

Data Mining (CS524)

Lab Assignment 3

Clustering

**Submitted by:
Paras Kumar
2016CSB1047**

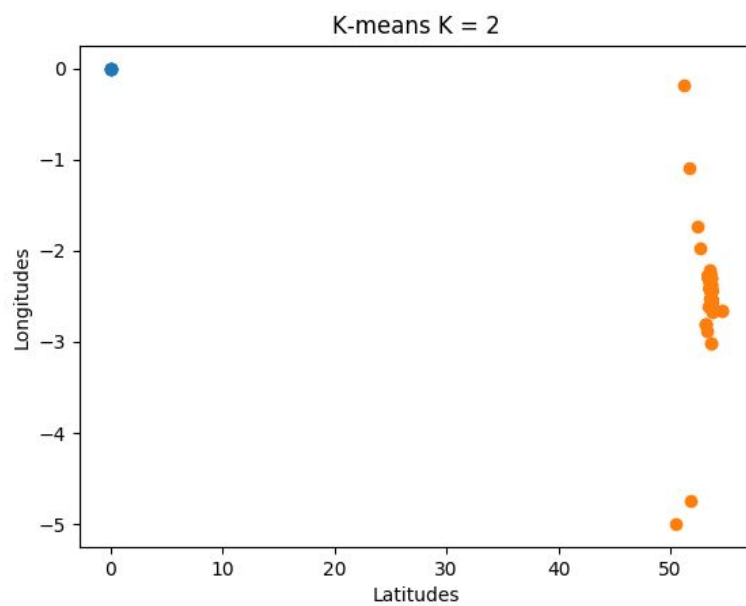
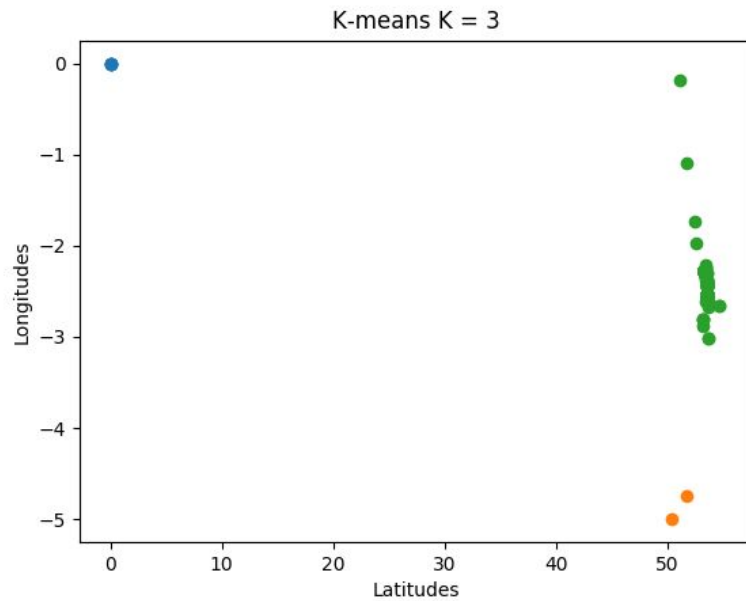
TASK 1

