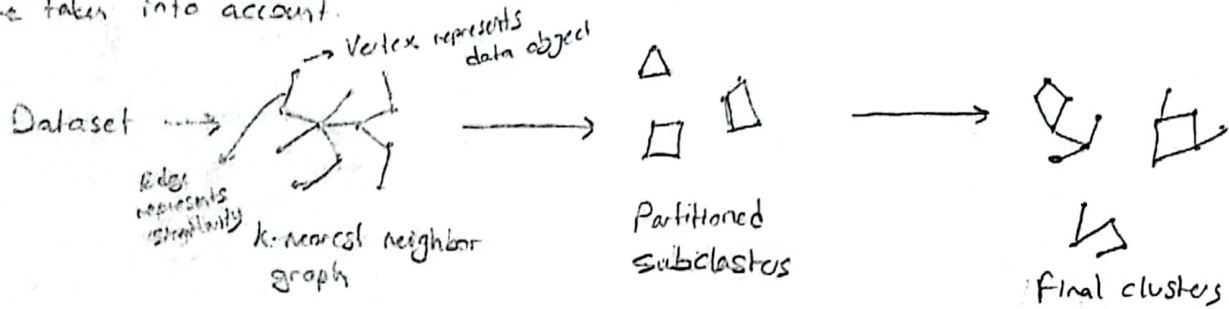


a) Chameleon is a hierarchical clustering algorithm that uses dynamic modeling to determine the similarity between pairs of clusters. It based on how well data points connected within a cluster and how close clusters are to each other. Two clusters are merged if they have high interconnectivity and they are close. So, chameleon does not rely on a fixed model provided by user; instead, it dynamically determines cluster similarity. It adapts to the internal characteristic of the clusters being merged. This adaptive approach helps discover natural and homogenous clusters.

It uses k-nearest neighbor graph method to create a sparse graph. Then, it employs a graph partitioning algorithm to divide k-nearest-neighbor graph into small subclusters. Subsequently, chameleon utilizes an agglomerative hierarchical clustering algorithm. It merges subclusters based on their similarity. To assesses the similarity between each pair of subclusters, their relative interconnectivity and relative closeness are taken into account.



Relative Interconnectivity of two clusters is calculated as;

$$RI(C_1, C_2) = \frac{|EC\{C_1, C_2\}|}{\frac{1}{2}(|ECC_1| + |ECC_2|)}$$

where $EC\{C_1, C_2\}$ is the edge cut,

ECC_x is the minimum sum of the cut edges that partition C_x into two roughly equal parts.

Relative closeness of two cluster is calculated as;

$$RC(C_1, C_2) = \frac{\overline{SECC}_{\{C_1, C_2\}}}{\frac{|C_1|}{|C_1| + |C_2|} \overline{SECC}_1 + \frac{|C_2|}{|C_1| + |C_2|} \overline{SECC}_2}$$

where $\overline{SECC}_{\{C_1, C_2\}}$ is the average weight of the edges that connect vertices in C_1 to vertices in C_2 ,

\overline{SECC}_x is the average weight of the edges that belong to min-cut bisector of cluster C_x .

```

1 procedure Chameleon Clustering (data-set, k, threshold)
2     Graph  $\leftarrow$  K-Nearest-Neighbor (data-set, k) // Create a graph using knn.
3     Clusters  $\leftarrow$  hMetis(G) // Partition the graph to clusters using any partition alg.
4     return Agglomerative Clustering (Clusters, threshold) // Apply Agglomerative Clustering:
5 endprocedure

1 procedure K-Nearest-Neighbor (data-set, k):
2     Graph  $\leftarrow$  {}
3     for i=0 to number-of-points:
4         Graph.addVertex(i) // Add each point to graph as vertex.
5     endfor
6     for i=0 to number-of-points:
7         Distances  $\leftarrow$  {}
8         for j=0 to number-of-points: // Calculate distance of each point with point i.
9             Distances[j]  $\leftarrow$  euclidian-distance (data-set[i], data-set[j])
10        endfor
11        SortedDistances  $\leftarrow$  sort (Distances.values()) // Sort the distances.
12        for vertex in SortedDistances[1:k+1]
13            Graph.addEdge (i, vertex) // Create edge for closest k neighbor.
14        endfor
15    endfor
16    return Graph
17 endprocedure.

1 procedure Agglomerative Clustering (Clusters, threshold):
2     while (number-of-clusters > 1): // Continue until there is no clusters to merge.
3         any-cluster  $\leftarrow$  false
4         for i=0 to number-of-clusters:
5             for j=0 to number-of-clusters: // Calculate each clusters similarities.
6                 if similarity (Clusters[i], Clusters[j]) > threshold: // if similarity bigger
7                     Clusters.append (Clusters[i] + Clusters[j]) // than threshold
8                     Clusters.remove (Clusters[i]) // merge the clusters
9                     Clusters.remove (Clusters[j])
10                any-cluster  $\leftarrow$  true
11            endif
12        endfor
13    endfor
14    if any-cluster is false: // if there is no cluster to merge, break the loop.
15        break
16    endif
17    endwhile
18    return Clusters
19 endprocedure

1 procedure similarity (x, y):
2     return RelativeInterconnectivity (x, y) * Relative Closeness (x, y)
3 endprocedure

