



COMP20008 Elements of Data Processing

Clustering and clustering visualisation



- Consultation session about the assignments
 - Today (March 20th), 2 pm
 - Room 7.02, Level 7, Doug McDonell (Building 169)
- Answers to workshop 2 will be available today.
- A mistake in the previous lecture!



- Clustering algorithms
 - K-means
 - Visualisation of clustering tendency (VAT)
- Next class
 - Hierarchical clustering (next class)



- For datasets with more than 4 dimensions
 - Difficult to visualise
- How can we determine what the significant groups/segments/communities are?
 - If we have this information
 - Can understand the data better
 - Apply separate interventions to each group (e.g. marketing campaign)

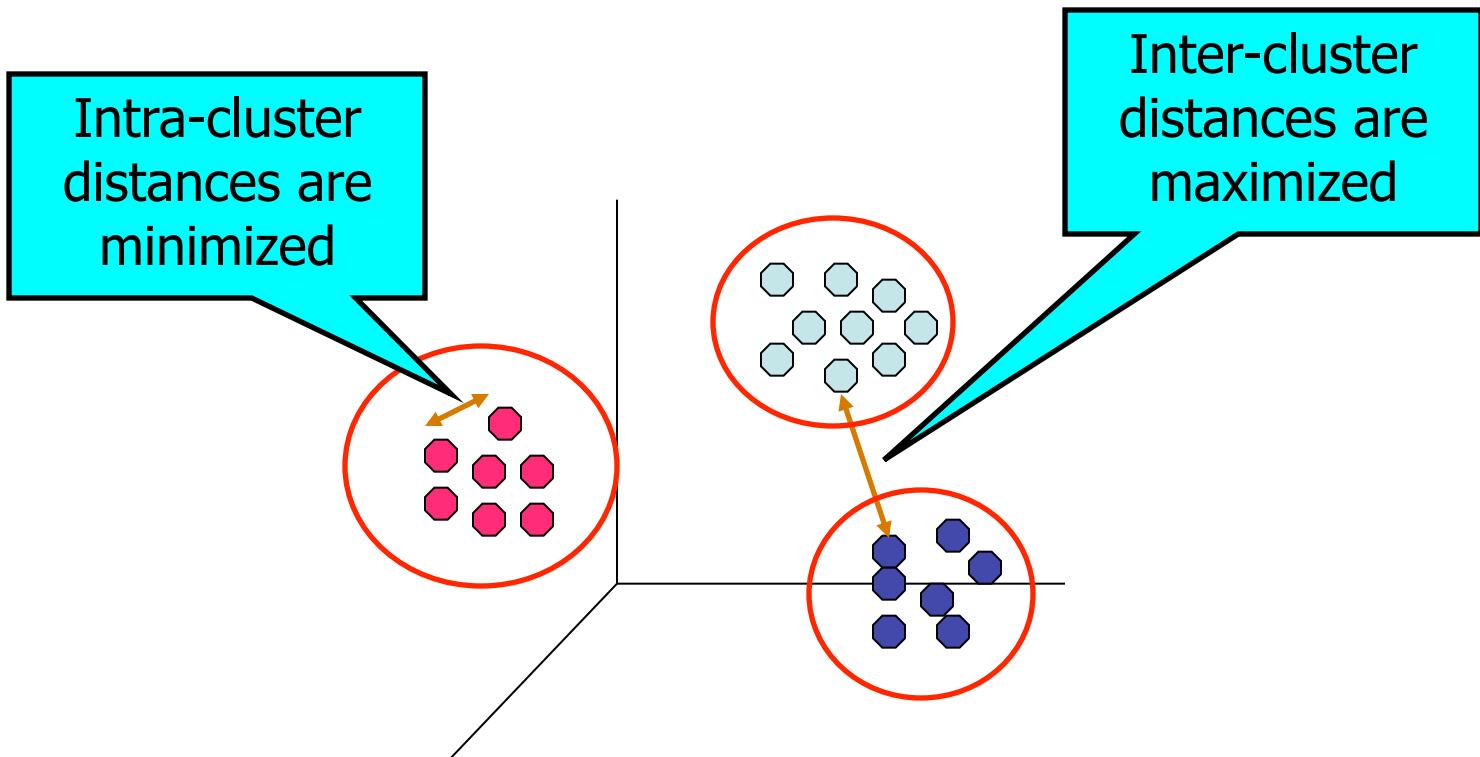


- A good clustering method will produce high quality clusters
 - Objects within same cluster are close together
 - Objects in different clusters are far apart
- Clustering is a major task in data analysis and visualisation, useful not just for outlier detection.
 - Market segmentation
 - Image analysis
 - Search engine result presentation
 -



What is Cluster Analysis?

- Figure below from Tan, Steinbach and Kumar 2004
- We will be looking at two classic clustering algorithms
 - K-means
 - Hierarchical clustering





- Clustering methods are typically distance based. Represent each object-instance as a vector and then can compute Euclidean distance between pairs of vectors.
- Commonly normalise each attribute into range [0,1] via a pre-processing step before computing distances

Given $X = (x_1, x_2, x_3, \dots, x_n)$ and $Y = (y_1, y_2, y_3, \dots, y_n)$

$$D(X, Y) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + (x_3 - y_3)^2 + \dots + (x_n - y_n)^2}$$



How to obtain the clusters?

- Need to assign each object to exactly one cluster
- Each cluster can be summarised by its centroid (the average of all its objects)



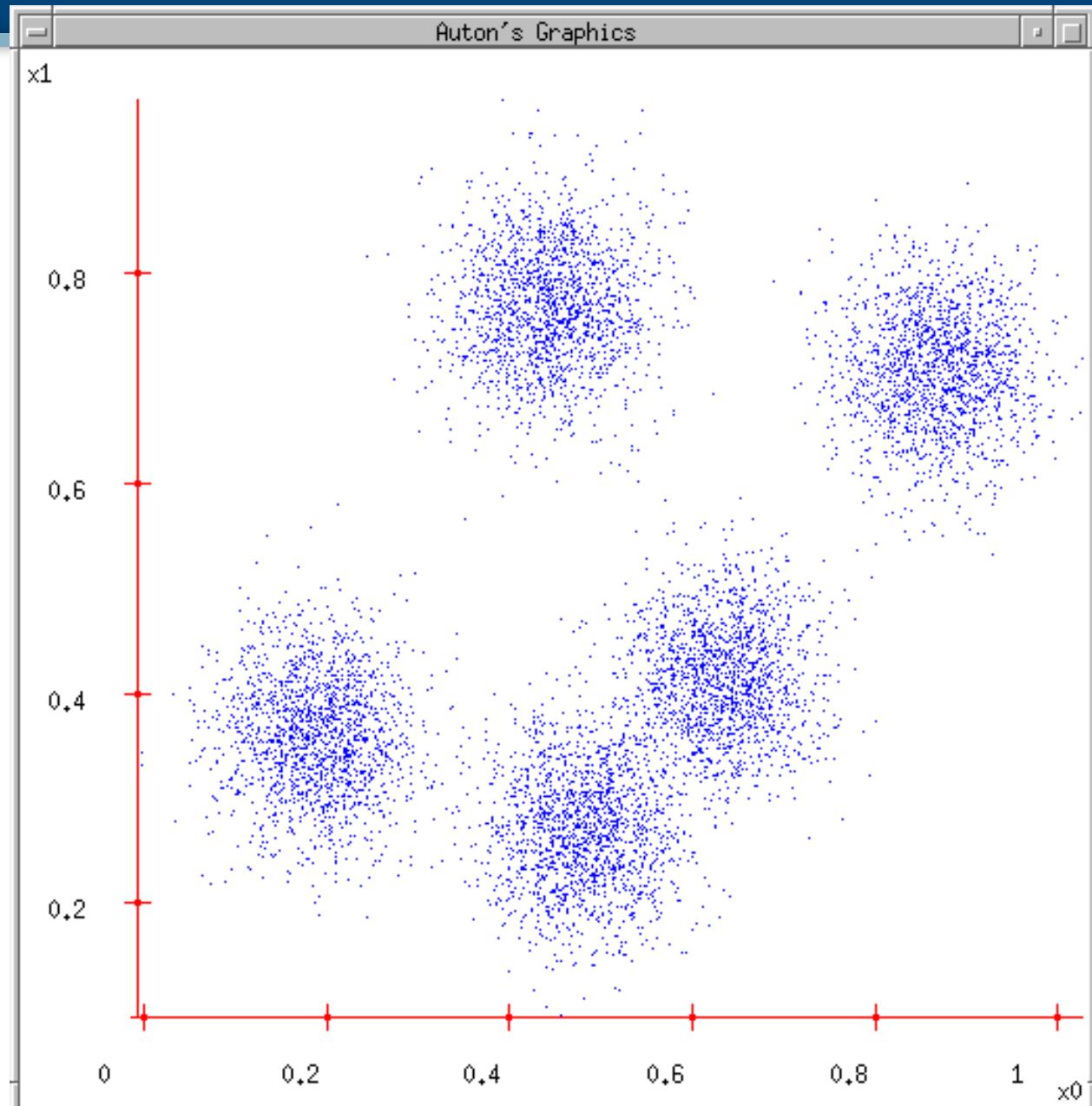
- Given parameter k , the *k-means* algorithm is implemented in four steps:
 1. Select k seed points as the initial cluster centres
 2. Assign each object to the cluster with the nearest seed point
 3. Compute new seed points as the centroids of the clusters of the current partitioning (the centroid is the center, i.e., ***mean point***, of the cluster)
 4. Go back to Step 2, stop when the assignment does not change



