

2. Clustering

2.2. Second dataset

Scatter Plot

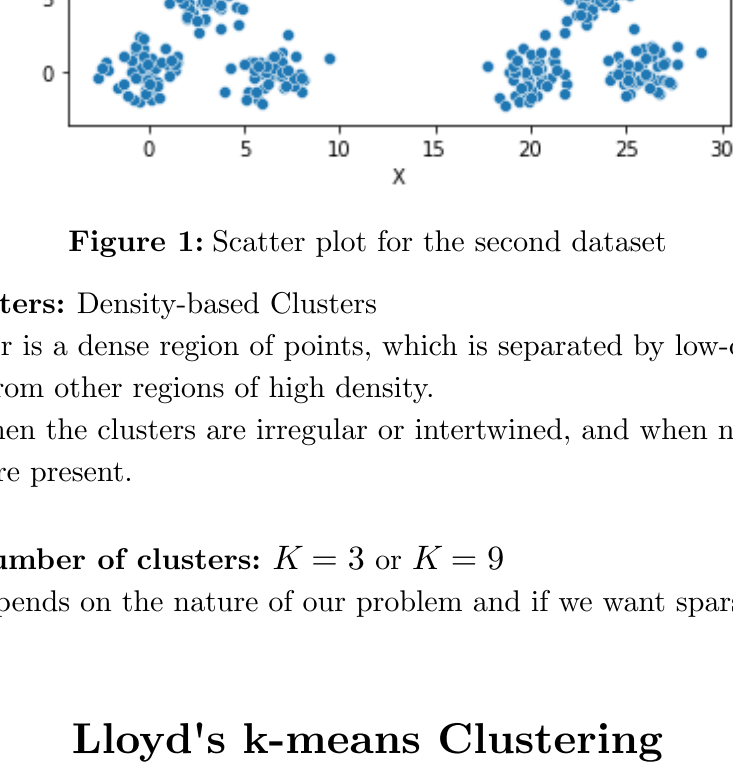


Figure 1: Scatter plot for the second dataset

Type of clusters: Density-based Clusters

- A cluster is a dense region of points, which is separated by low-density regions, from other regions of high density.
- Used when the clusters are irregular or intertwined, and when noise and outliers are present.

Predicted number of clusters: $K = 3$ or $K = 9$

This solely depends on the nature of our problem and if we want sparse or dense clusters.