```


b) Chameleon clustering offers several advantages that make it a versatile algorithm for cluster analysis. Firstly, it is not sensitive to the shape and size of clusters, allowing it to effectively identify clusters with irregular shapes and varying sizes. Another notable advantage is its ability to autonomously determine the number of clusters without requiring user input, alleviating the need for prior knowledge about the dataset. Moreover, Chameleon incorporates both relative closeness and relative interconnectivity concurrently, enhancing its robustness against noise and outliers and contributing to the overall quality of the clusters.

However, Chameleon does come with several drawbacks. One significant disadvantage is its computational expense, as it sequentially applies the K-Nearest Neighbor, Partitioning Algorithm, and Agglomerative Clustering Algorithm, coupled with pairwise distance calculations. Additionally, finding optimal parameters such as the threshold for merging clusters and the value of k for the Nearest Neighbor Algorithm can be a non-trivial task, potentially requiring extensive experimentation. Lastly, the choice of distance measure method can significantly impact the success of the algorithm, adding an extra layer of complexity to its implementation and performance.

Algorithm	Time	Space	Performance
K-Means	$O(knd)$	$O(kd)$	Works well for spherical shapes
DBSCAN	$O(n^2)$	$O(n)$	It is efficient for clusters with arbitrary shapes and sizes. It is sensitive to density of data points.
Hierarchical	$O(n^3)$	$O(n^2)$	It can capture clusters at different scales. It is sensitive to noise and outliers.
BIRCH	$O(n)$		Fast algorithm with good quality of clustering. However, it does not produce natural clusters.
Chameleon	$O(n^2)$	$O(n^2)$	It can identify clusters with different shapes and sizes. It provides good quality clusters.

c) As we stated in question a, Chameleon Clustering consist of 3 main steps and its time complexity depends on these steps. We can calculate them separately to find overall time complexity.

Start with K-Nearest Neighbor procedure:

The most time consuming process in this procedure is calculating distance between each data point. Nested for loop takes $O(n^2)$ if we assume that calculating distance between two point has $O(1)$ time complexity.

Continue with hMetis procedure:

Its time complexity is $O(V+E)$.

End with Agglomerative procedure.

Agglomerative clustering procedure contains 3 nested loops which run n time if partitioning algorithm partitioned until each cluster contains only one point. So, the overall time complexity of this procedure is $O(n^3)$, if we assume that appending an element to cluster or removing an element from it has $O(1)$ time complexity.

To sum up, Chameleon Clustering Algorithm has $O(n^2 + V + E + n^3)$ time complexity. However, n^3 dominates the time complexity. So, overall worst time complexity of this clustering algorithm is $O(n^3)$.