K-means

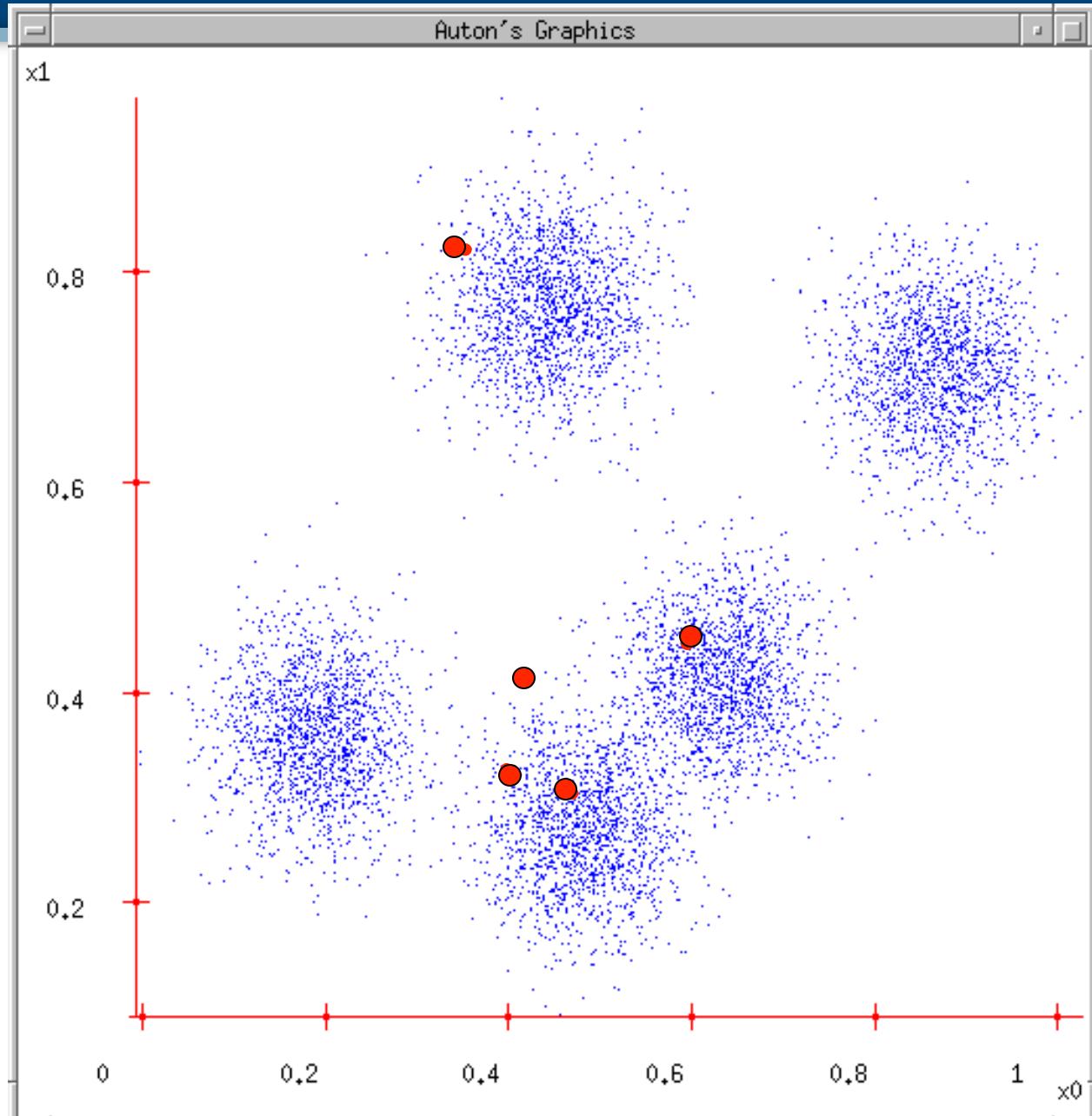
1. Ask user how many clusters they'd like.
(e.g. K=5)

(Example from Andrew Moore
[http://www.autonlab.org/tutorials/
kmeans11.pdf](http://www.autonlab.org/tutorials/kmeans11.pdf))