Lloyd's k-means Clustering

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

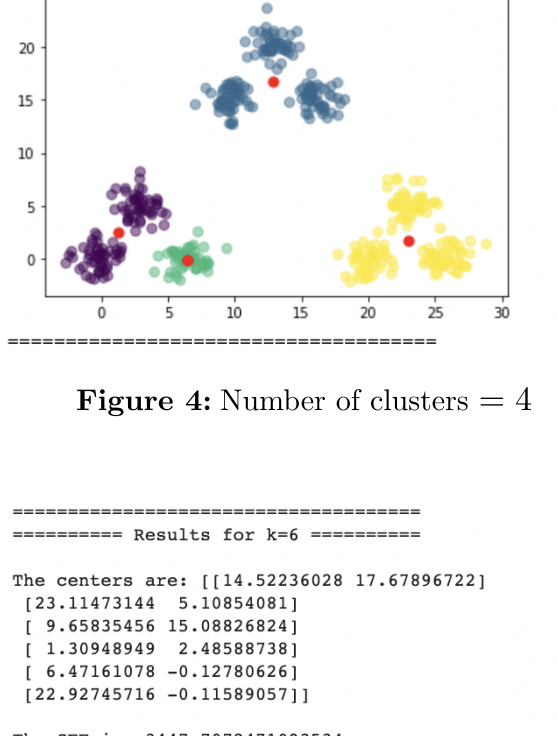


Figure 2: Number of clusters = 2

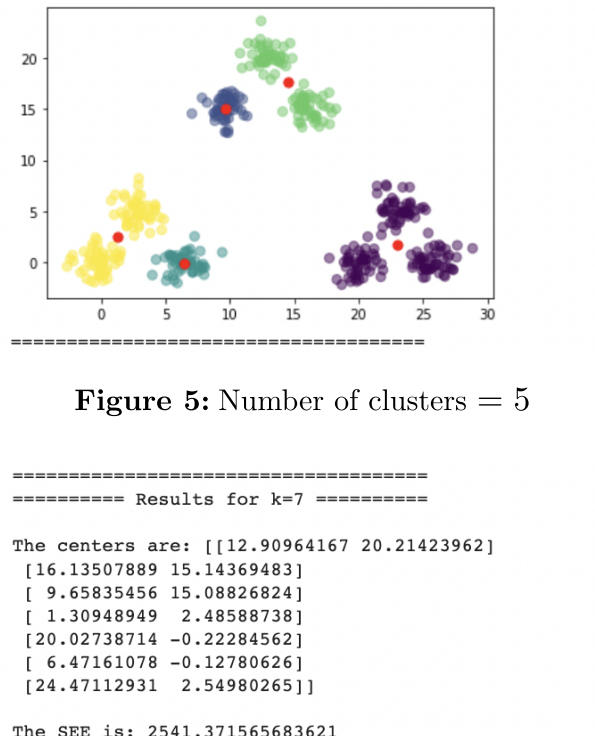


Figure 3: Number of clusters = 3

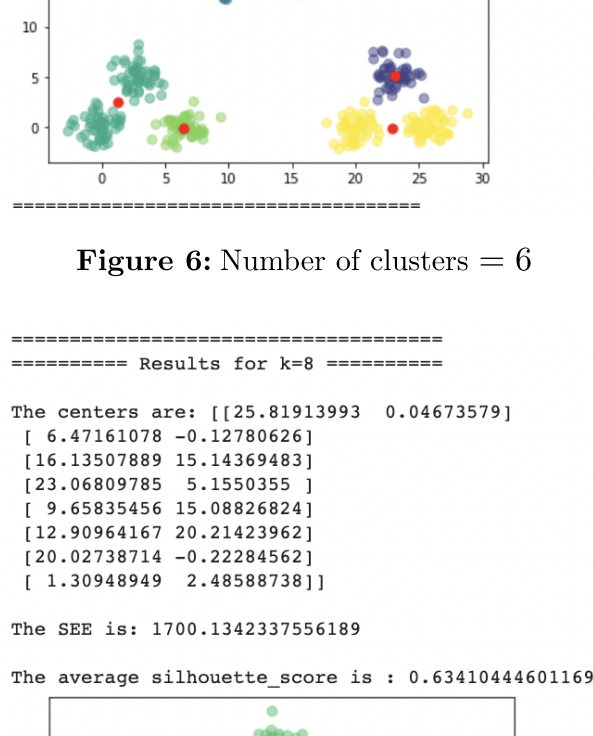


Figure 4: Number of clusters = 4

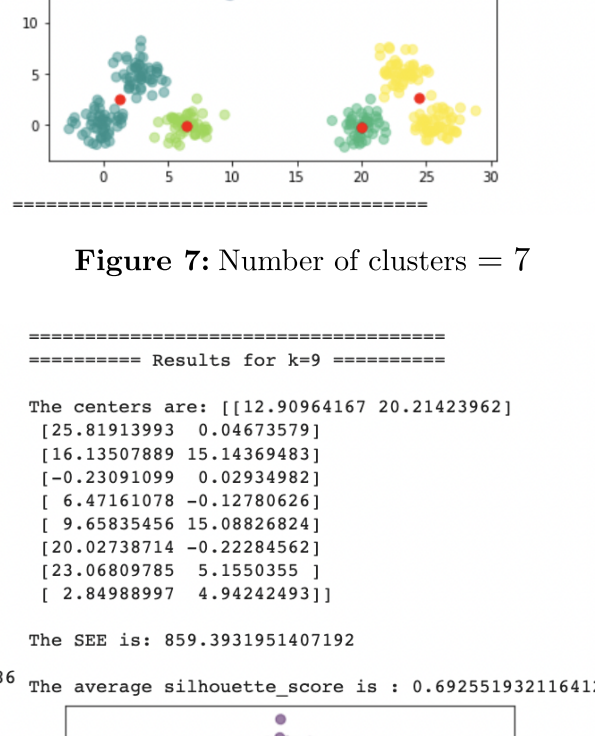


Figure 5: Number of clusters = 5

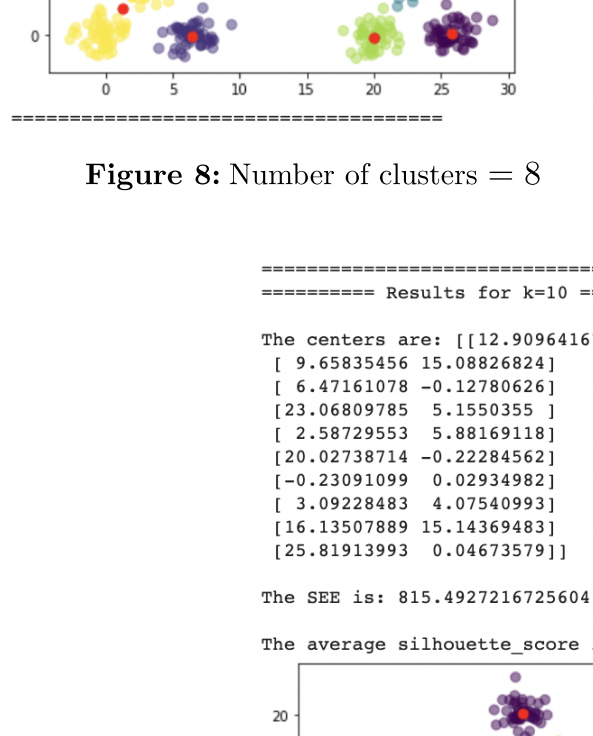