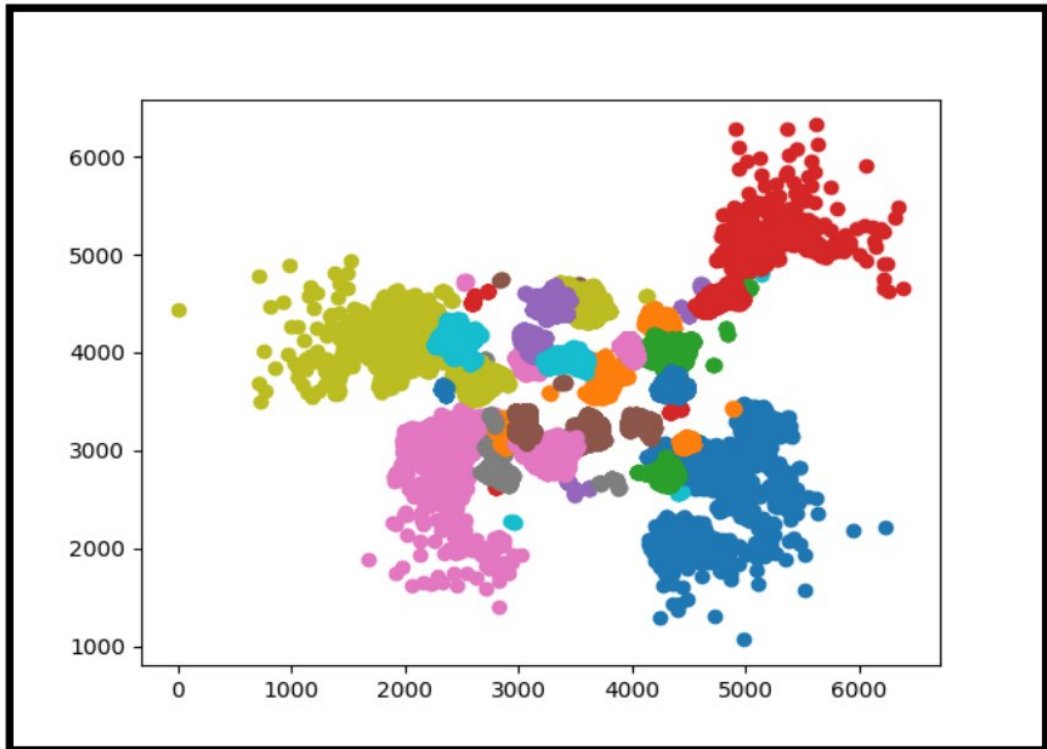
Exclude because we are not responsible

Note: Time complexity of an algorithm can vary according to implementation. Above results is calculated as my pseudocode for this algorithm.

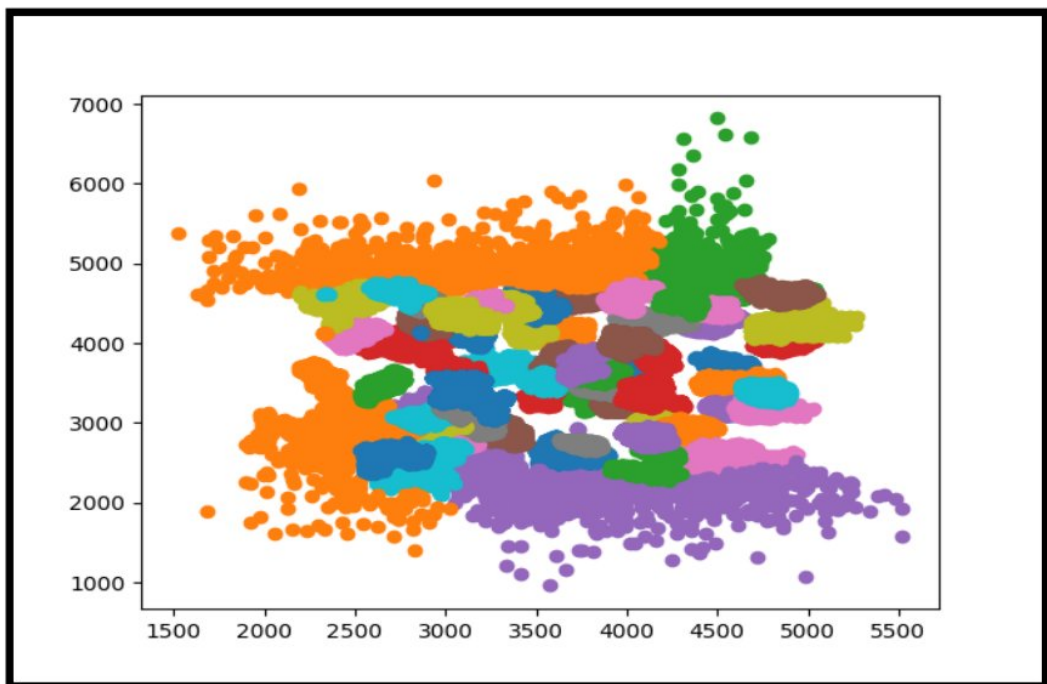
d) This is my video link: <https://clipchamp.com/watch/Utk6BNyl2hF>

e) I applied this euclidian tests to 105.000 2D points.