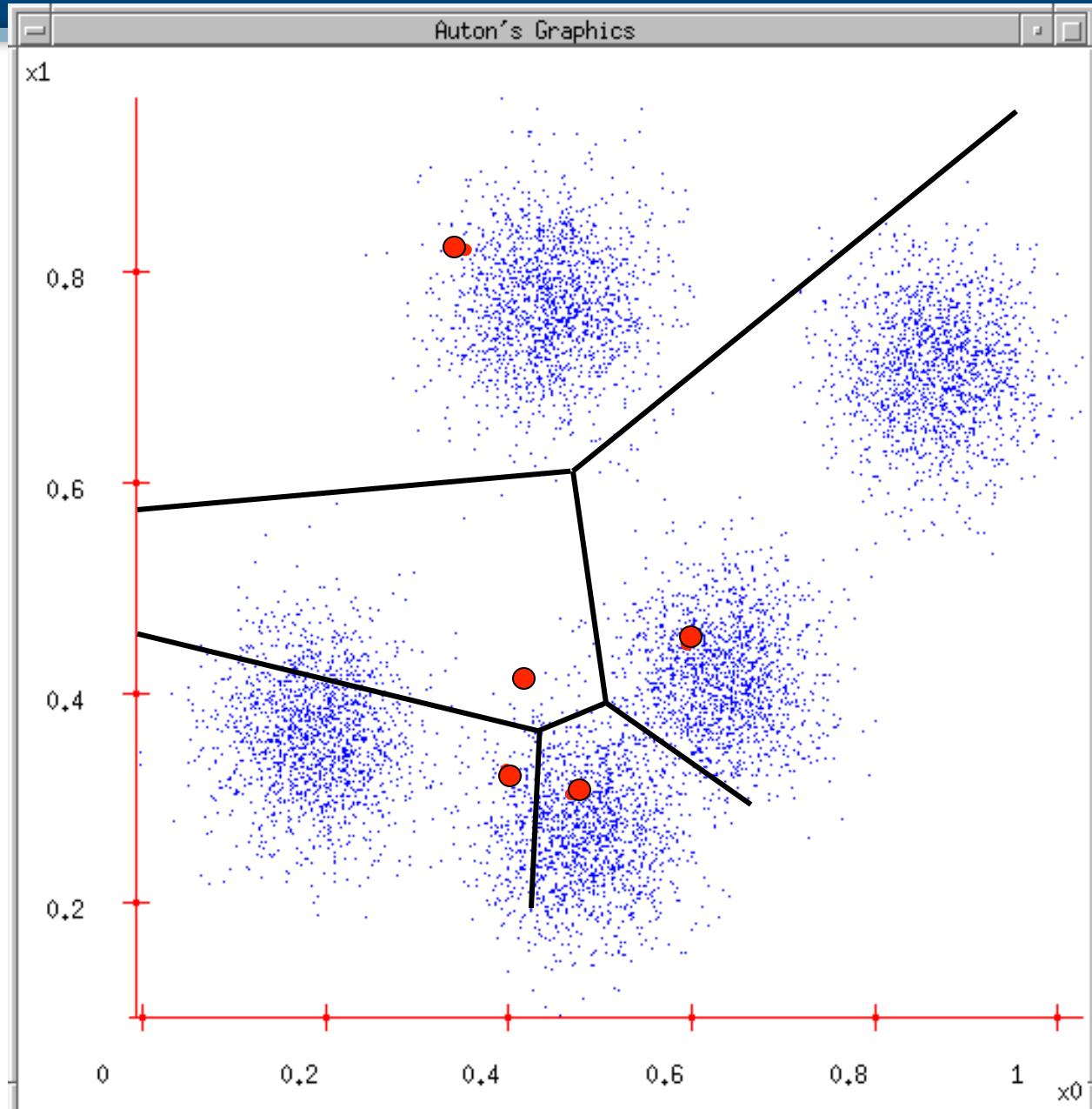


1. Ask user how many clusters they'd like.
(e.g. K=5)
2. Randomly guess K cluster Center locations



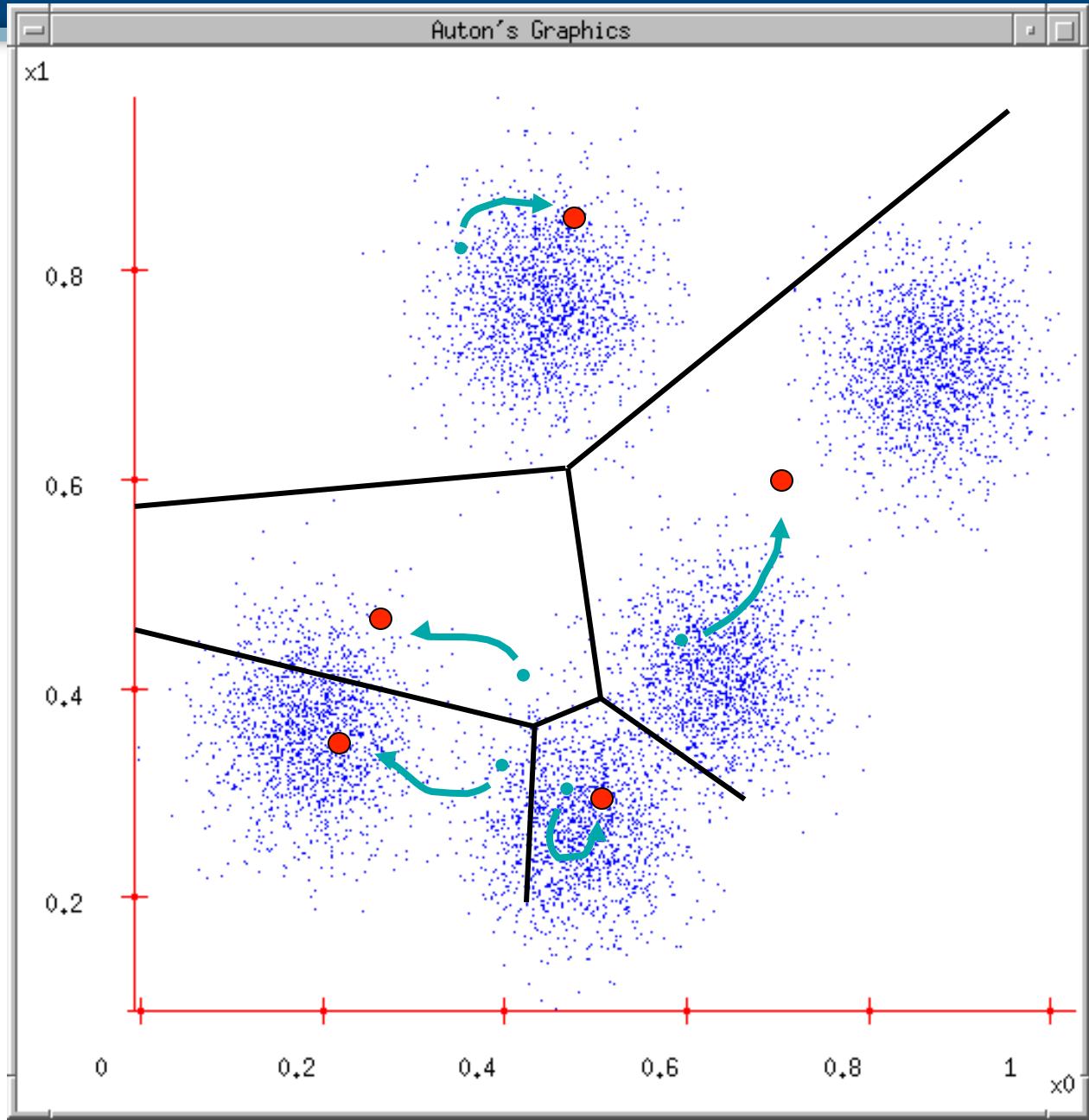


1. Ask user how many clusters they'd like.
(e.g. K=5)
2. Randomly guess K cluster Center locations
3. Each datapoint finds out which Center it's closest to. (Thus each Center "owns" a set of datapoints)



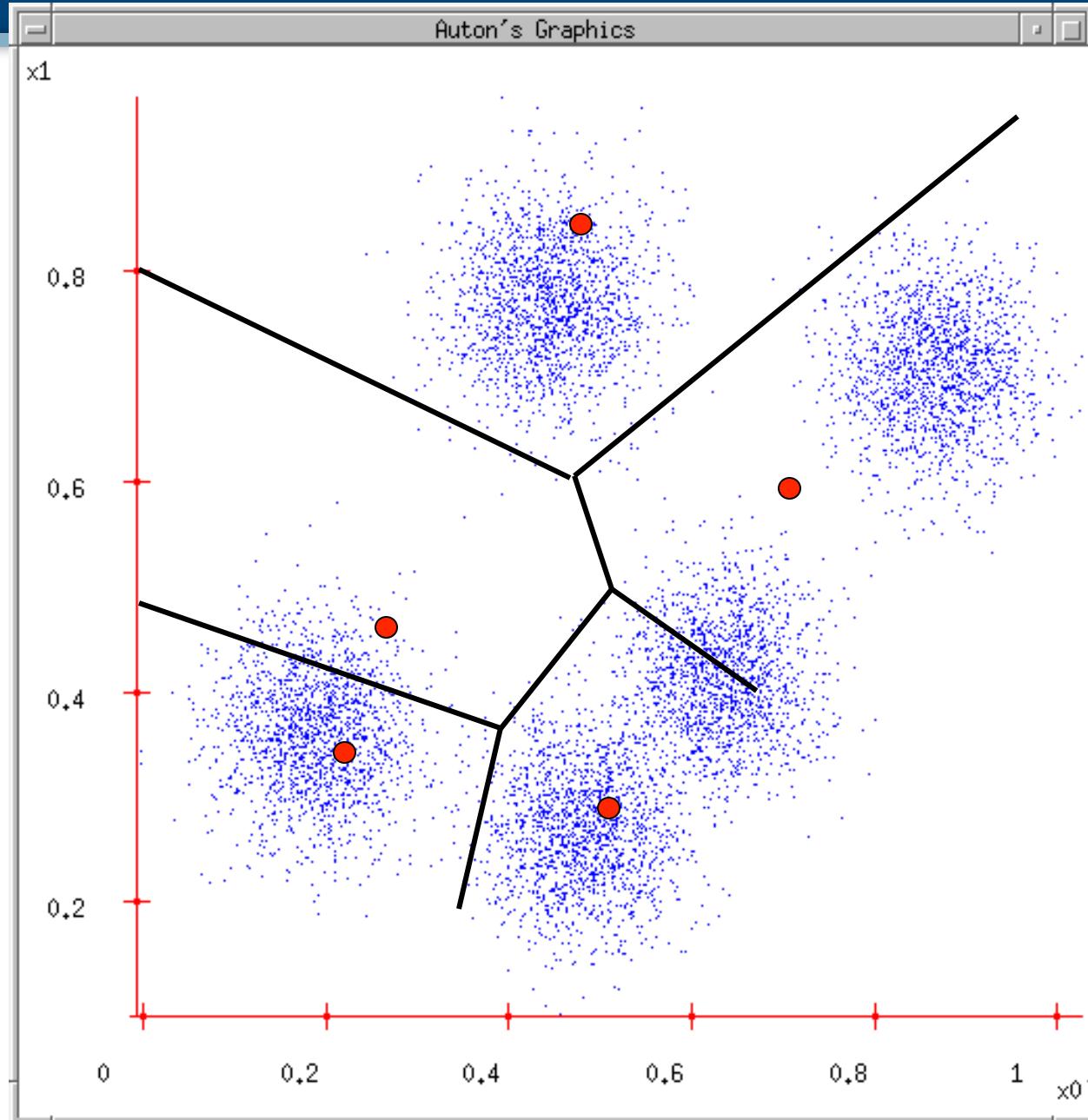


1. Ask user how many clusters they'd like.
(e.g. $k=5$)
2. Randomly guess k cluster Center locations
3. Each datapoint finds out which Center it's closest to.
4. Each Center finds the centroid of the points it owns





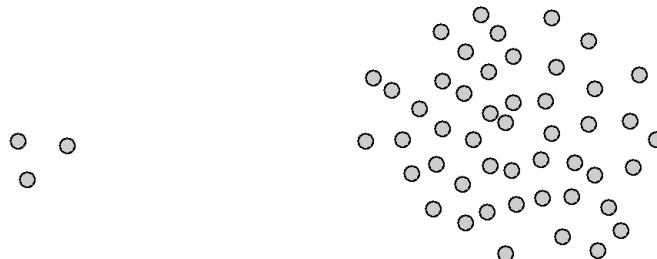
1. Ask user how many clusters they'd like.
(e.g. $k=5$)
2. Randomly guess k cluster Center locations
3. Each datapoint finds out which Center it's closest to.
4. Each Center finds the centroid of the points it owns
5. New Centers => new boundaries
6. Repeat until no change