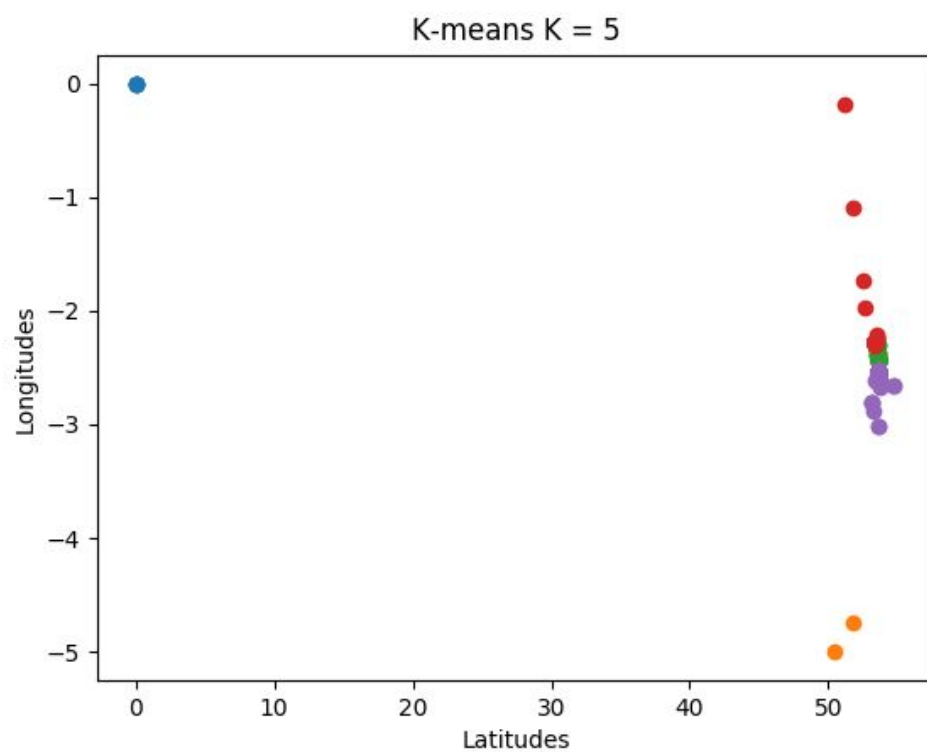
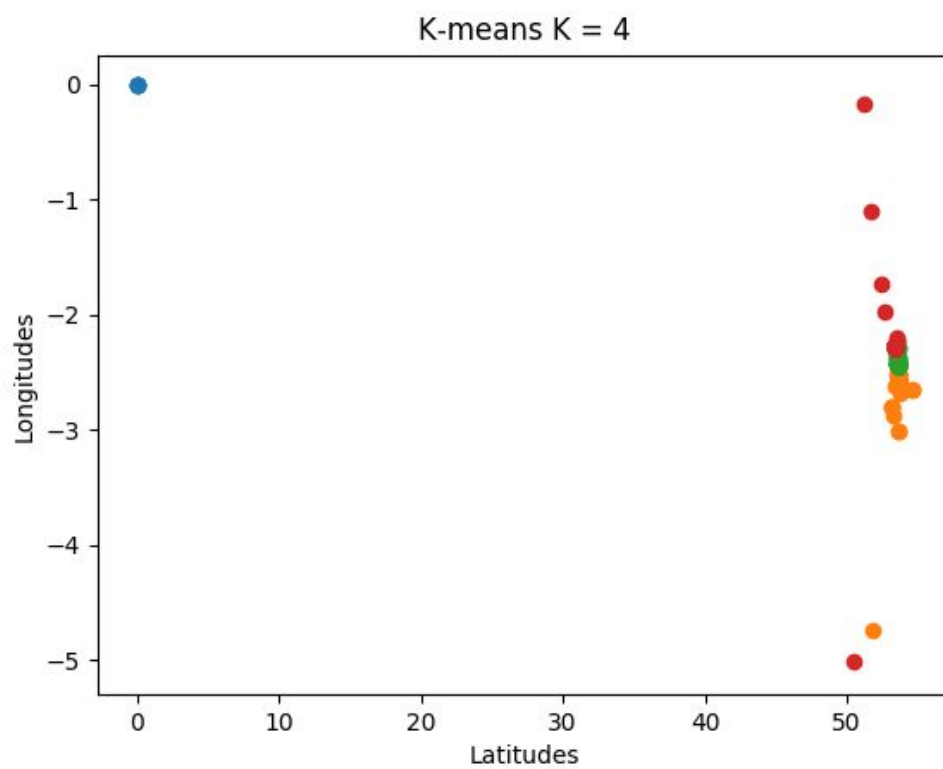
All the four clustering algorithms are implemented in python3 and submitted along with this report. I have followed the algorithms given in the book and in case of Agglomerative clustering, I used distance matrix. The code is modularized and easy to understand.

Note 1: I have used Euclidean distance in all the algorithms.

TASK 2

Clusters obtained by K-means algorithm:



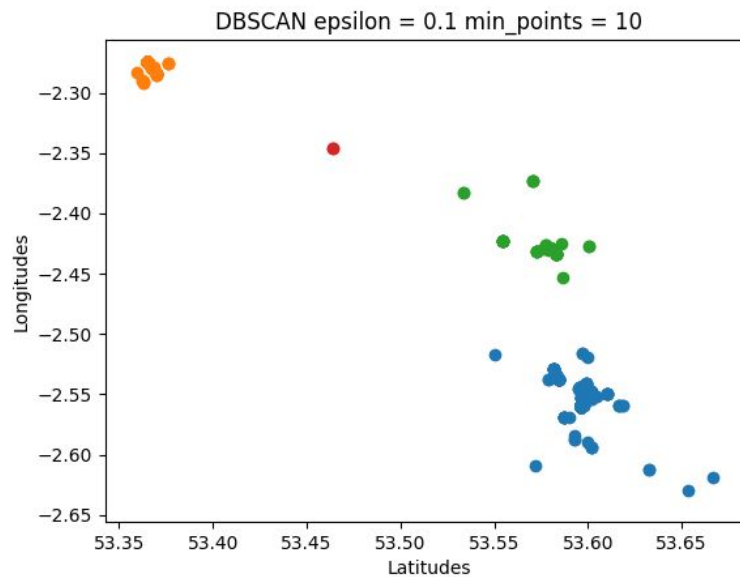
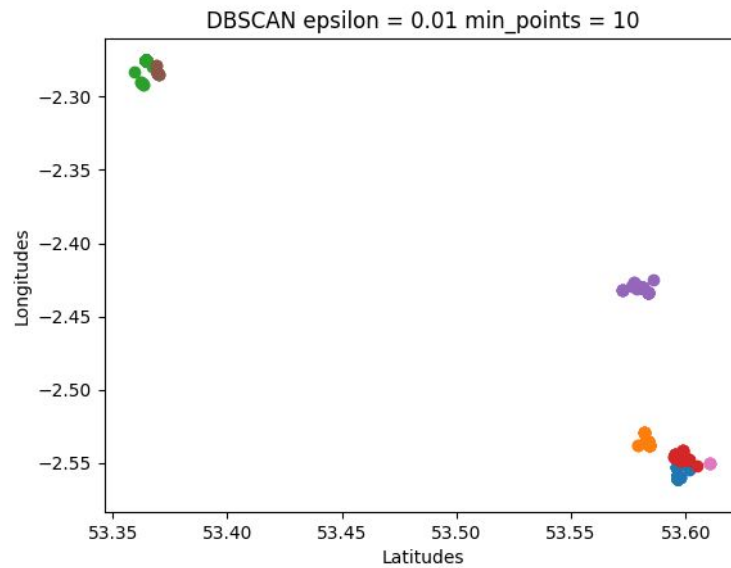