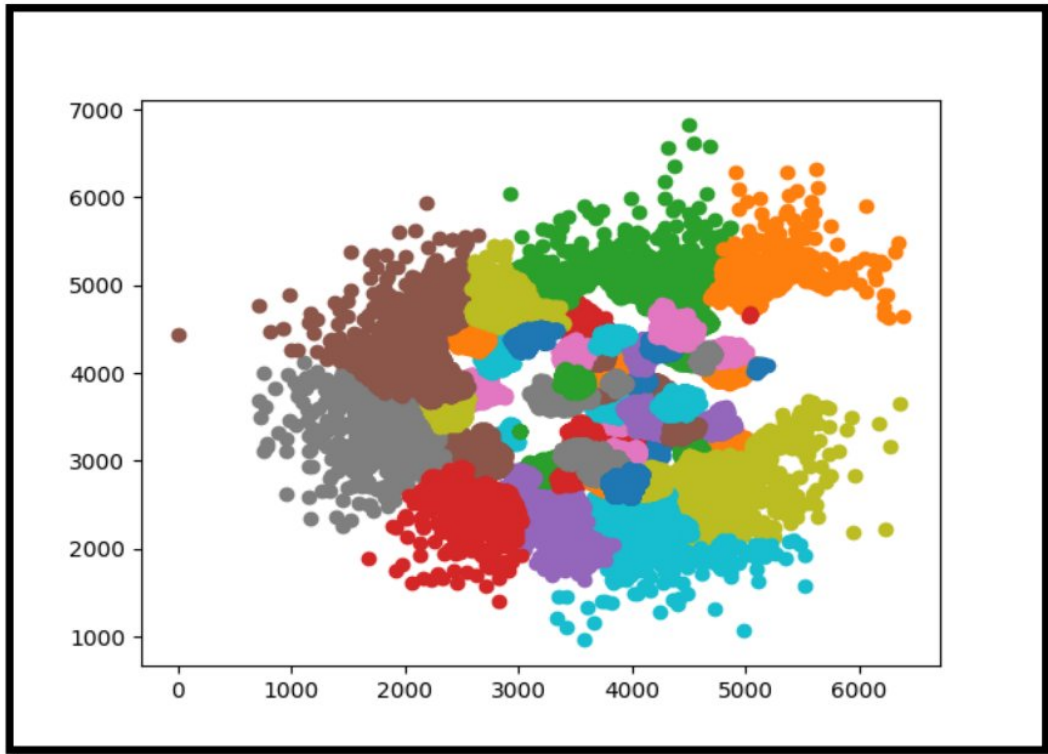
k = 3, threshold = 0.7



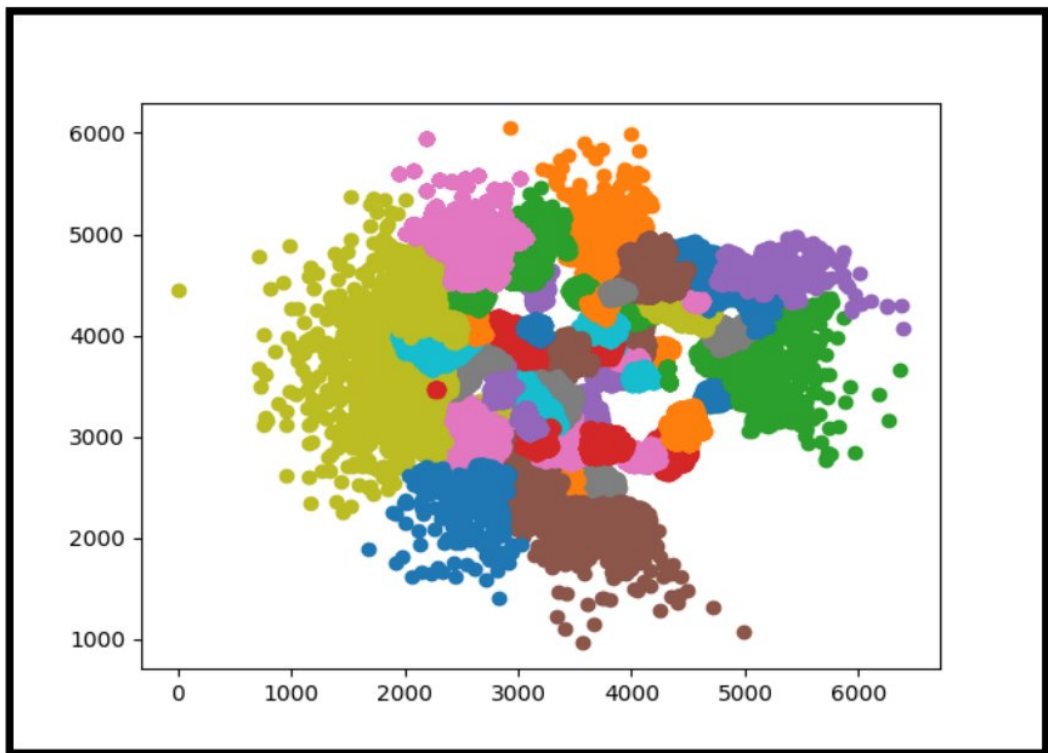
k = 4, threshold = 0.7