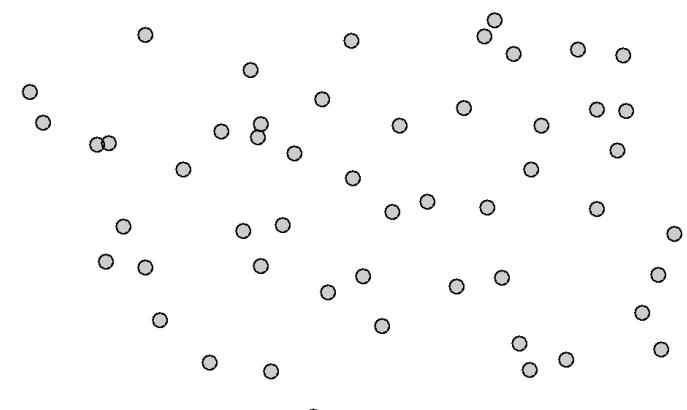


For which dataset does k-means require less number of iterations?
 $k=2$

Dataset 1



Dataset 2





- Typically choose the initial seed points randomly
 - Different runs of the algorithms will produce different results
- Closeness measured by Euclidean distance (Can also use other distance functions)
- Algorithm can be shown to converge (to a local optimum), typically doesn't require many iterations



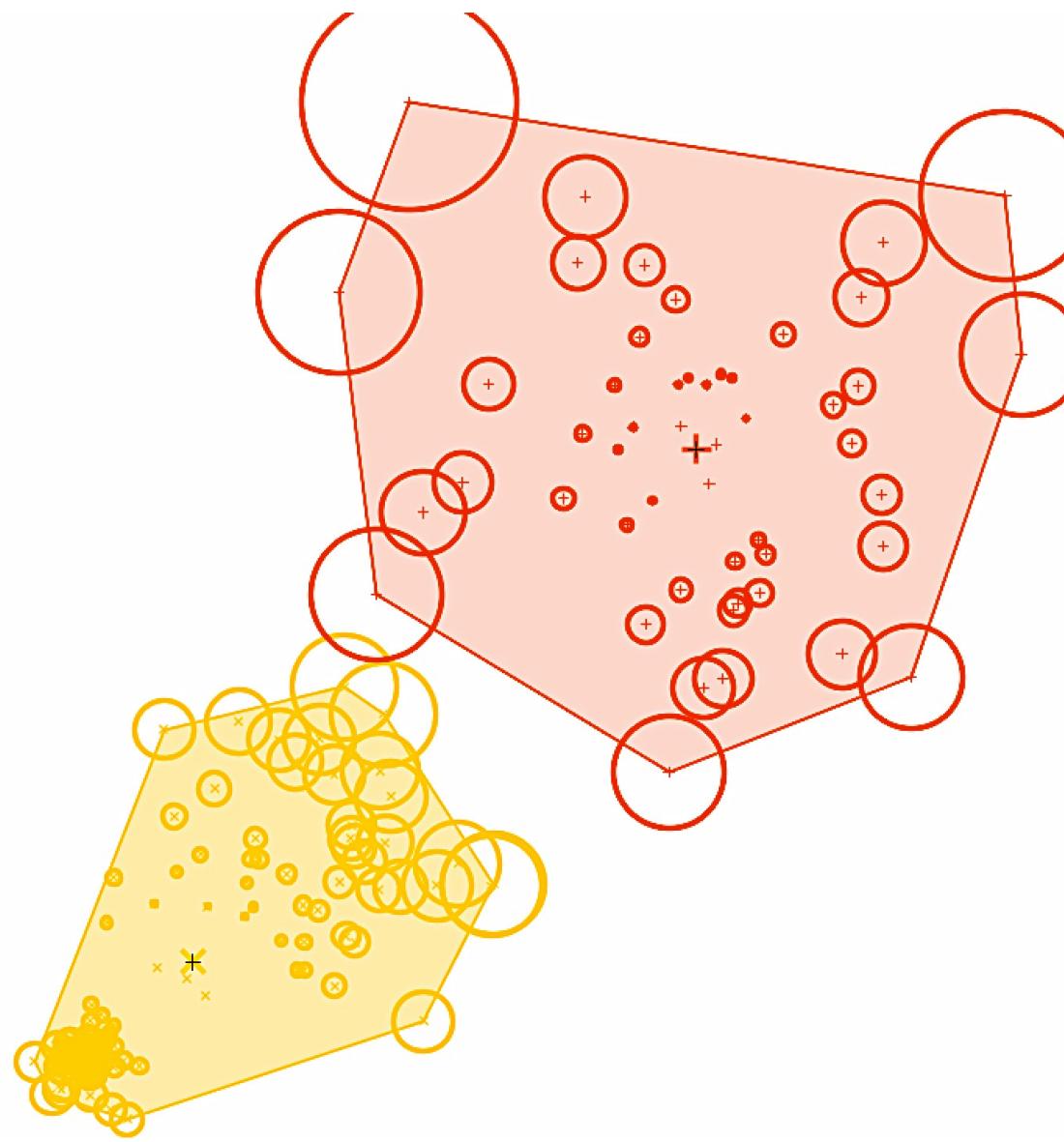
- [http://home.deib.polimi.it/matteucc/Clustering/tutorial_html/
AppletKM.html](http://home.deib.polimi.it/matteucc/Clustering/tutorial_html/AppletKM.html)



- An outlier is expected to be far away from any groups of normal objects
- Each object is associated with exactly one cluster and its outlier score is equal to the distance from its cluster centre.



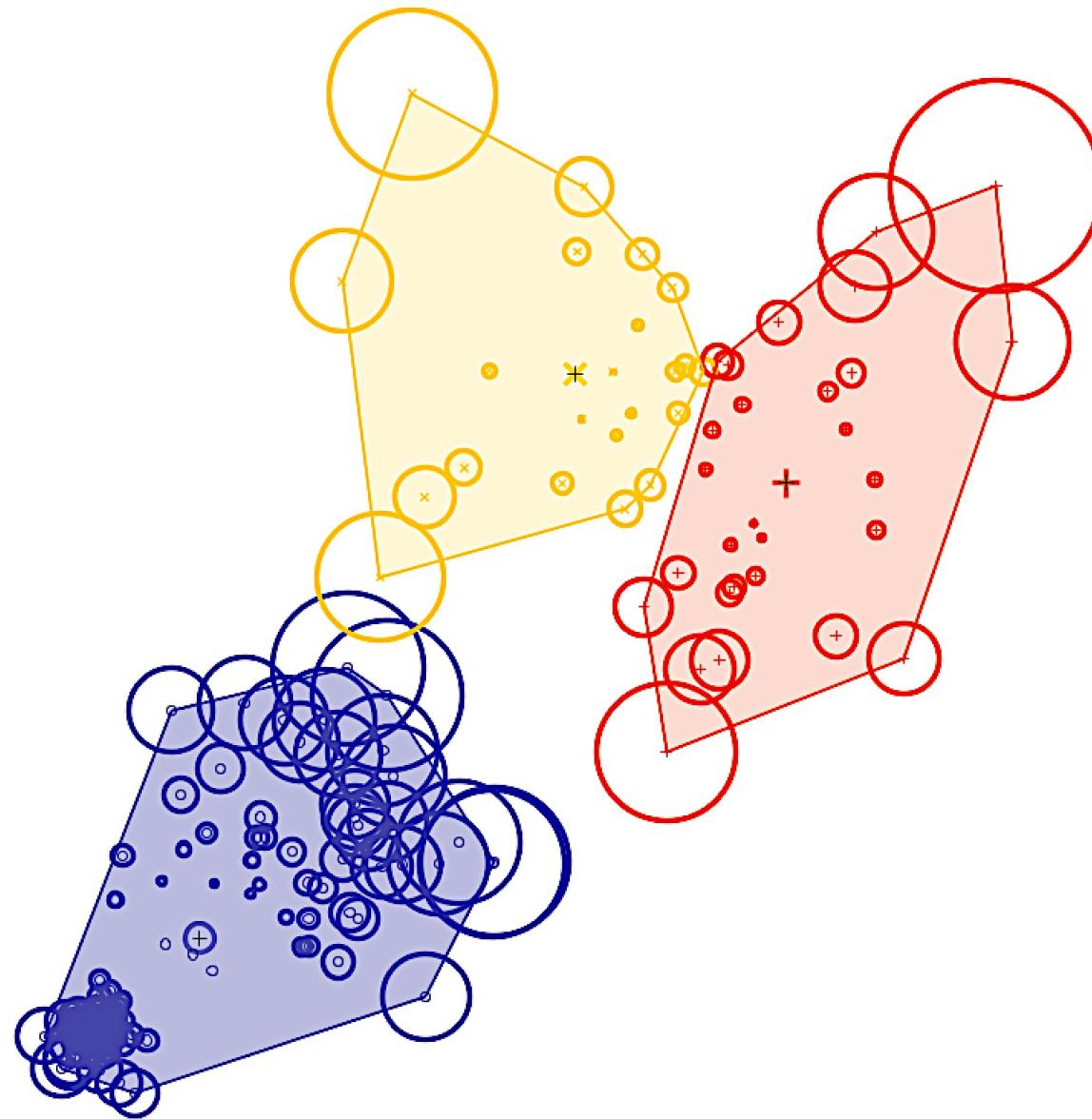
Clustering based outlier detection. 2 clusters