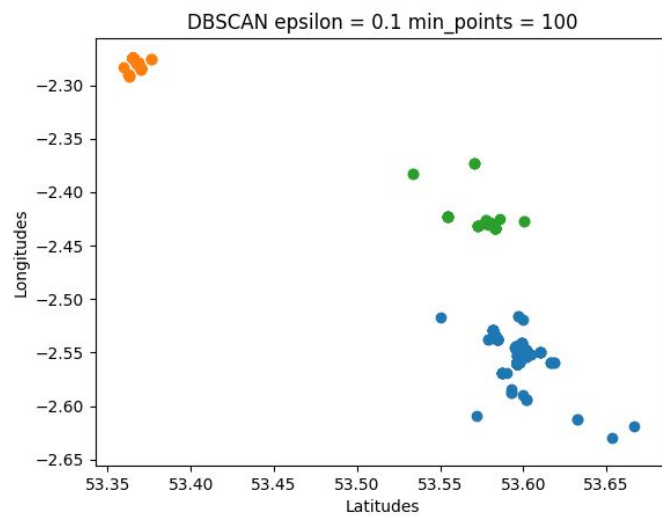
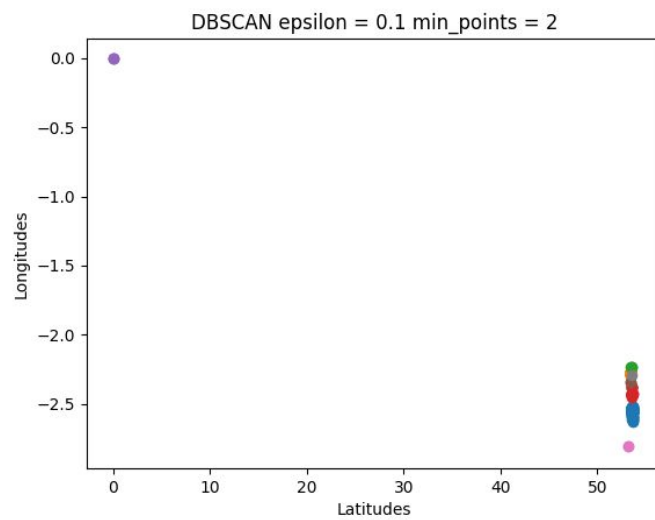
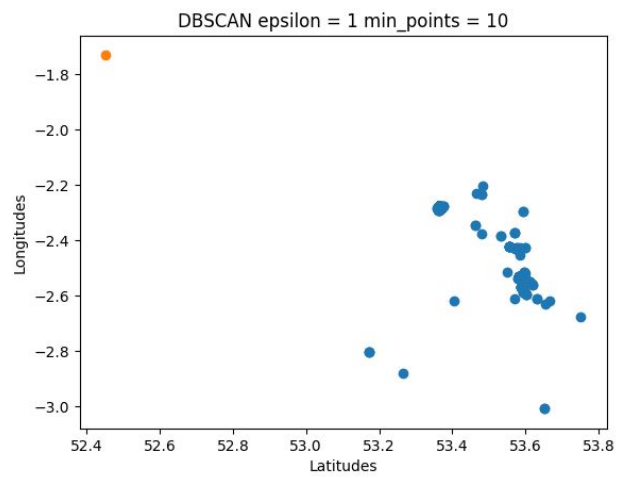
K-means Analysis

We observe that for $K = 5$ we get best clustering results with K-means. We also observe the effect of outliers on clusters and also K-means can only find convex clusters, so it leads to some errors which can be fixed by using a density based clustering algorithm like DBSCAN.

