

Mining massive Datasets WS 2017/18

Problem Set 2

Rudolf Chrispens, Marvin, Daniela Schacherer

November 5, 2017

Exercise 01

See hand written solution.

Exercise 02

In order to also maintain information about which points are in which cluster one could store each Point p as

(*vector sum of points in cluster, number of points in cluster, list of data points in the cluster*).

Pseudocode for step 3 and 4 of subroutine *merge* from lecture 3:

- Find best pair, merge those two points/clusters and compute new cluster center
d, (p,q), (ip, iq) = bestPair
sum = p[0] + q[0]
count = p[1] + q[1]
points = p[2].append(q[2])
newCenter = (**sum**, **count**, **points**)
- Filter p and q from the inCluster
- Re-number centroid index in inCluster
- Add new cluster to the outCluster

Exercise 04 - Pseudocode

Step 1: Take a small sample of the data and cluster it in main memory.

In principle, any clustering method could be used, but as CURE is designed to handle oddly shaped clusters, it is often advisable to use a hierarchical method in which clusters are merged when they have a close pair of points.

```

#Take a Sample of the whole Dataset
#use hierarchical clustering algorithm to cluster the sampleData
    sampleSize;
    sampleData = textFile.takeSample(false, sampleSize);
    hierarchicalCluster = Hcluster( sampleData );

```

Step 2: Select a small set of points from each cluster to be representative points. These points should be chosen to be as far from one another as possible, using the method described in Section 7.3.2.

```

#Select small set of points of each cluster.
#They should be as far as possible from each other
FOR EACH hCluster in hierarchicalCluster DO
    #Pick the first point at random
        repPoints.add( hCluster.takeSample(false, 1) );
    WHILE repPoints.count() < hierarchicalCluster.count() DO
        #Add the point whose minimum distance from the selected points is as largest;
        furthestPoint;
    FOR EACH point in hCluster DO
        If( furthestPoint.squared_distance( repPoints ) < Point.squared_distance( repPoints ) ) DO
            furthestPoint = point;
        repPoints.add( furthestPoint);
END;

```

Step 3: Move each of the representative points a fixed fraction of the distance between its location and the centroid of its cluster. Perhaps 20% is a good fraction to choose. Note that this step requires a Euclidean space, since otherwise, there might not be any notion of a line between two points.

```

FOR EACH repPoint in repPoints DO
    Center = clusterCenters(repPoints);
    repPoint.moveTo(Center,0.2);

```

Step 4: The next phase of CURE is to merge two clusters if they have a pair of representative points, one from each cluster, that are sufficiently close. The user may pick the distance that defines ?close.? This merging step can repeat, until there are no more sufficiently close clusters.

```

Threshold;
FOR EACH cluster in repPoints DO
    IF (closestClusterFound = closest_cluster(repPoints, cluster, Threshold):
        mergeCluster(closestClusterFound, cluster):

```

Step 5: The last step of CURE is point assignment. Each point p is brought from secondary storage and compared with the representative points. We assign p to the cluster of the representative point that is closest to p .

```
FOR EACH point in Dataset DO
    nearestCluster(point, repPoints).add(point);
```