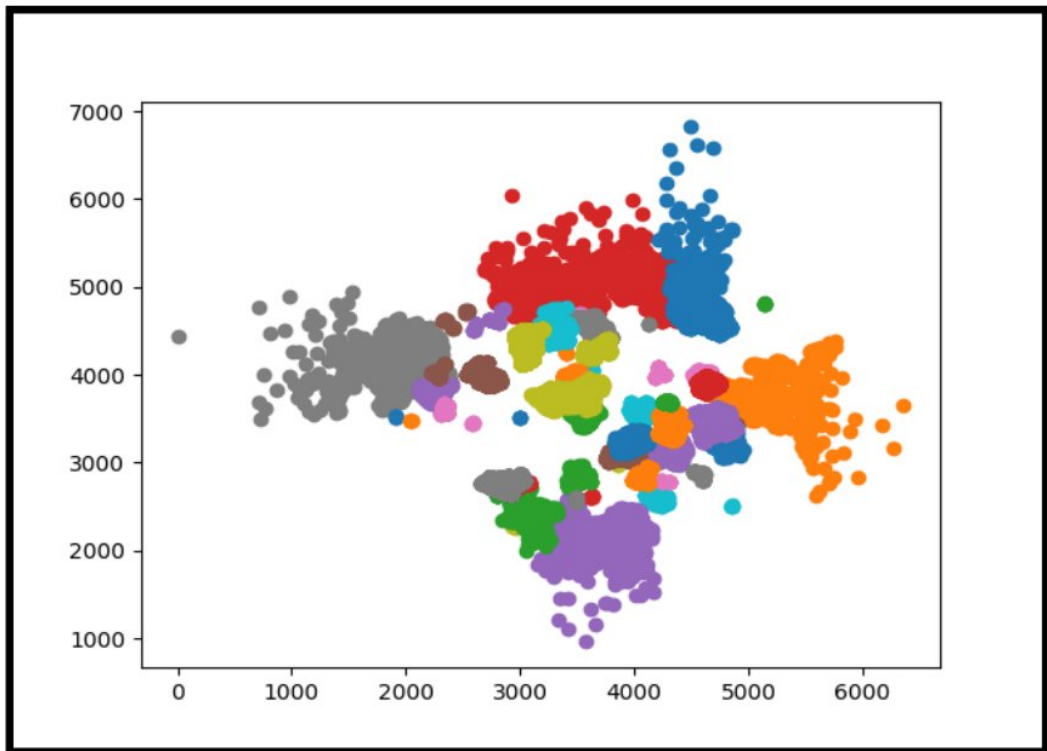
k = 5, threshold = 0.7



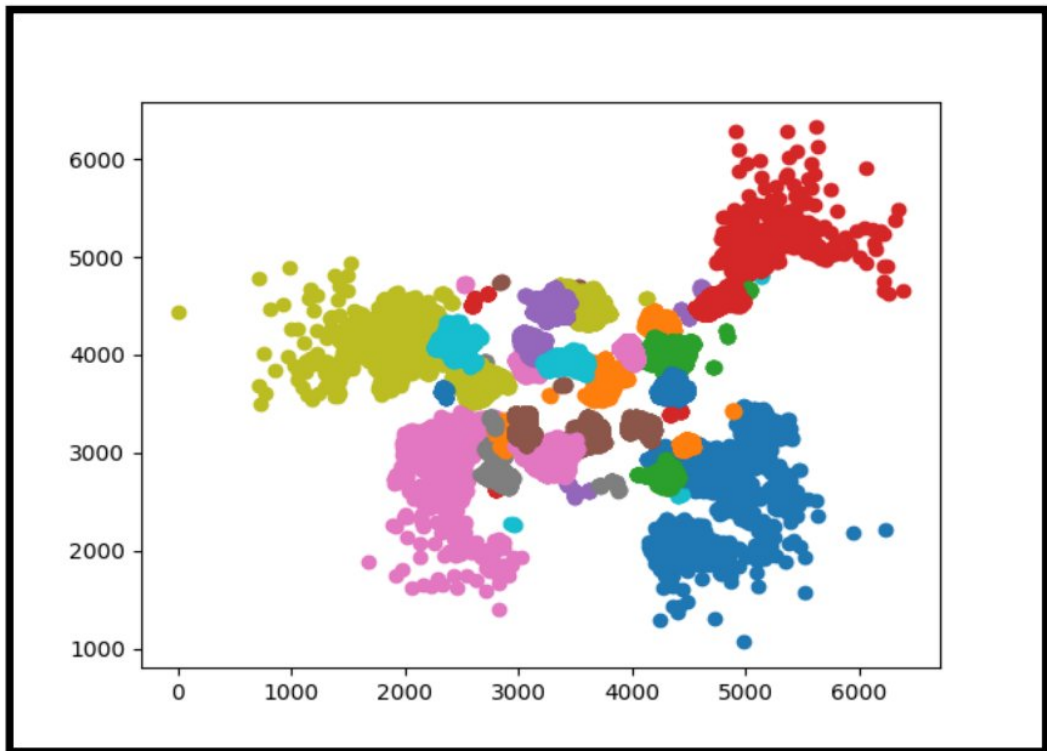
k = 6, threshold = 0.70



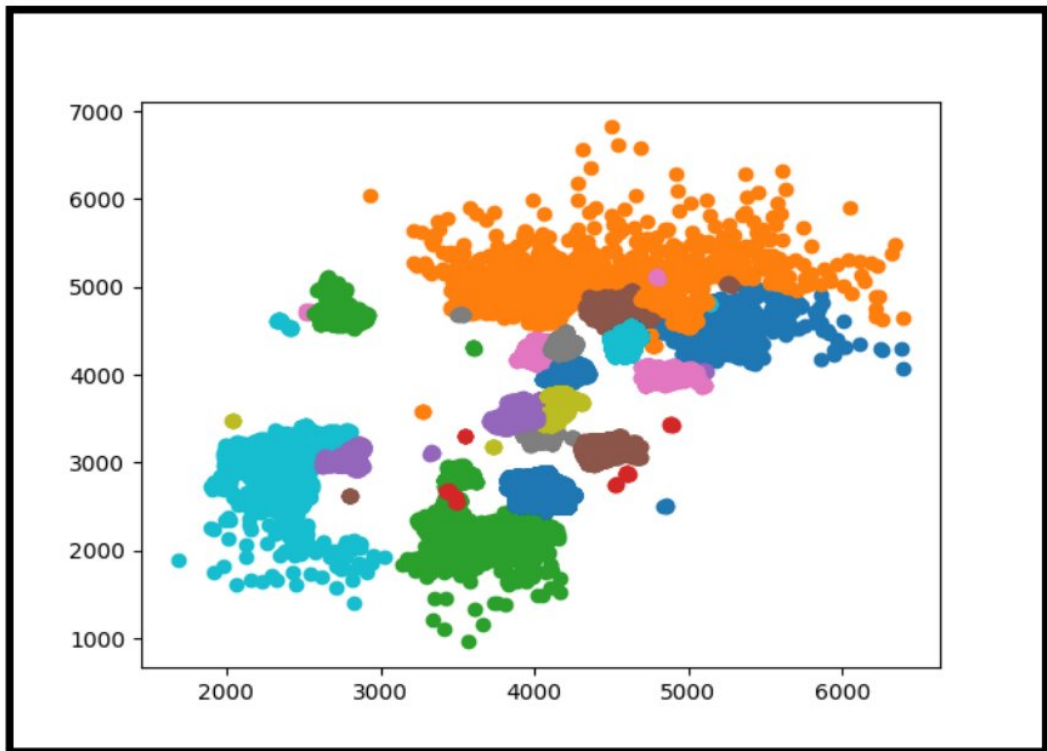