TASK 3

Clustering results with DBSCAN:





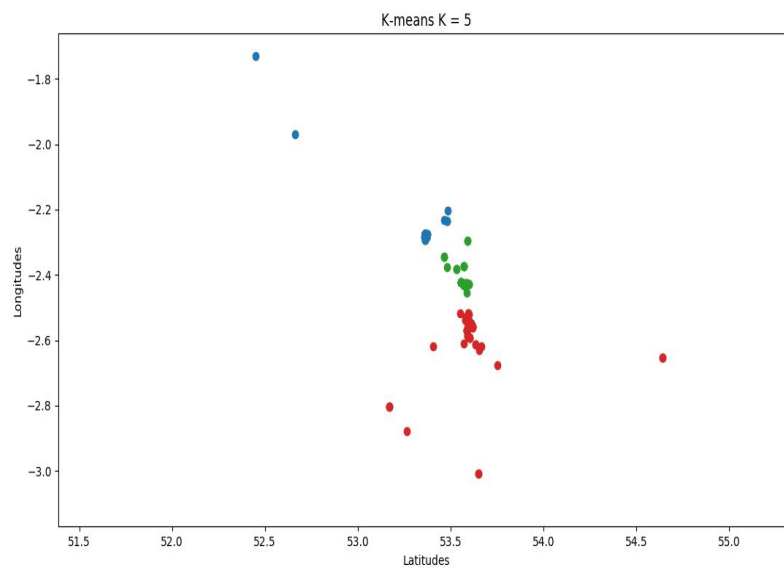
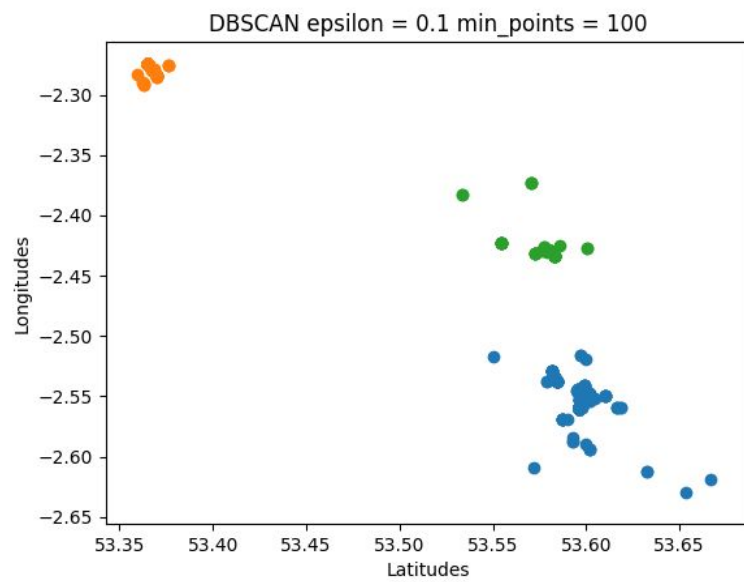
In these graphs I have not shown the noise points so that we can observe the important clusters in the data. DBSCAN produced the best results out of the four algorithms. It filters out noise points and doesn't require number of clusters as input. The best result is obtained for $\epsilon = 0.1$ and $\text{min_points} = 10$.

The clusters formed are very sensitive to epsilon values, with its increase the number of clusters decreases.

We also observe that for very low min_points value, like 2 some outliers also start forming clusters. But, increasing it from 10 to 100 leads to very little change in the clusters.

TASK 4

Comparison of DBSCAN and K-means

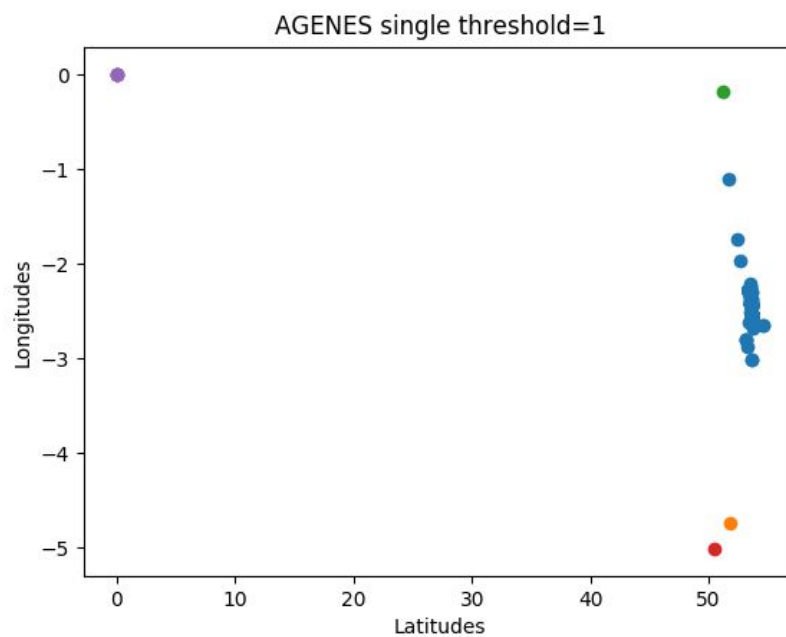
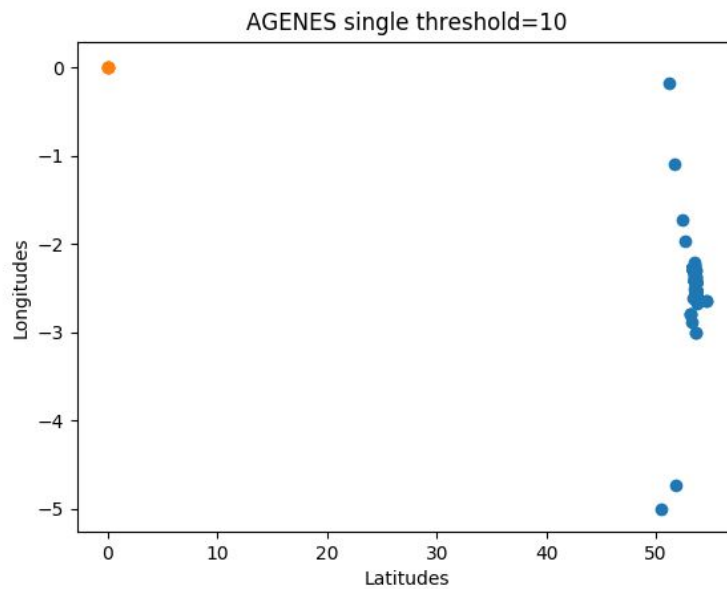