Figure 6: Number of clusters = 6

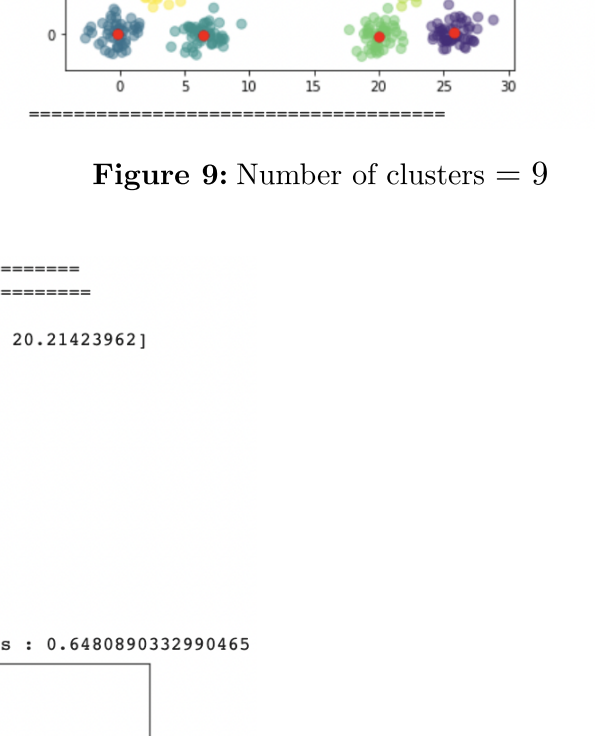


Figure 7: Number of clusters = 7

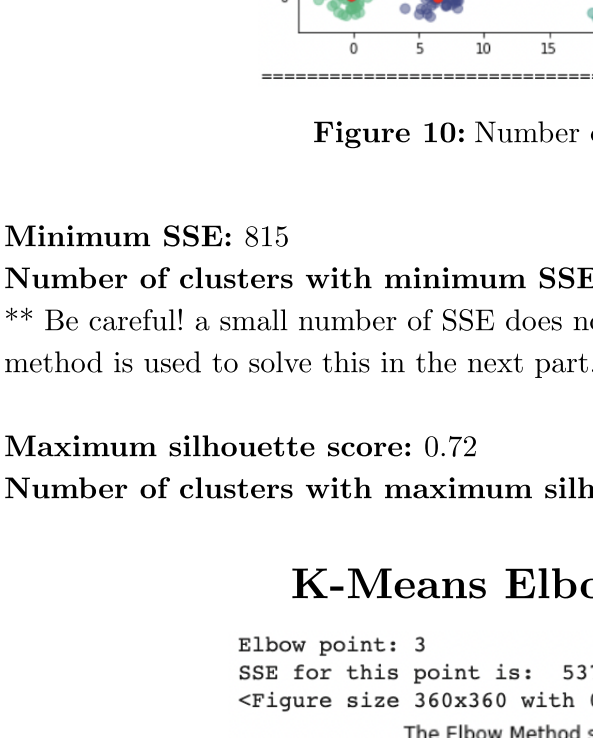


Figure 8: Number of clusters = 8

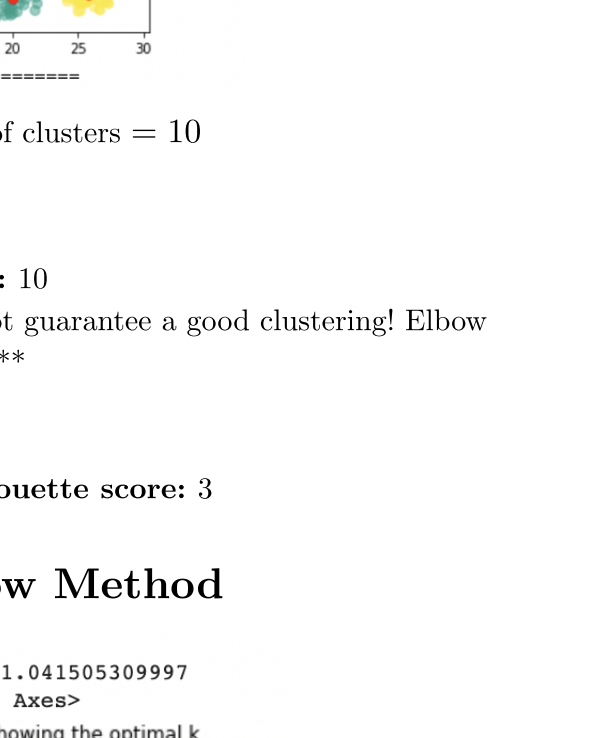


Figure 9: Number of clusters = 9

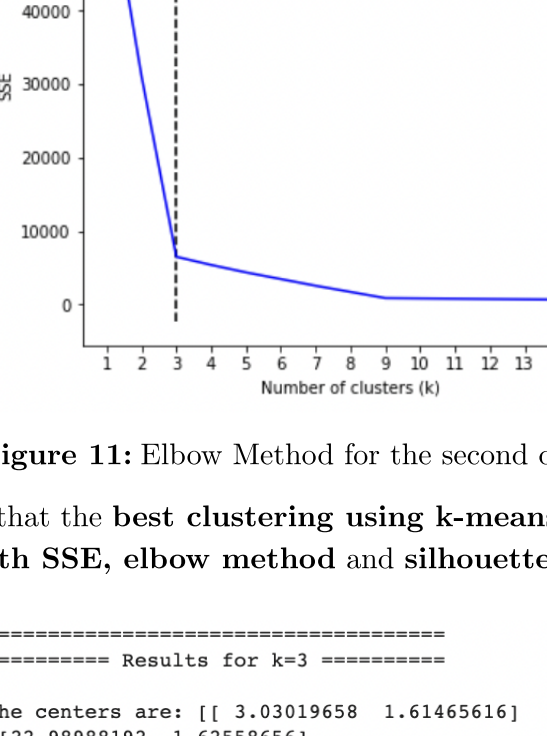


Figure 10: Number of clusters = 10

Minimum SSE: 815

Number of clusters with minimum SSE: 10

