
COMP9318 (2009 s1) Tutorial 2

Wei WANG

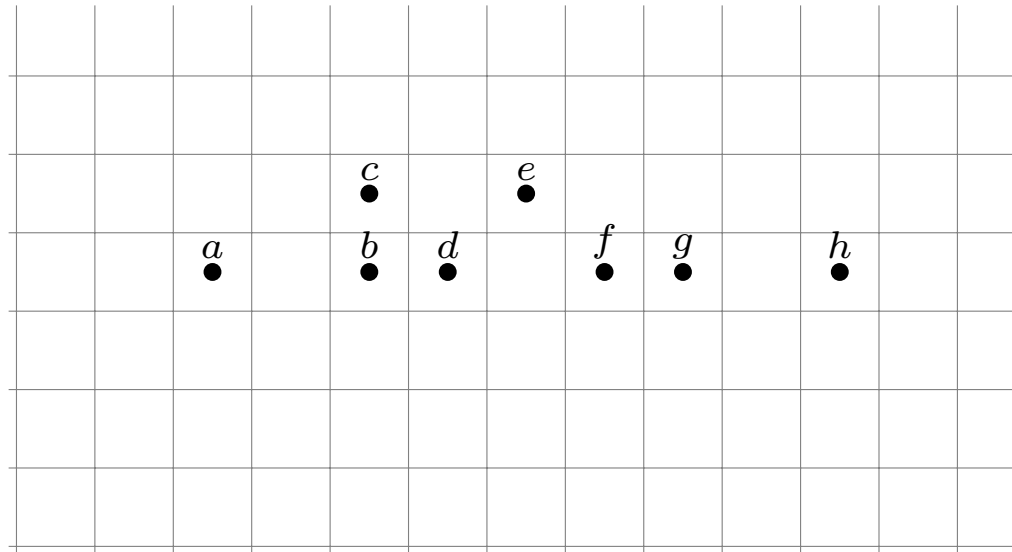
The University of New South Wales

weiw@cse.unsw.edu.au

① Clustering

Q1

Consider eight tuples represented as points in the two dimensional space as follows:



Assume that (1) each point lies within the center of the grid; (2) the grid is a uniform partition of the data space; and (3) each grid is a square with sides of length 1.

We consider applying DBSCAN clustering algorithm on this

dataset. Specifically, we set the minimum number of points (excluding the point in the center) within the ϵ -neighborhood ($MinPts$) to be 3, the radius of the ϵ -neighborhood (ϵ) to be 2, and we adopt the **Manhattan Distance** metric in the computation (i.e., a point p is within the ϵ -neighborhood of point o if and only if the Manhattan distance between p and o is no larger than ϵ).

- ① What is the Manhattan distance between point a and e ?
- ② List all the *core objects*.
- ③ What is the clustering result of the DBSCAN algorithm on this dataset if points are accessed following the alphabetical order? You need to write out all the cluster (with points that belonging to the cluster), as well as the outliers (if any).

SOLUTION TO Q1

- ① The Manhattan distance between a and e is 5.
- ② Consider each object and list number of points (excluding itself) in the neighborhood.

a	1
b	3
c	3
d	4
e	3
f	3
g	2
h	1

Therefore, the core objects are $\{b, c, d, e, f\}$

- ③ DBScan:
→ Start from a , a is not a core object, skip it.

→ Process b , b is a core object, then recursively grow the cluster of b .

The final cluster will be $\{b, a, c, d, e, f, g\}$.