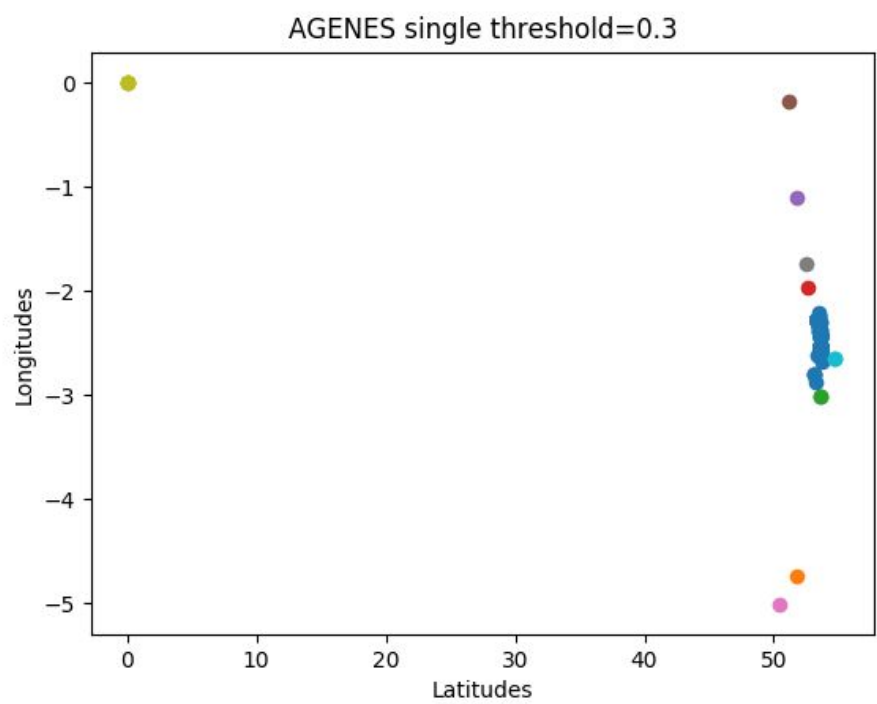
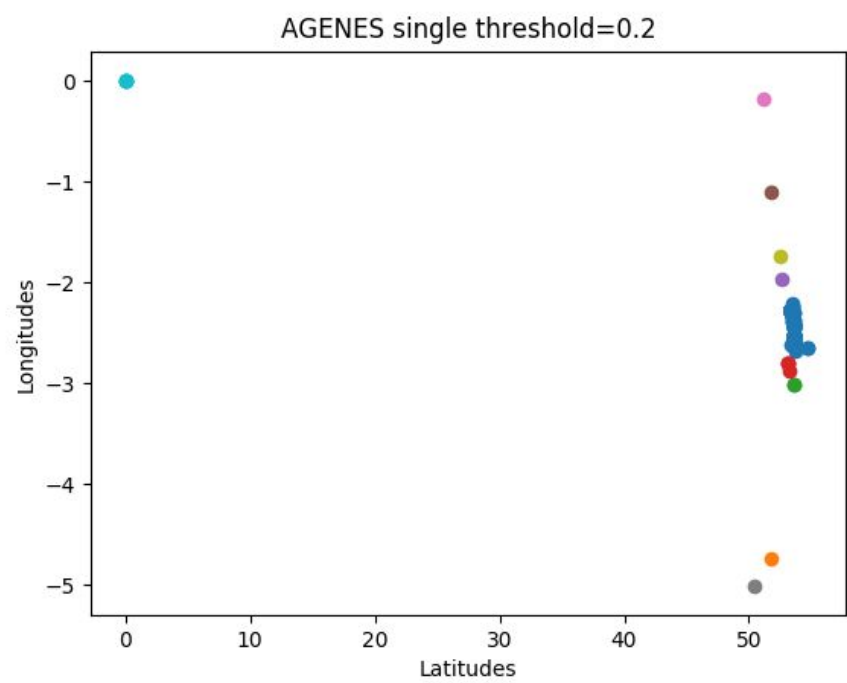
Comparison of DBSCAN and K-means