** 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: 0.72

Number of clusters with maximum silhouette score: 3

K-Means Elbow Method

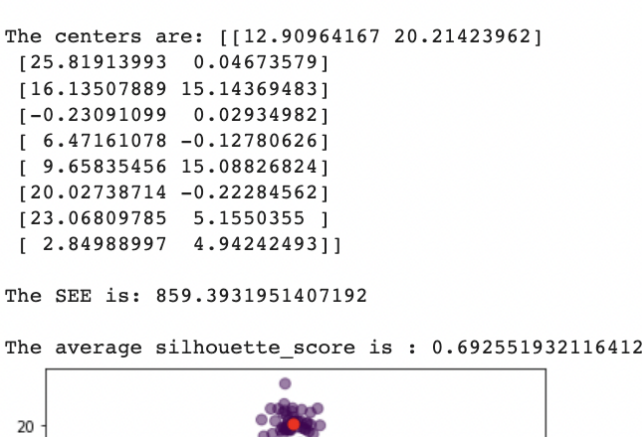


Figure 11: Elbow Method for the second dataset

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

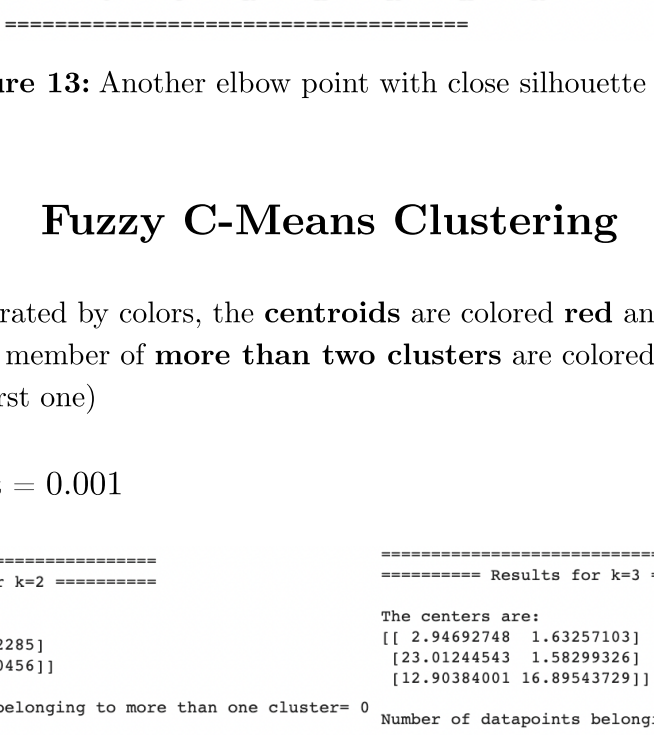


Figure 12: Best clustering for the second dataset using SSE, elbow method and silhouette score

It is worth mentioning we have **another elbow point** at $k = 9$:

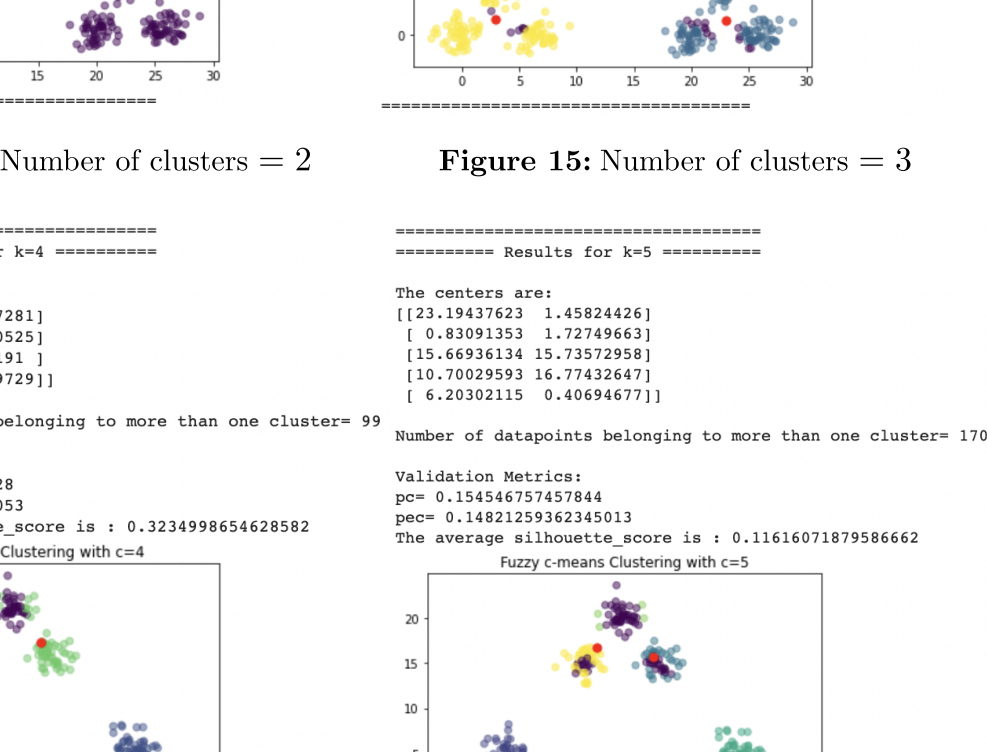


Figure 13: Another elbow point with close silhouette score

Fuzzy C-Means Clustering

Clusters are separated by colors, the centroids are colored red and the data points that are a member of **more than two clusters** are colored purple.

(except for the first one)