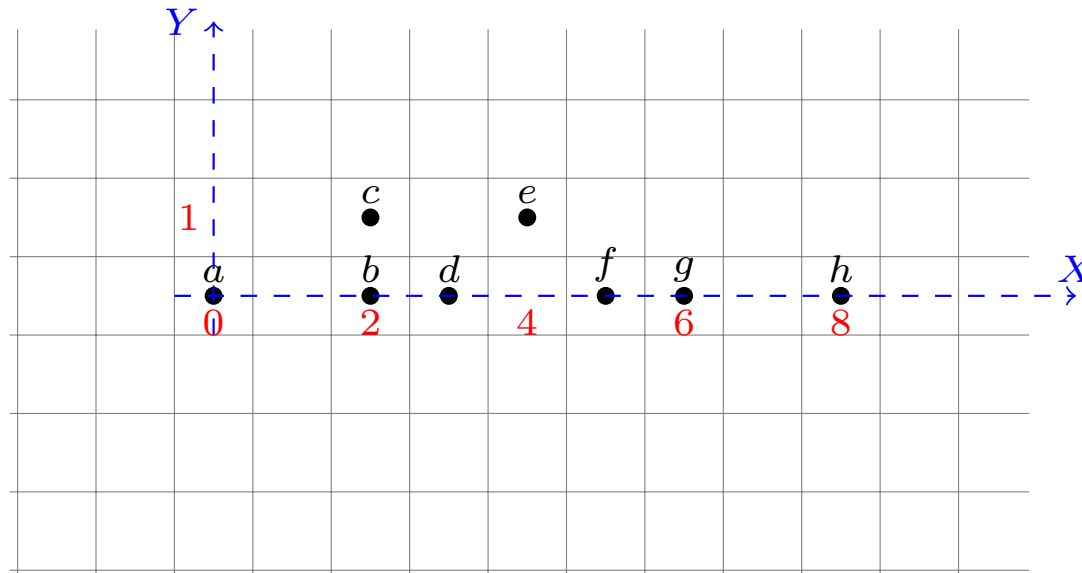
Therefore, the final answer is 1 cluster ($\{a, b, a, c, d, e, f, g\}$) and 1 outlier (h).

Q2

Consider the same dataset as Q1. What is the result of applying average link clustering algorithm on the dataset (still using the Manhattan distance)?

SOLUTION TO Q2

- ① Let's move the origin of the coordinate to a (it is easy to see that the clustering result does not depend on the origin of the coordinate system) and every point is now assign an coordinate.



a	$(0, 0)$
b	$(2, 0)$
c	$(2, 1)$
d	$(3, 0)$
e	$(4, 1)$
f	$(5, 0)$
g	$(6, 0)$
h	$(8, 0)$

② Average Link:

- Initially, every point is assigned into a distinct cluster.
- The closest pair of points (under L_1 distance) is one of (b, c) , (b, d) , and (f, g) . Assume we take (b, c) . (You may break the tie in any consistent way). We only need to update the distance between (b, c) and d , which is $\frac{1+2}{2} = 1.5$.
- The next pair to merge is (f, g) . After the merge, (f, g) will have the same distance of $\frac{3+2}{2} = 2.5$ to d , e , or h .

-
- The next pair to merge is (b, c) and d . Afterwards, we calculate the new distance between clusters as

$(a) - (b, c, d)$	2.67
$(b, c, d) - e$	2.33
$(b, c, d) - (f, g)$	3.50

- The next pair to merge is (b, c, d) and e . Afterwards, we calculate the new distance between clusters as

$(a) - (b, c, d, e)$	3.25
$(b, c, d, e) - (f, g)$	3.25

- The next pair to merge is (f, g) and h . Afterwards, we calculate the new distance between clusters as

$(b, c, d, e) - (f, g, h)$	4.08
----------------------------	------

- The next pair to merge is a and (b, c, d, e) .

Q3

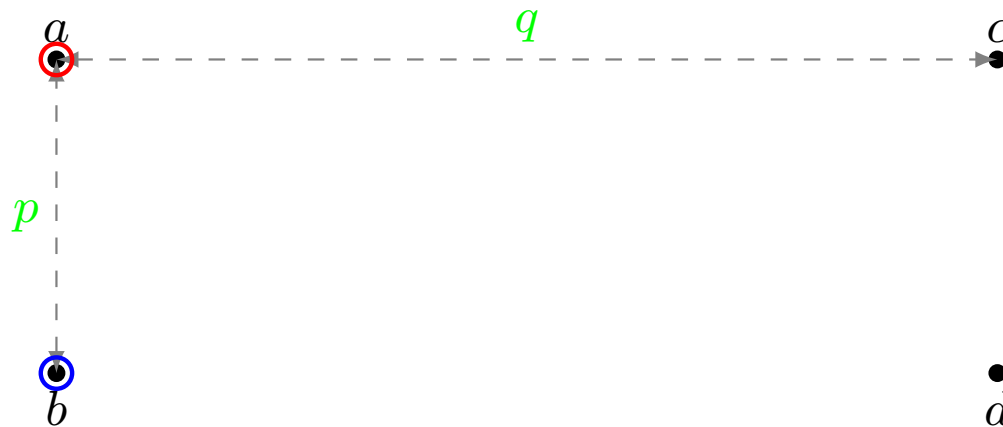
Consider the k-means clustering algorithm.