Comparing the best clusters of both the algorithms (zoomed graph in case of K-means) we observe that DBSCAN produces higher quality clusters. The clusters in case of DBSCAN are more dense and the noise points are removed, so it gives us a better visualization of the data. Also, since the clusters are not spherical so K-means fails to produce good quality clusters.

TASK 5

Single Linkage Agglomerative clustering





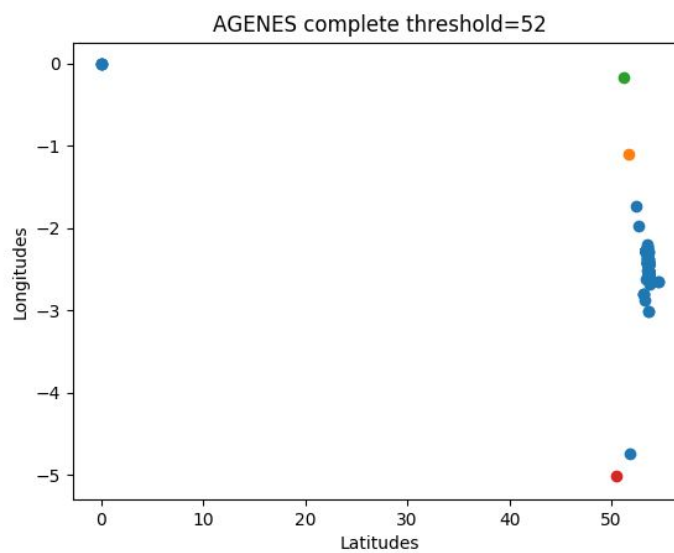
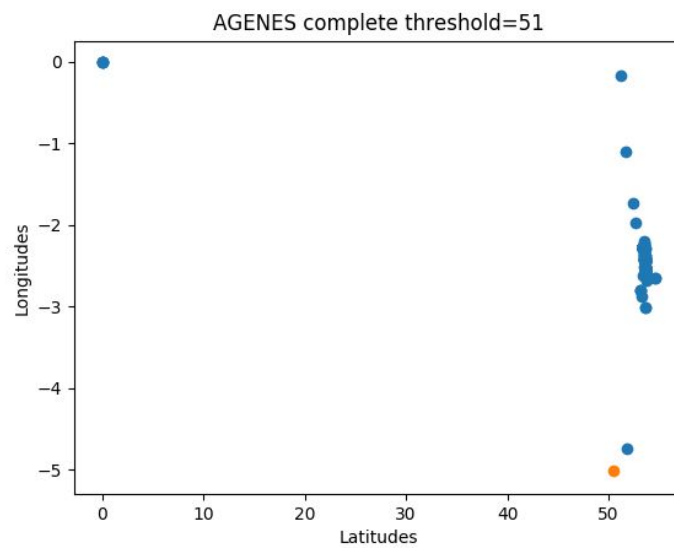
Single linkage agglomerative clustering