k = 3, threshold = 0.75



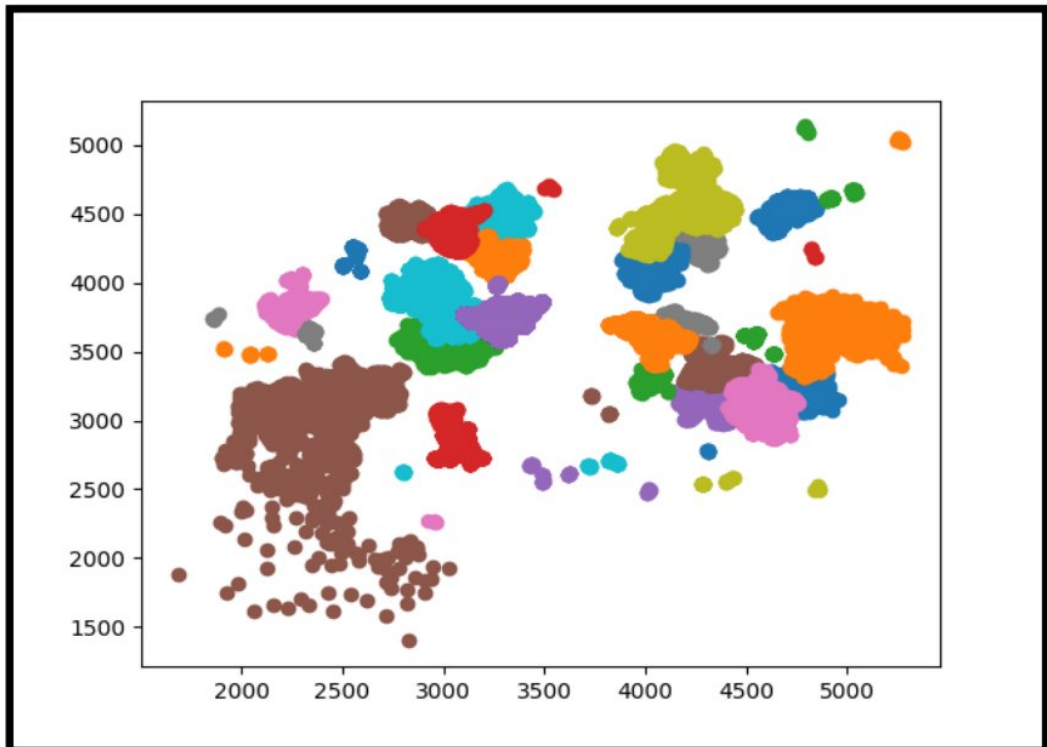
k = 3, threshold = 0.7



k = 3, threshold = 0.65