cutoff_ coefficient = 0.001

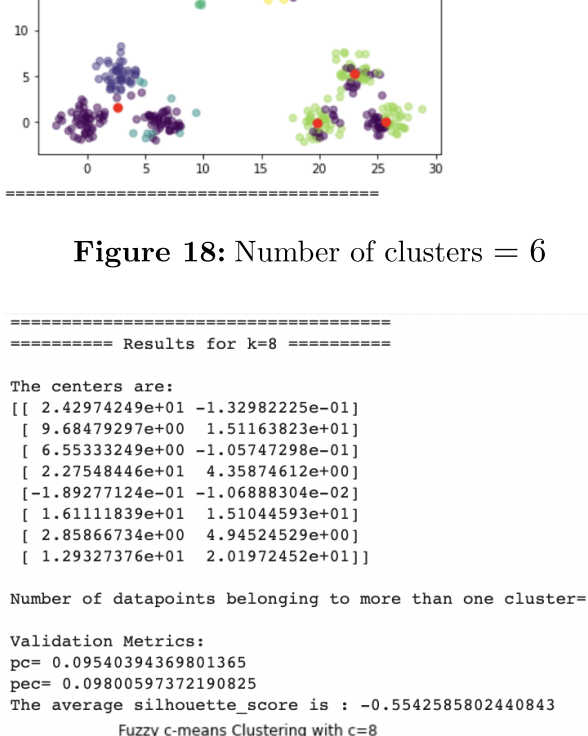


Figure 14: Number of clusters = 2

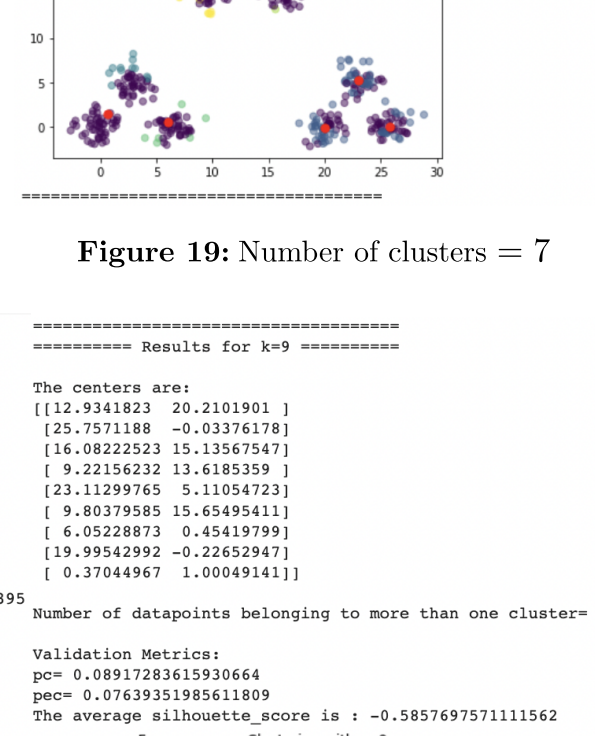


Figure 15: Number of clusters = 3

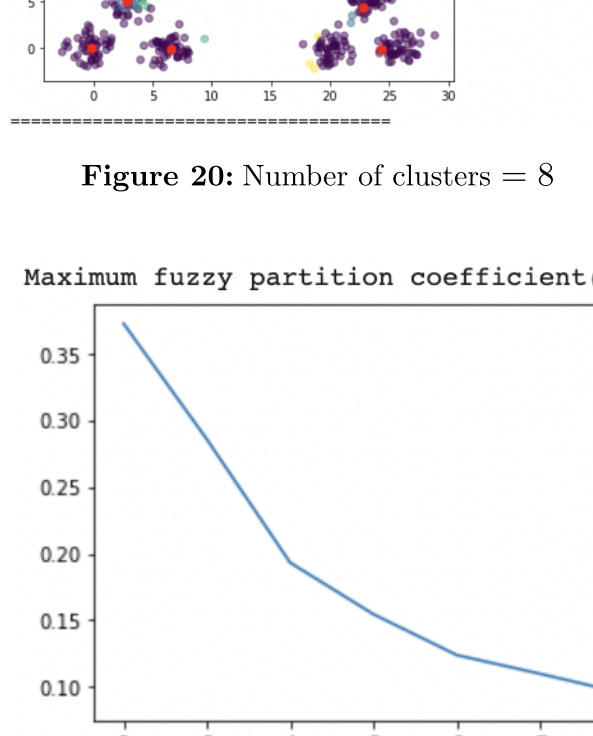


Figure 16: Number of clusters = 4

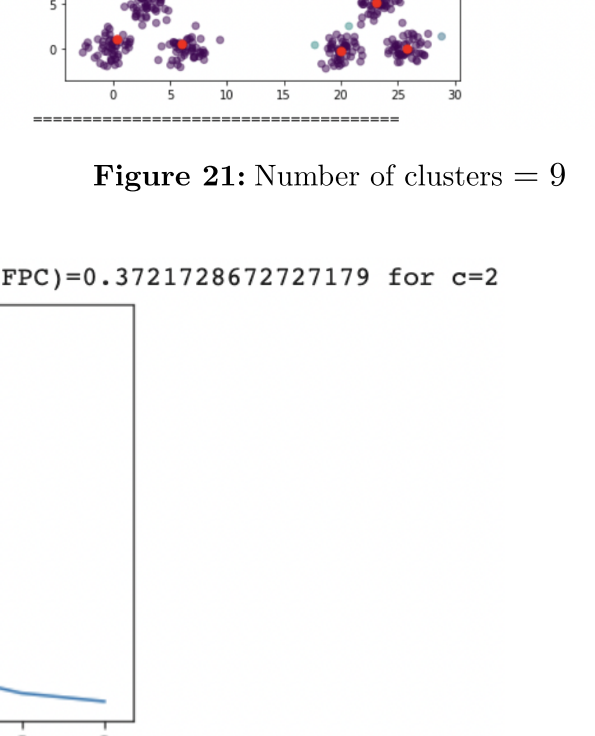


Figure 17: Number of clusters = 5

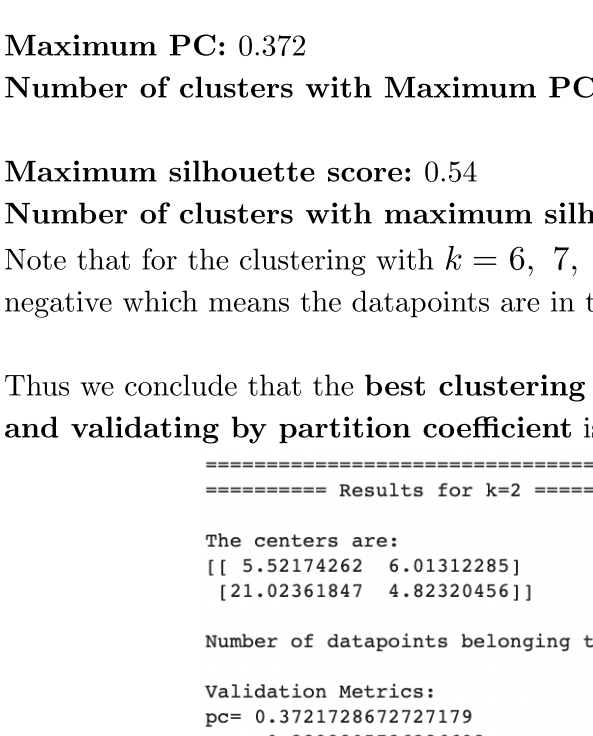


Figure 18: Number of clusters = 6

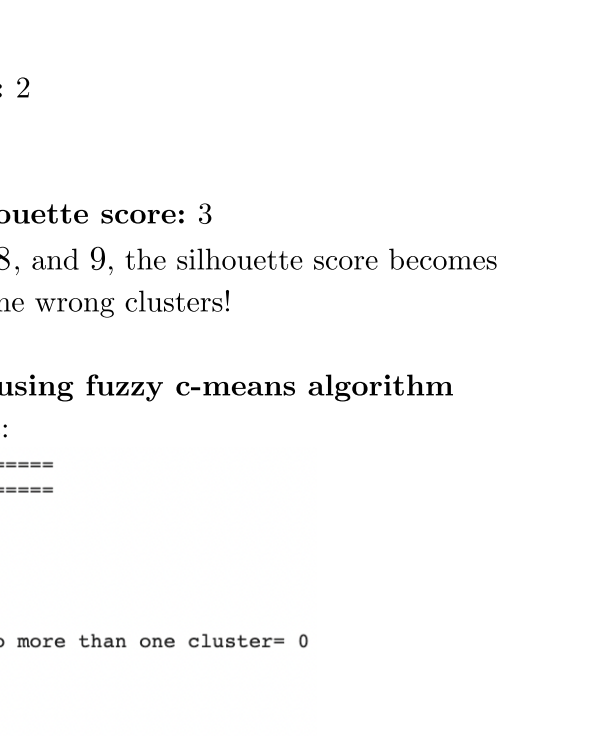