3 clusters

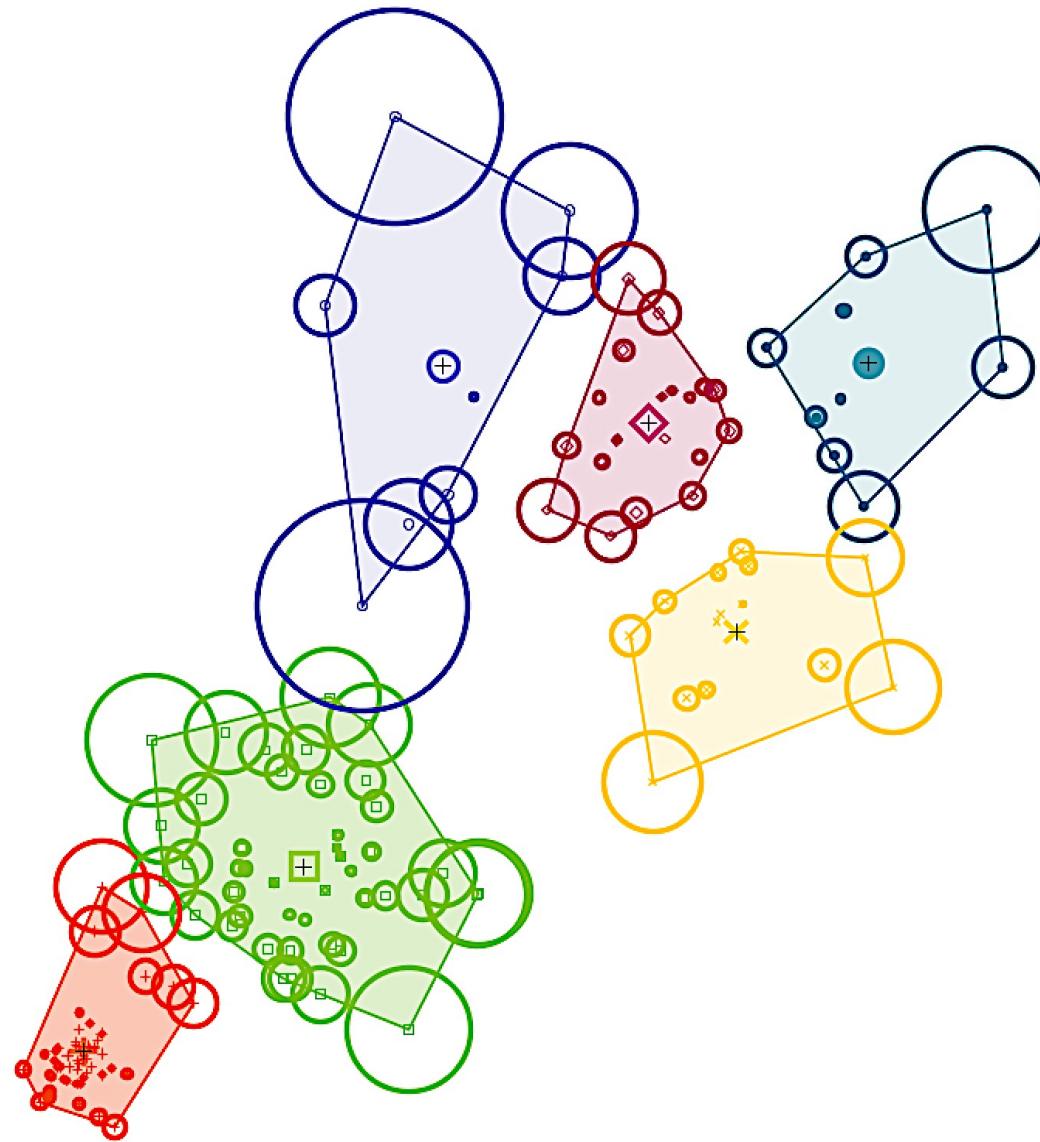
MELBOURNE





6 clusters

MELBOURNE





9 clusters



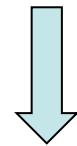


- How many clusters are in the data? How big are they? What is the likely membership of objects in each cluster?
 - The k parameter for k-means
- One solution: *visually determine the clustering structure by inspecting a heat map*
 - Represent datasets in an $n*n$ image format
 - Applicable for many different types of object data



Background: What is a dissimilarity matrix?

Object id	Feature1	Feature2	Feature3
1	5	10	15
2	10	5	10
3	20	20	20



Compute all pairwise distances between objects. This gives a dissimilarity matrix.

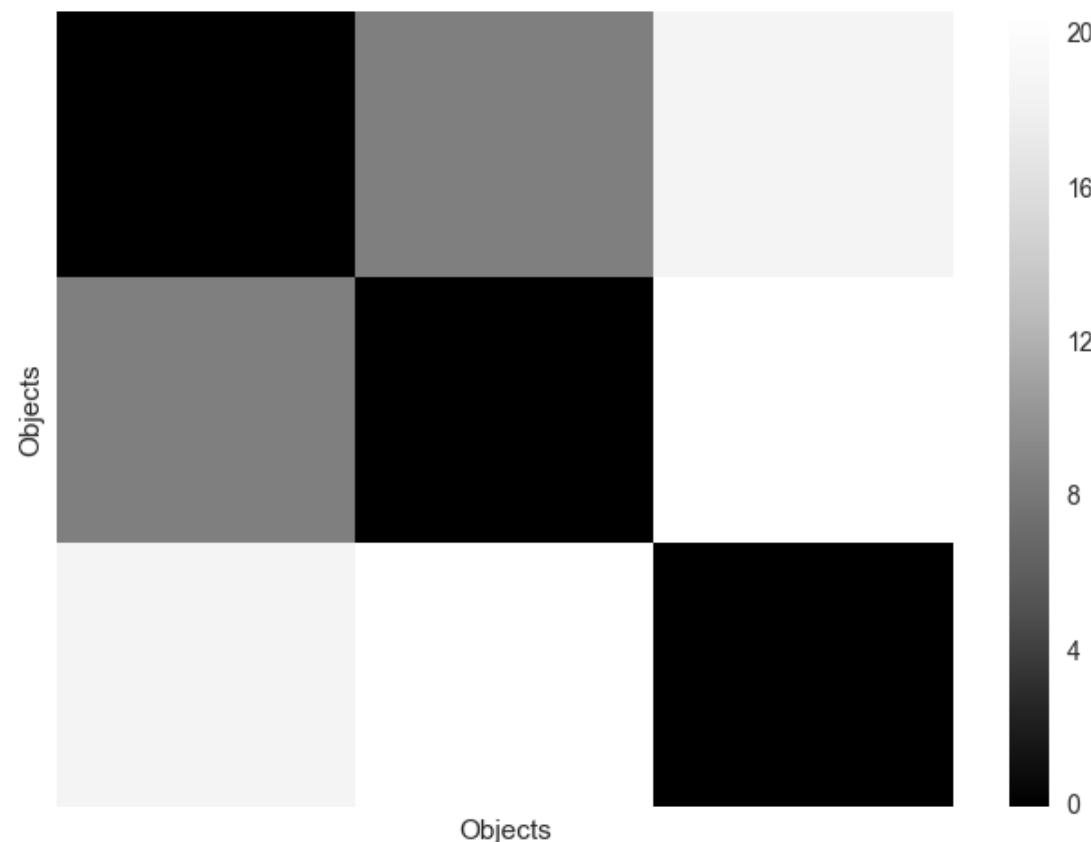
Object	1	2	3
1	0	8.7	18.7
2	8.7	0	20.6
3	18.7	20.6	0



Visualising a dissimilarity matrix

- We can visualise a dissimilarity matrix as a *heat map*, where the colour of each cell indicates that cell's value

