

Homework 4

1)

Input: k = number of clusters | datapoints = set of data points

Output: data points will be labeled with their cluster id

```
k_means(k, datapoints):
    centroids[k] //array of size k for storing centroids
    for i = 0, i < k, ++i:
        centroids[i] = random point from datapoints
    bool significant_change = true
    while(significant_change)
        significant_change = false
        for each x in datapoints
            smallest_dist = dist(x, centroids[0])
            smallest_dist_id = 0
            cluster_id = 0
            for each c in centroids
                if(dist(x, c) < smallest_dist)
                    smallest_dist = dist(x, c)
                    smallest_dist_id = cluster_id
                cluster_id += 1
            if(x.label != smallest_id)
                significant_change = true
            x.label = smallest_id
    end
```

2)

Note: Will be using a KDTree for quick nearest neighbour lookup

Input: k = number of clusters | datapoints = set of data points

Output: data points will be labeled with their cluster id

```
agglomerative(k, datapoints):
    kdtree kdt = kdtree(datapoints)
    while kdt.root.num_children > k:
        node1 = datapoints[0]
        node2 = datapoints[1]
        shortest_dist = dist(node1, node2)
        for each x in datapoints
            y = kdt.nearestNeighbour(x)
            if(dist(x, y) < shortest_dist)
                node1 = x; node2 = y
                shortest_dist = dist(x, y)
        z = merge(node1, node2)
        kdtree.remove(node1)
```

```

        kdtree.remove(node2)
        kdtree.add(z)
    end

```

3)

Input: k = number of clusters | datapoints = set of data points
 Output: data points will be labeled with their cluster id

```

divisive(k, datapoints):
    num_clusters = 2
    k_means(2, datapoints) //initial split of the system
    while num_clusters < k
        x = random cluster id
        particular_cluster = datapoints.cluster(x)
        k_means(2, particular_cluster)
    end

```

4)

Algorithm 18.1: Density-based Clustering Algorithm

```

dbscan ( $\mathcal{D}$ ,  $\epsilon$ ,  $minpts$ ) :
1  foreach  $x \in \mathcal{D}$  do
2      Compute  $N_\epsilon(x)$ 
3      Classify  $x$  as core, border, or noise
4   $id = 0$ 
5  foreach  $x \in \mathcal{D}$ , such that  $x$  is core and unmarked do
6       $id = id + 1$ 
7      DensityConnected( $x, id$ )
8  return Clustering  $\{\mathcal{D}_i\}_{i=1}^{id}$ , where  $\mathcal{D}_i = \{x \in \mathcal{D} : x \text{ has label } i\}$ 

DensityConnected ( $x, id$ ):
9  Mark  $x$  with current cluster  $id$ 
10 foreach  $y \in N_\epsilon(x)$  do
11     Mark  $y$  with current cluster  $id$ 
12     if  $y$  is core then
13         DensityConnected( $y, id$ )

```

Above is the pseudo-code for the DBSCAN Algorithm given in the slides. I will explain what is happening in this algorithm, step by step.

First – D is the set of all datapoints, ϵ is the radius for finding neighbors, $minpts$ is the minimum number of close neighbors required to be a core.

Second – The first foreach loop iterates through all the points, calculates its neighbors by finding all other points within the radius ϵ and labels them core, border or noise.

Core is for points that have $minpts$ amount of points within ϵ .

Border is for points that have fewer than $minpts$ but fall within the range of a core.

Noise is for points that have fewer than $minpts$ and don't fall within range of a core.

Third – We then perform another foreach loop going through all points in D that are core and not labeled. Each one is given an id and sent to the `DensityConnected` function. This function's purpose is to label all points within core and border distance of the initial point with the same label. This process will label all core and border points that are close, giving us our density clusters.

Finally – We return the clusters.

Bonus:

C1 – $\{(2, 10), (2, 5), (8, 4)\}$: Centroid – (4, 6.33)

C2 – $\{(5, 8), (7, 5), (6, 4)\}$: Centroid – (6, 5.67)

C3 – $\{(1, 2), (4, 9)\}$: Centroid – (2.5, 5.5)