- ① Construct a simple example (with at most four data points) which shows that the clustering result of k-means algorithm can be arbitrarily worse than the optimal clustering results in terms of the cost. (The cost of a clustering is measured as the sum of square distance to the cluster center)
- ② Why the k-means algorithm updates the new cluster center for cluster C_j as the $(\frac{1}{|C_j|} \sum_{i=1}^{|C_j|} x_i, \frac{1}{|C_j|} \sum_{i=1}^{|C_j|} y_i)$?

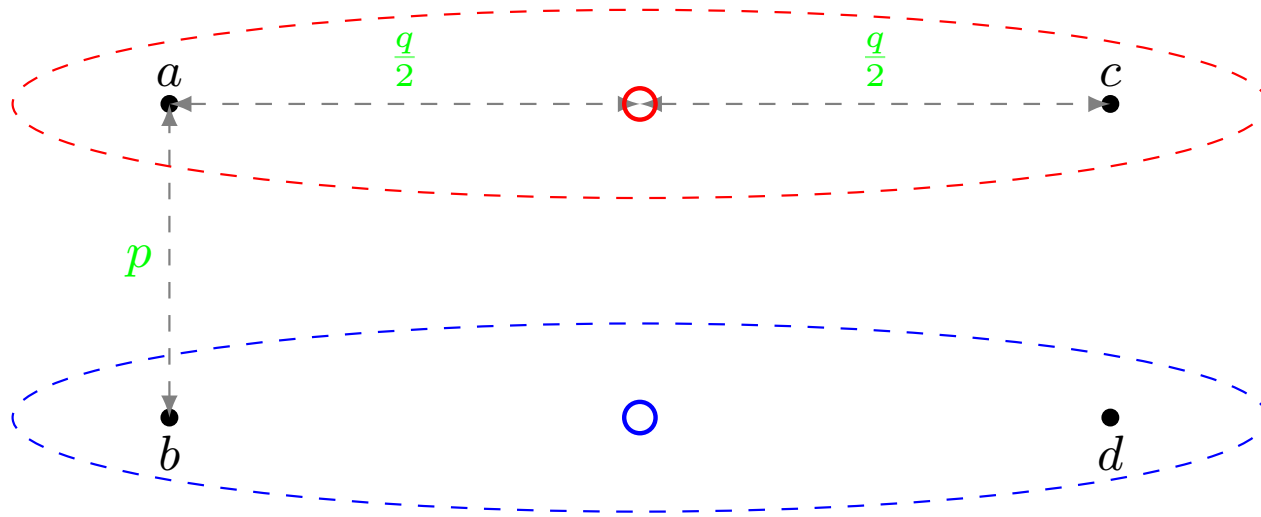
SOLUTION TO Q3

Consider the k-means clustering algorithm.

- ① See the figure below for $k = 2$ and $p < q$ and let the initial cluster centers be a and b .



The final clusters will be:



The cost of this clustering result is

$$((q/2)^2 + (q/2)^2) + ((q/2)^2 + (q/2)^2) = q^2$$

The optimal clustering is to group a and b in one cluster and c and d in another cluster. The total cost is

$$((p/2)^2 + (p/2)^2) + ((p/2)^2 + (p/2)^2) = p^2$$

The cost ratio is $(\frac{q}{p})^2$ and can be arbitrarily large.

② Let the new cluster center be (o_x, o_y) . The total cost for the current

cluster C_j is

$$\begin{aligned} Cost_j &= \sum_{1 \leq i \leq |C_j|} ((x_i - o_x)^2 + (y_i - o_y)^2) \\ &= \sum_{1 \leq i \leq |C_j|} (x_i - o_x)^2 + \sum_{1 \leq i \leq |C_j|} (y_i - o_y)^2 \\ &= XCost_j + YCost_j \end{aligned}$$

It is obvious that $XCost_j$ and $YCost_j$ is independent of each other if we want to minimize the overall cost.

To minimize $XCost_j$, we can set o_x such that $\frac{dXCost_j}{do_x} = 0$

$$\begin{aligned} \frac{dXCost_j}{do_x} &= \frac{d}{do_x} \sum_{1 \leq i \leq |C_j|} (x_i^2) - 2o_x \sum_{1 \leq i \leq |C_j|} x_i + \sum_{1 \leq i \leq |C_j|} o_x^2 \\ &= 2|C_j|o_x - 2 \sum_{1 \leq i \leq |C_j|} x_i \end{aligned}$$

Obviously, when $o_x = \frac{1}{|C_j|} \sum x_i$, the above equation equals 0, and it can be verified that the cost is minimized.