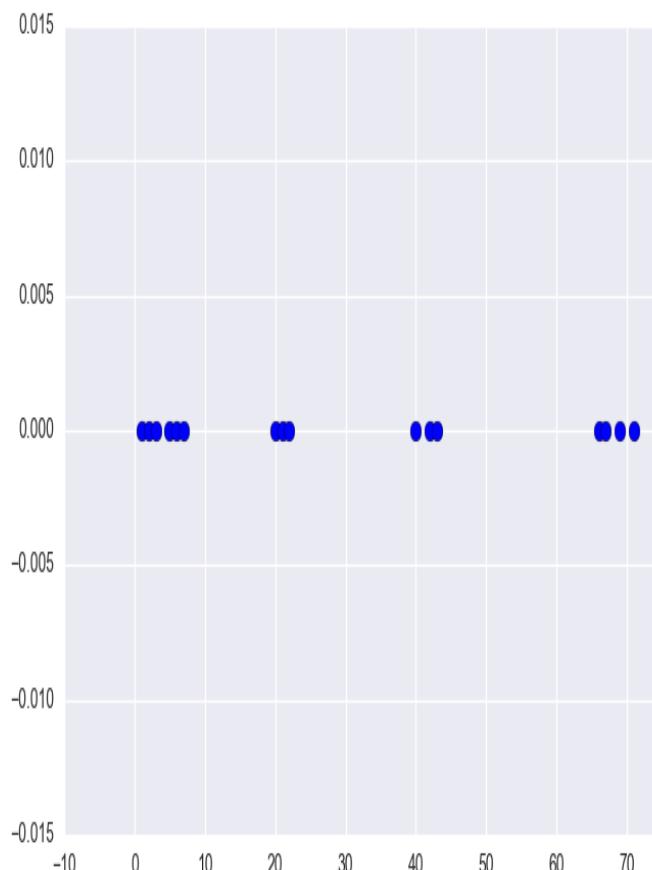
Object	1	2	3
1	0	8.7	18.7
2	8.7	0	20.6
3	18.7	20.6	0



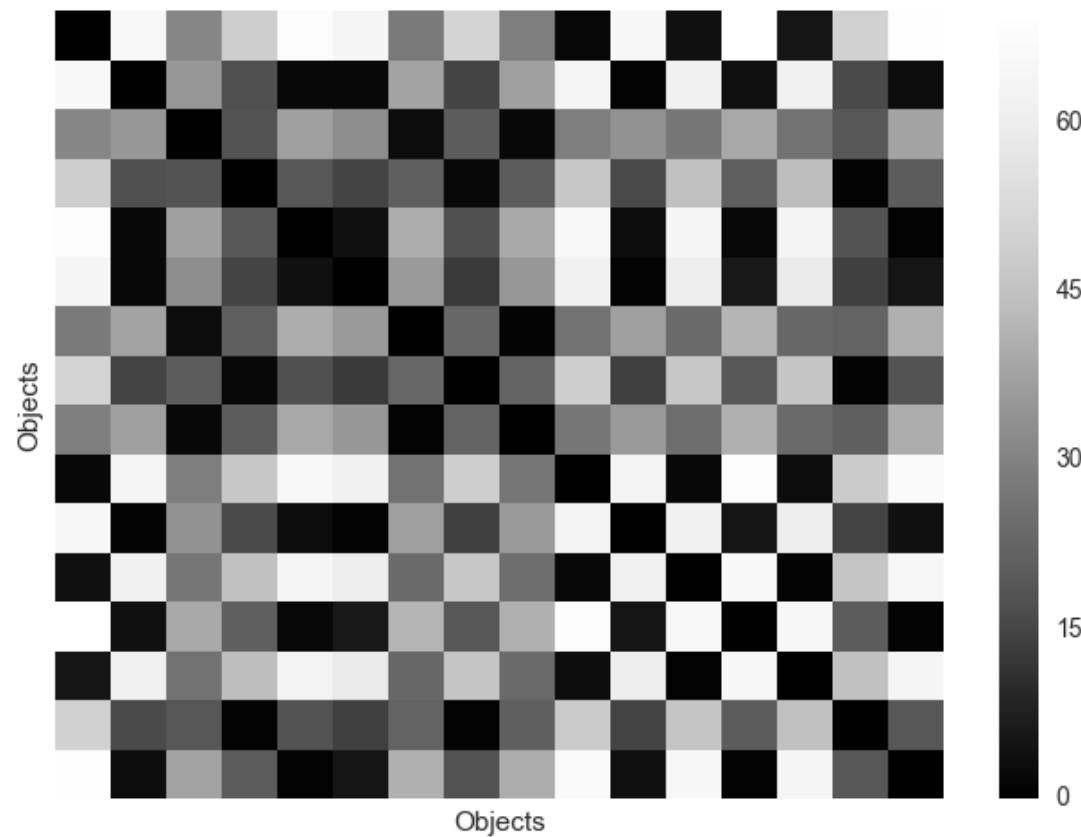
- The diagonal of D is all zeros
- D is symmetric about its leading diagonal
 - $D(i,j)=D(j,i)$ for all i and j
 - Objects follow the same order along rows and columns
- In general, visualising the (raw) dissimilarity matrix may not reveal enough useful information
 - Further processing is needed



Reordering a Dissimilarity matrix cont.

