

---

# COMP9318 (2009 s1) Tutorial 2

Wei WANG

The University of New South Wales

[weiw@cse.unsw.edu.au](mailto:weiw@cse.unsw.edu.au)

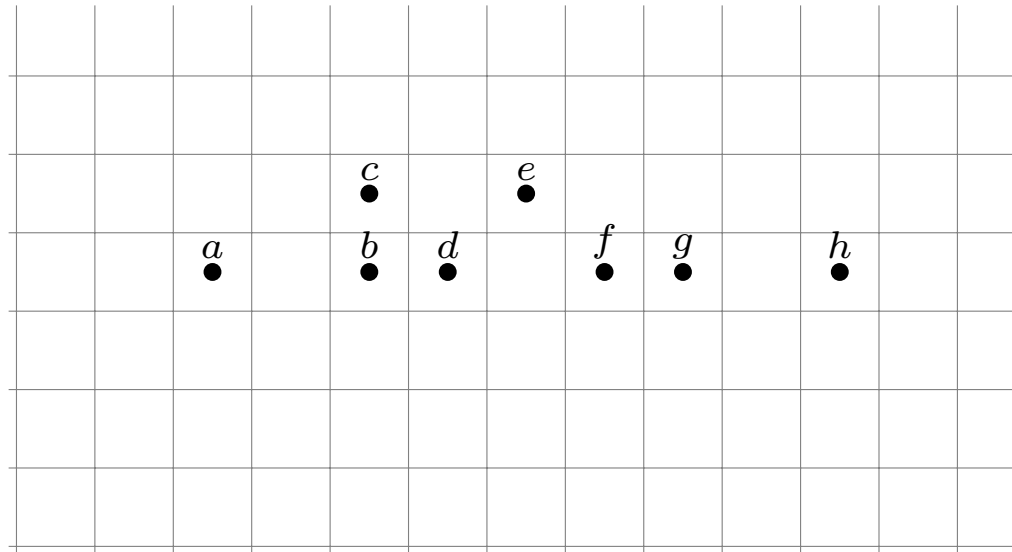
---

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  $\epsilon$ -neighborhood ( $MinPts$ ) to be 3, the radius of the  $\epsilon$ -neighborhood ( $\epsilon$ ) to be 2, and we adopt the **Manhattan Distance** metric in the computation (i.e., a point  $p$  is within the  $\epsilon$ -neighborhood of point  $o$  if and only if the Manhattan distance between  $p$  and  $o$  is no larger than  $\epsilon$ ).

- 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).

---

## 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)?

---

## 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  $(\frac{1}{|C_j|} \sum_{i=1}^{|C_j|} x_i, \frac{1}{|C_j|} \sum_{i=1}^{|C_j|} y_i)$ ?