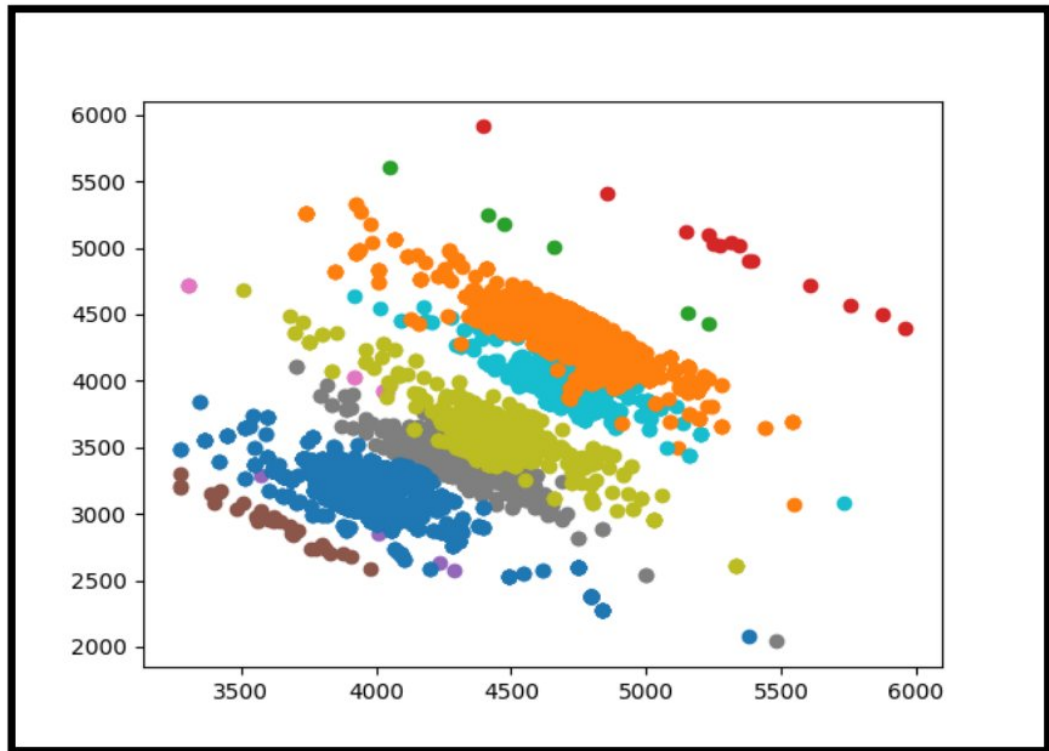


k = 3, threshold = 0.6

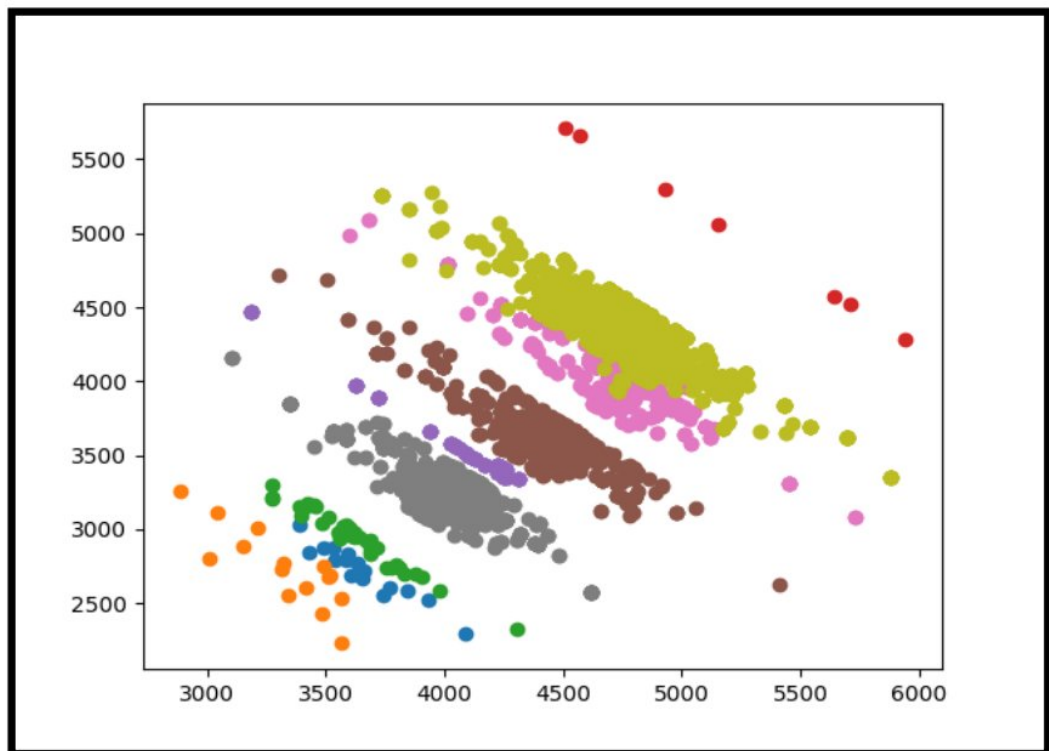


f) I applied this manhattan tests to 10.000 2D points.