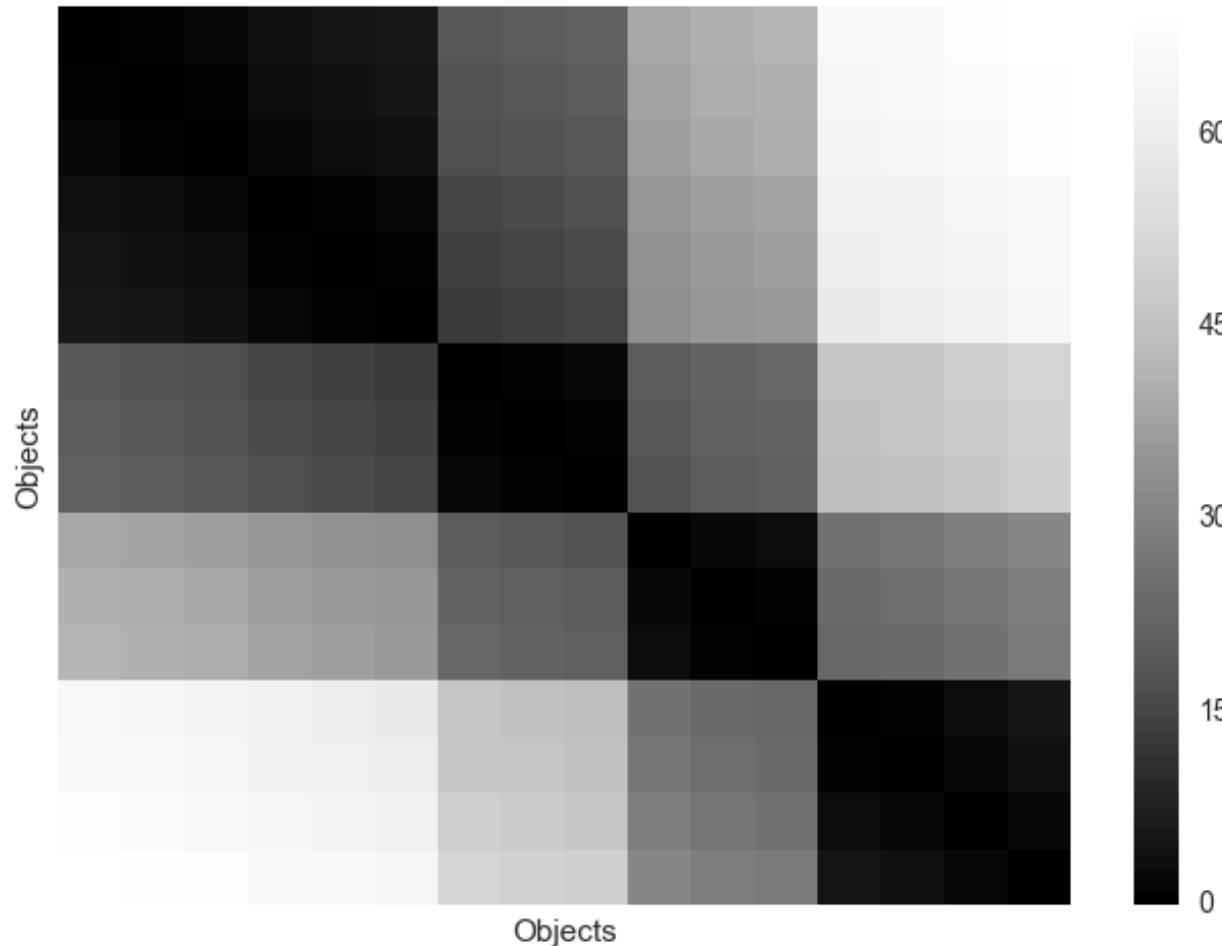


Example dataset
with 16 objects



Random order of the 16 objects
for the dissimilarity matrix

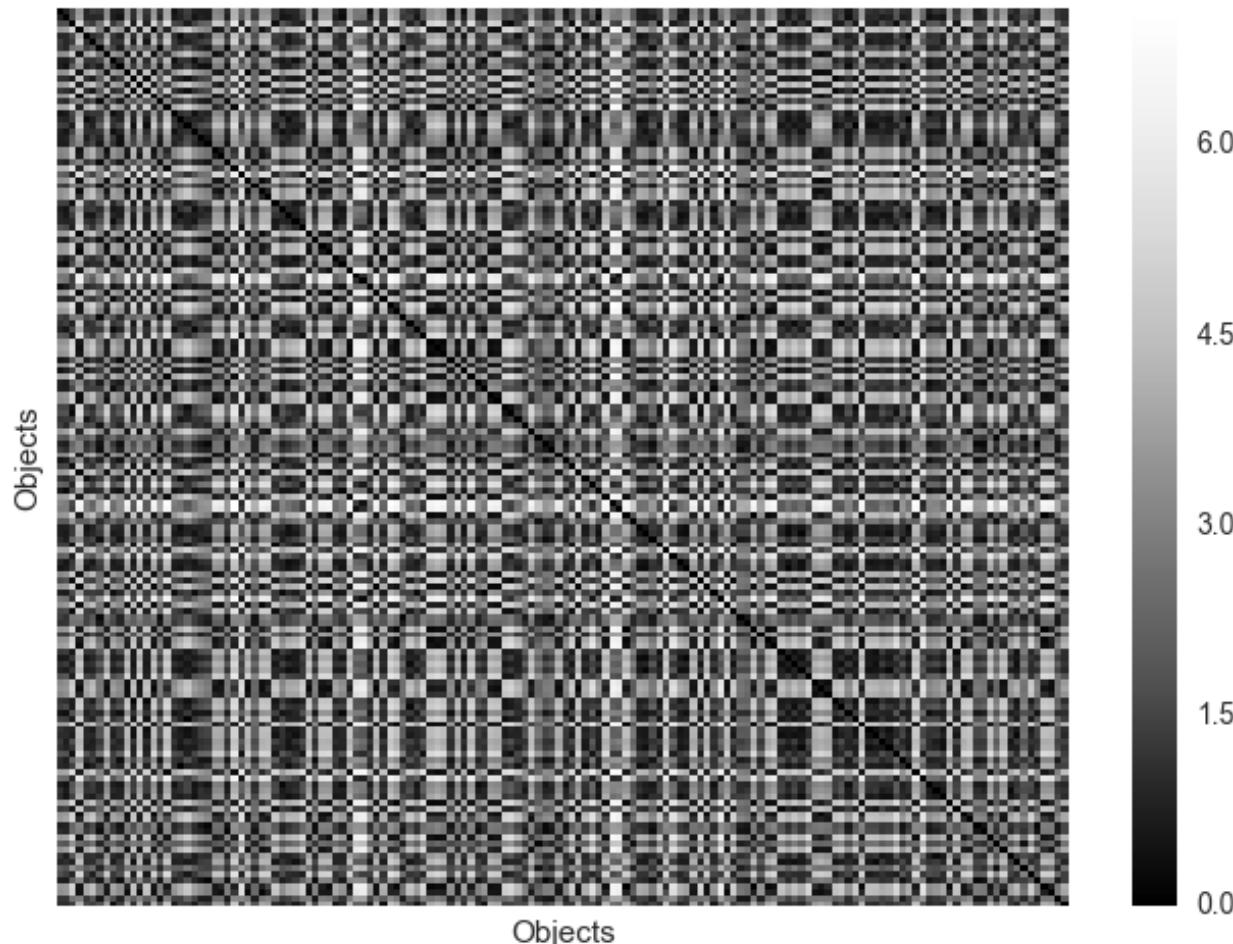
Reordering the matrix reveals the clusters



A better ordering of the 16 objects. Nearby objects in the ordering are similar to each other, producing large dark blocks. We can see four clusters along the diagonal.



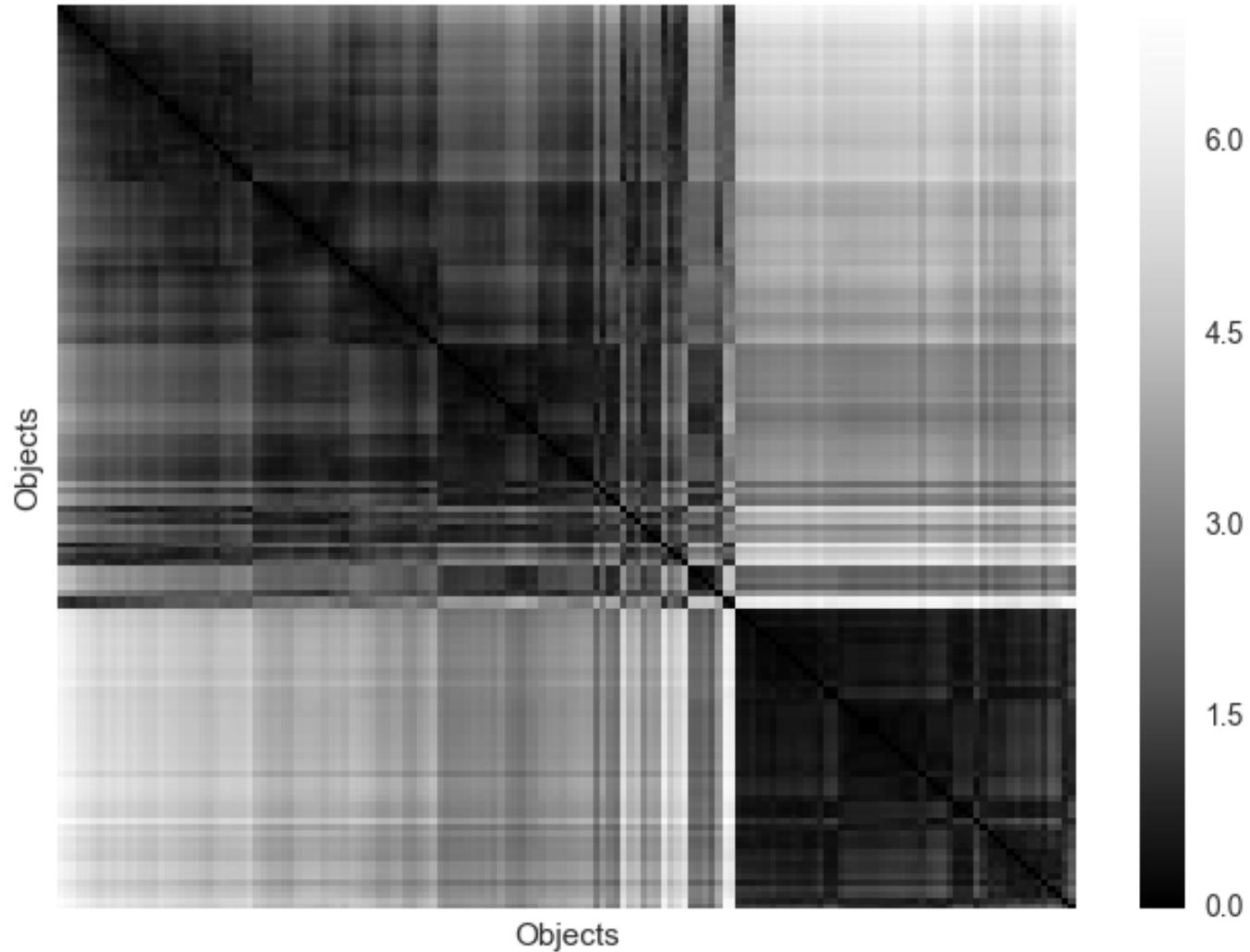
- A good VAT image suggests both the number of and approximate members of object clusters.
- A diagonal dark block appears in the VAT image only when a tight group exists in the data (low within-cluster dissimilarities)



Random order for 150 objects:
Where are the clusters???