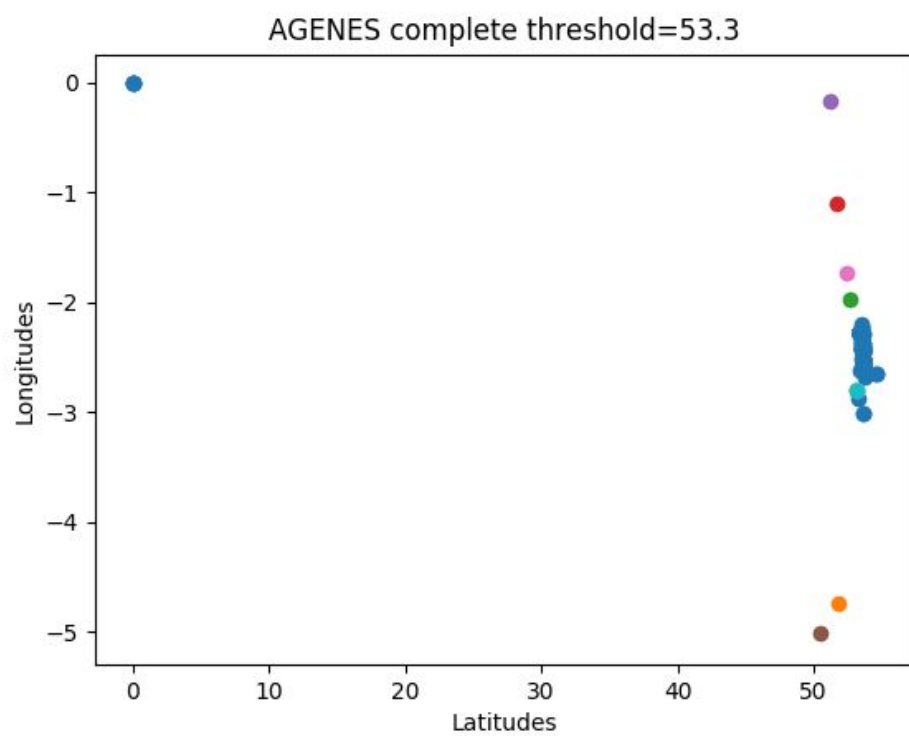
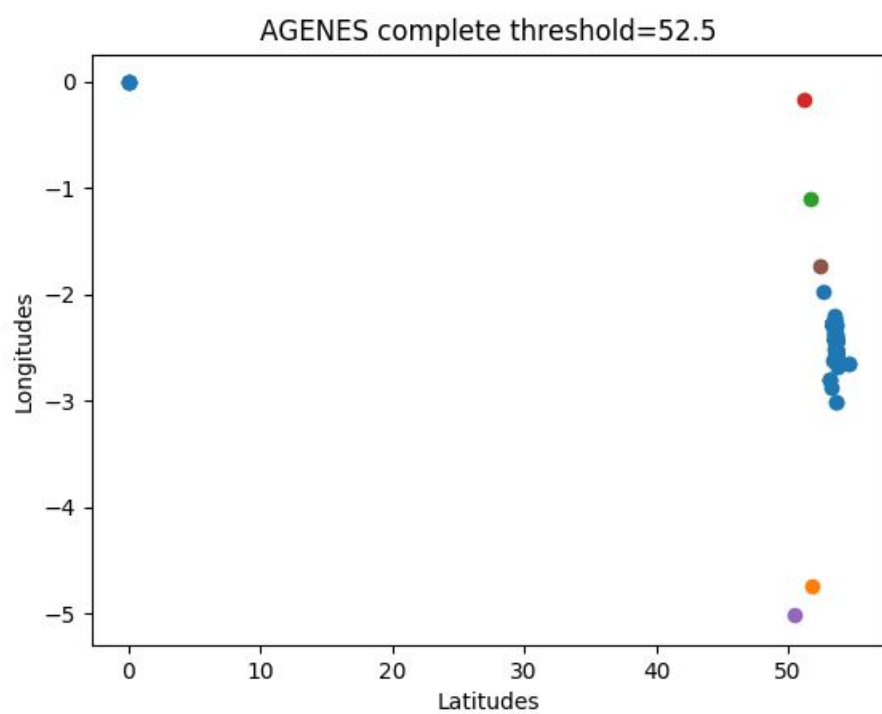
Number of cluster changing with threshold

Clusters = 25 val= 0.030657312374048208
Clusters = 24 val= 0.0341728571969048
Clusters = 23 val= 0.03426762771187951
Clusters = 22 val= 0.03609376390735594
Clusters = 21 val= 0.03852844742524672
Clusters = 20 val= 0.044564868068918434
Clusters = 19 val= 0.05290915423629441
Clusters = 18 val= 0.06429120904913865
Clusters = 17 val= 0.08029192044658032
Clusters = 16 val= 0.09890103786108798
Clusters = 15 val= 0.10341916134353342
Clusters = 14 val= 0.11174072745870499
Clusters = 13 val= 0.12170888746924169
Clusters = 12 val= 0.16583583251215558
Clusters = 11 val= 0.29610169022482846
Clusters = 10 val= 0.319532770133206
Clusters = 9 val= 0.34618331110843586
Clusters = 8 val= 0.7624995130818116
Clusters = 7 val= 0.8896737167034904
Clusters = 6 val= 0.9542529433111546
Clusters = 5 val= 1.088190669290084
Clusters = 4 val= 1.3766891173863463
Clusters = 3 val= 2.3724108890124858
Clusters = 2 val= 50.69529529983724