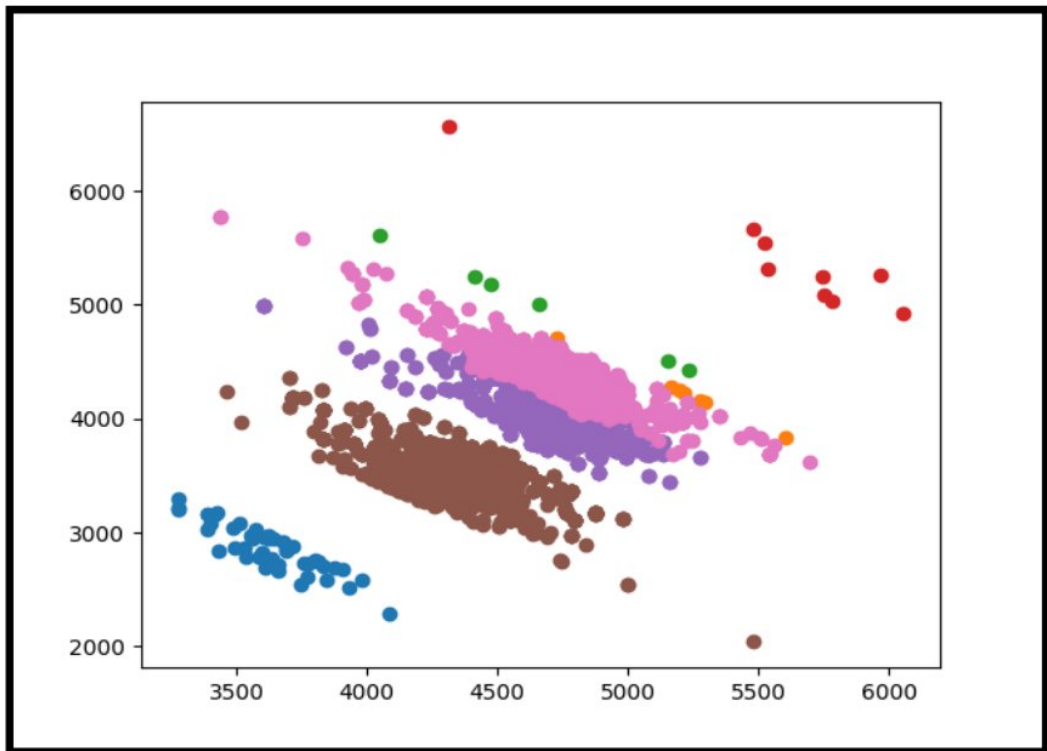
k = 3, threshold = 0.7



k = 5, threshold = 0.7



k = 3, threshold = 0.6



k = 4, threshold = 0.65