Figure 19: Number of clusters = 7

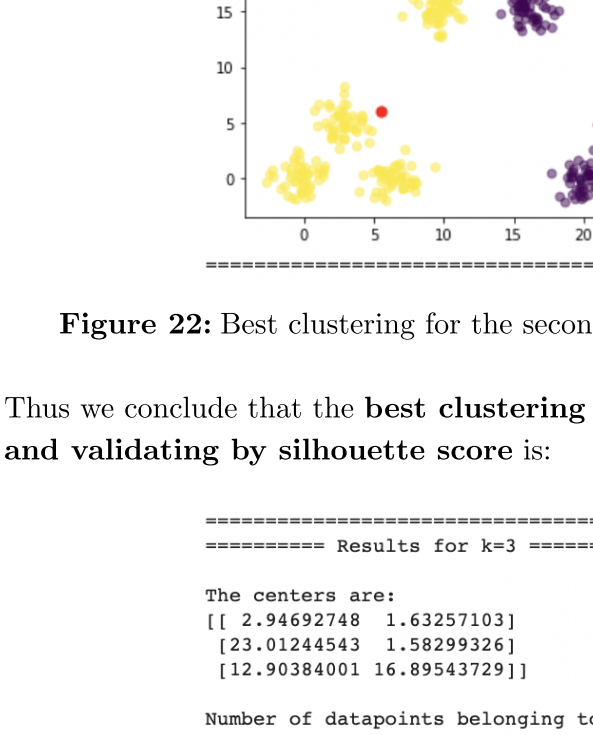


Figure 20: Number of clusters = 8

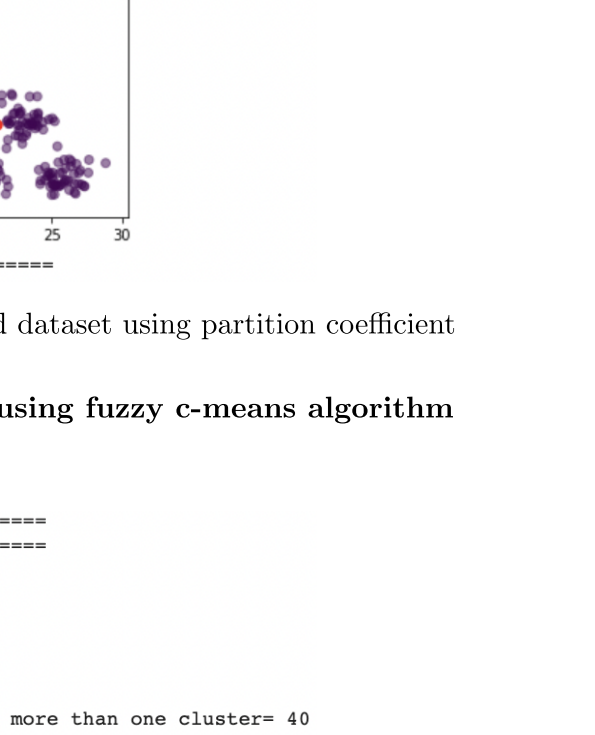


Figure 21: Number of clusters = 9

Maximum fuzzy partition coefficient(FPC)=0.3721728672727179 for c=2

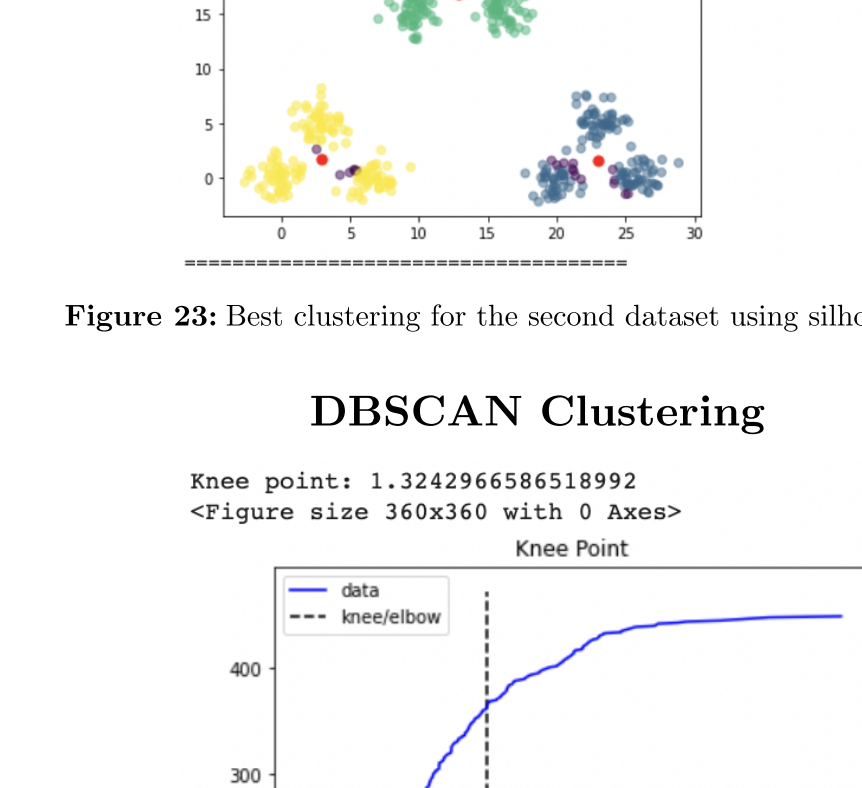


Figure 24: Best eps for DBSCAN using kNN algorithm

As explained in the previous section, we find the best *eps* parameter for the DBSCAN using the kNN algorithm. Thus the **best eps is 1.324**.