TASK 5

Complete Linkage Agglomerative clustering



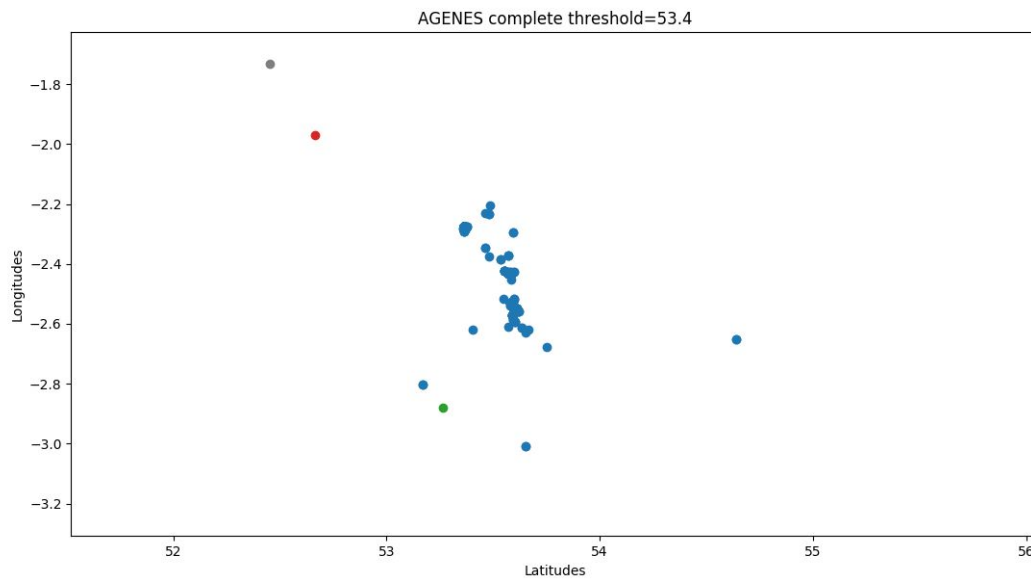
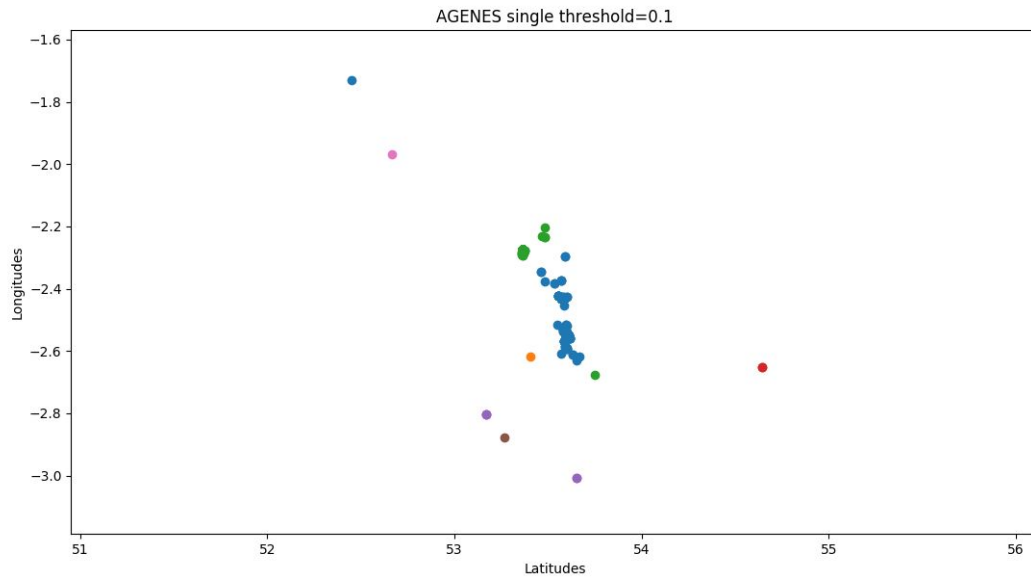


Complete linkage agglomerative clustering

Number of cluster changing with threshold

Clusters = 25 val= 53.41317201627914
Clusters = 24 val= 53.41317201627914
Clusters = 23 val= 53.41317201627914
Clusters = 22 val= 53.41317201627914
Clusters = 21 val= 53.41317201627914
Clusters = 20 val= 53.41317201627914
Clusters = 19 val= 53.4122594242068
Clusters = 18 val= 53.4122594242068
Clusters = 17 val= 53.41212886438529
Clusters = 16 val= 53.41212886438529
Clusters = 15 val= 53.41105255131923
Clusters = 14 val= 53.41105255131923
Clusters = 13 val= 53.40812883928181
Clusters = 12 val= 53.40812883928181
Clusters = 11 val= 53.34386362420635
Clusters = 10 val= 53.245603960901875
Clusters = 9 val= 53.245603960901875
Clusters = 8 val= 53.245603960901875
Clusters = 7 val= 52.70110505946175
Clusters = 6 val= 52.479841547936296
Clusters = 5 val= 52.015916287355836
Clusters = 4 val= 51.74962911351704
Clusters = 3 val= 51.16303404489612
Clusters = 2 val= 50.69529529983724

Comparison of single link and complete link



Comparison of single link and complete link

We observe that complete link gives better results since it tends to minimize the cluster diameter and the original clusters are also compact. The single link fails to identify the clusters as shown in the zoomed graphs.

In general both these algorithms are slow and performs worse than DBSCAN.

TASK 6

Since, the program for this task was running very slow, I randomly selected 1/5th of original data. Also, I observed a problem with the given distance function, it would always result in those 5 users which have min number of transactions because the distance function is not normalized. I have produced the following output following the instructions given in the assignment.

Note: DBSCAN doesn't not work well for this problem as there are very few points and they are spread apart, they all will be treated as outliers or we can interpret them as a single cluster.

Top 5 closest users and clustering of their check-in data using DBSCAN

