80.74148881477217

The average silhouette_score is : 0.493438919978942

===== Results for k=3 =====

The centers are: [[0.10425992 0.10045201]

[0.15276675 -1.18531547]

[0.10859584 0.08671815]]

The SEE is: 80.74148881477217

The average silhouette_score is : 0.493438919978942

===== Results for k=4 =====

The centers are: [[0.11927991 0.71182715]

[0.13013541 -0.87919281]]

The SEE is: 146.8885600533514

The average silhouette_score is : 0.5362906097948957

===== Results for k=5 =====

The centers are: [[0.15170482 -1.18531547]

[0.10859584 0.08671815]]

The SEE is: 80.74148881477217

The average silhouette_score is : 0.493438919978942

===== Results for k=6 =====

The centers are: [[0.10425992 0.10045201]

[0.15276675 -1.18531547]

[0.10859584 0.08671815]]

The SEE is: 80.74148881477217

The average silhouette_score is : 0.493438919978942

===== Results for k=7 =====

The centers are: [[0.11927991 0.71182715]

[0.13013541 -0.87919281]]

The SEE is: 146.8885600533514

The average silhouette_score is : 0.5362906097948957

===== Results for k=8 =====

The centers are: [[0.15170482 -1.18531547]

[0.10859584 0.08671815]]

The SEE is: 80.74148881477217

The average silhouette_score is : 0.493438919978942

===== Results for k=9 =====

The centers are: [[0.10425992 0.10045201]

[0.15276675 -1.18531547]

[0.10859584 0.08671815]]

The SEE is: 80.74148881477217

The average silhouette_score is : 0.493438919978942

===== Results for k=10 =====

The centers are: [[0.11927991 0.71182715]

[0.13013541 -0.87919281]]

The SEE is: 146.8885600533514

The average silhouette_score is : 0.5362906097948957

===== Results for k=11 =====

The centers are: [[0.15170482 -1.18531547]

[0.10859584 0.08671815]]

The SEE is: 80.74148881477217

The average silhouette_score is : 0.493438919978942

===== Results for k=12 =====

The centers are: [[0.10425992 0.10045201]

[0.15276675 -1.18531547]

[0.10859584 0.08671815]]

The SEE is: 80.74148881477217

The average silhouette_score is : 0.493438919978942

===== Results for k=13 =====

The centers are: [[0.11927991 0.71182715]

[0.13013541 -0.87919281]]

The SEE is: 146.8885600533514

The average silhouette_score is : 0.5362906097948957

===== Results for k=14 =====

The centers are: [[0.15170482 -1.18531547]

[0.10859584 0.08671815]]

The SEE is: 80.74148881477217

The average silhouette_score is : 0.493438919978942

===== Results for k=15 =====

The centers are: [[0.10425992 0.10045201]

[0.15276675 -1.18531547]

[0.10859584 0.08671815]]

The SEE is: 80.74148881477217

The average silhouette_score is : 0.493438919978942

===== Results for k=16 =====

The centers are: [[0.11927991 0.71182715]

[0.13013541 -0.87919281]]

The SEE is: 146.8885600533514

The average silhouette_score is : 0.5362906097948957

===== Results for k=17 =====

The centers are: [[0.15170482 -1.18531547]

[0.10859584 0.08671815]]

The SEE is: 80.74148881477217

The average silhouette_score is : 0.493438919978942

===== Results for k=18 =====

The centers are: [[0.10425992 0.10045201]

[0.15276675 -1.18531547]

[0.10859584 0.08671815]]

The SEE is: 80.74148881477217

The average silhouette_score is : 0.493438919978942

===== Results for k=19 =====

The centers are: [[0.11927991 0.71182715]

[0.13013541 -0.87919281]]

The SEE is: 146.8885600533514

The average silhouette_score is : 0.5362906097948957

===== Results for k=20 =====

The centers are: [[0.15170482 -1.18531547]

[0.10859584 0.08671815]]

The SEE is: 80.74148881477217

The average silhouette_score is : 0.493438919978942

===== Results for k=21 =====

The centers are: [[0.10425992 0.10045201]

[0.15276675 -1.18531547]

[0.10859584 0.08671815]]

The SEE is: 80.74148881477217

The average silhouette_score is : 0.493438919978942

===== Results for k=22 =====