Then we try a range of numbers on the *min_samples* parameter and maximize the silhouette score. Doing so we achieve the below clustering:

Purple dots are the noise.

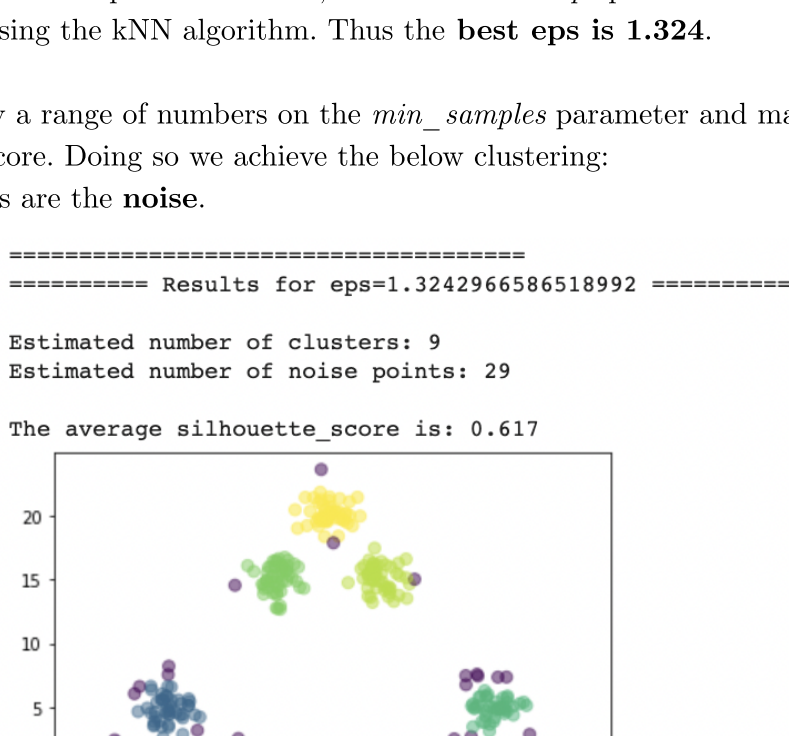


Figure 22: Best clustering for the second dataset using partition coefficient

Thus we conclude that the **DBSCAN algorithm succeeds** wonderfully to cluster this dataset. Which is probably because it is density-based.