Dissimilarity matrix for Iris data (reordered)

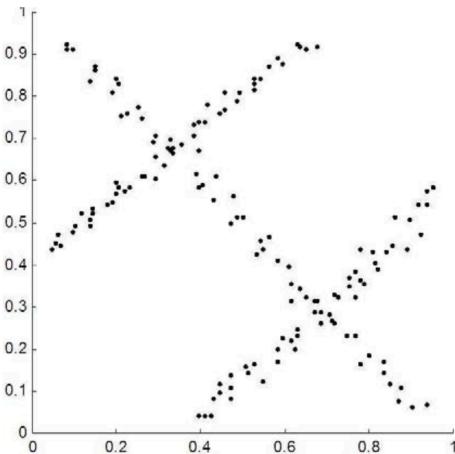
MELBOURNE



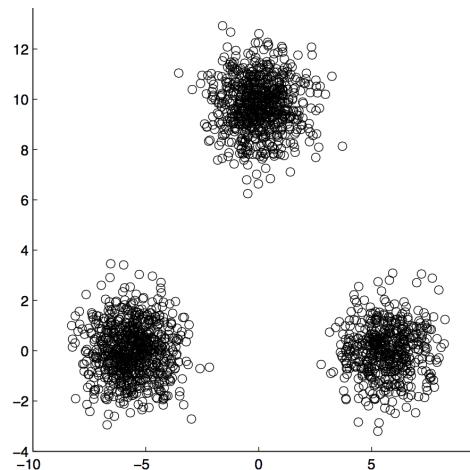


- Match the datasets with the VAT images

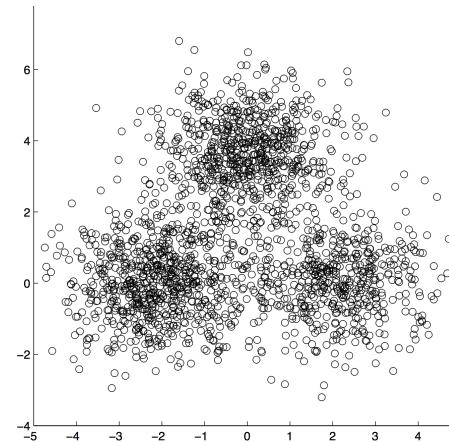
1



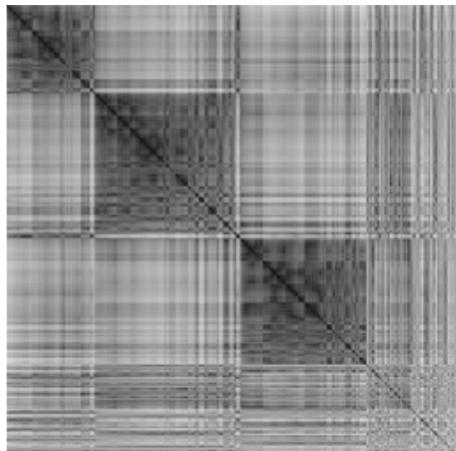
2



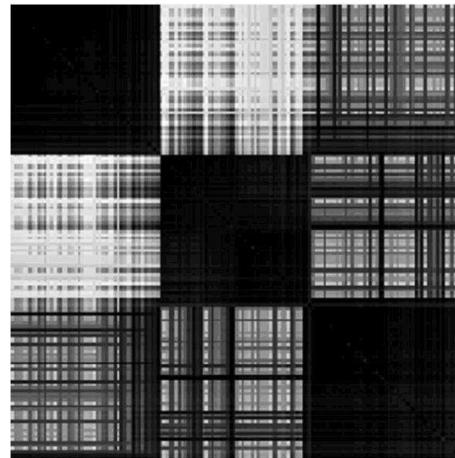
3



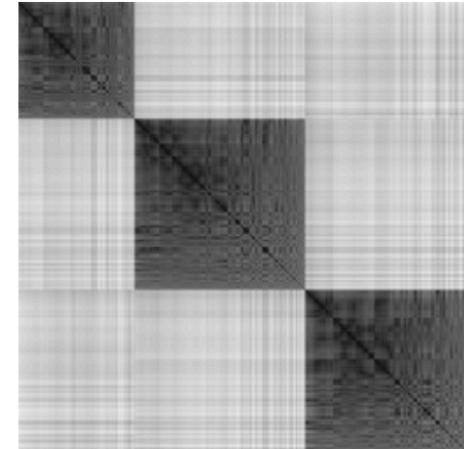
a



b



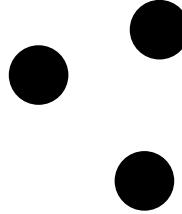
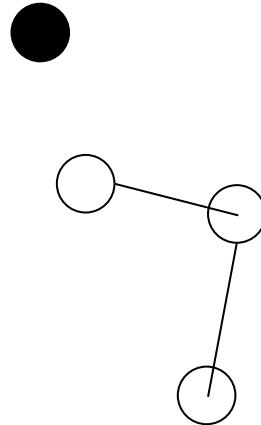
c







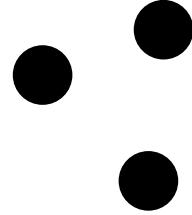
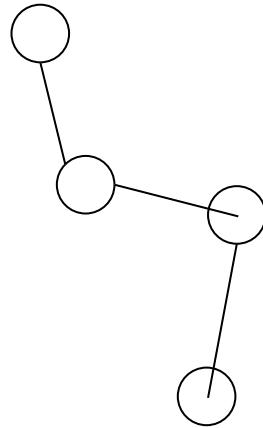


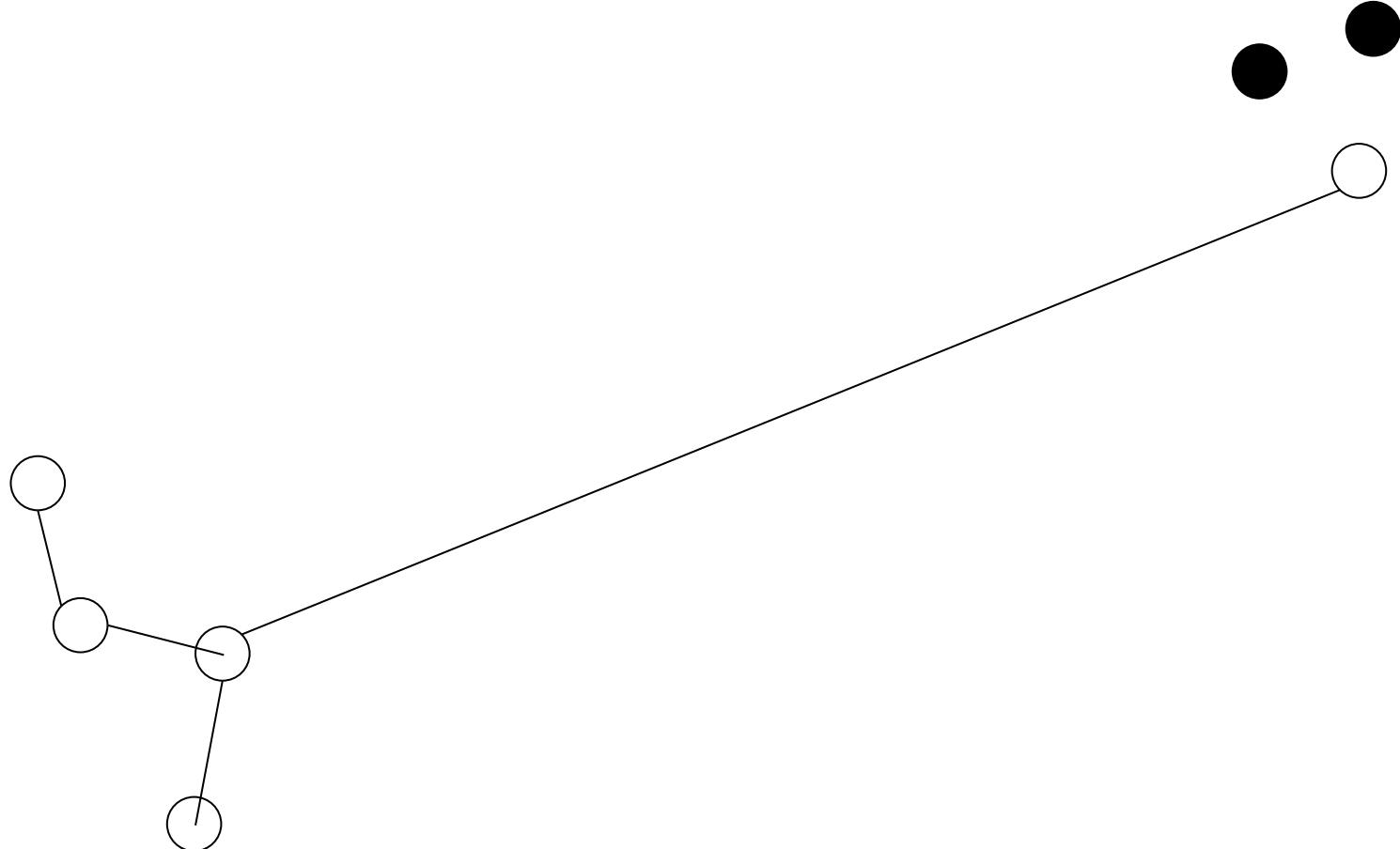


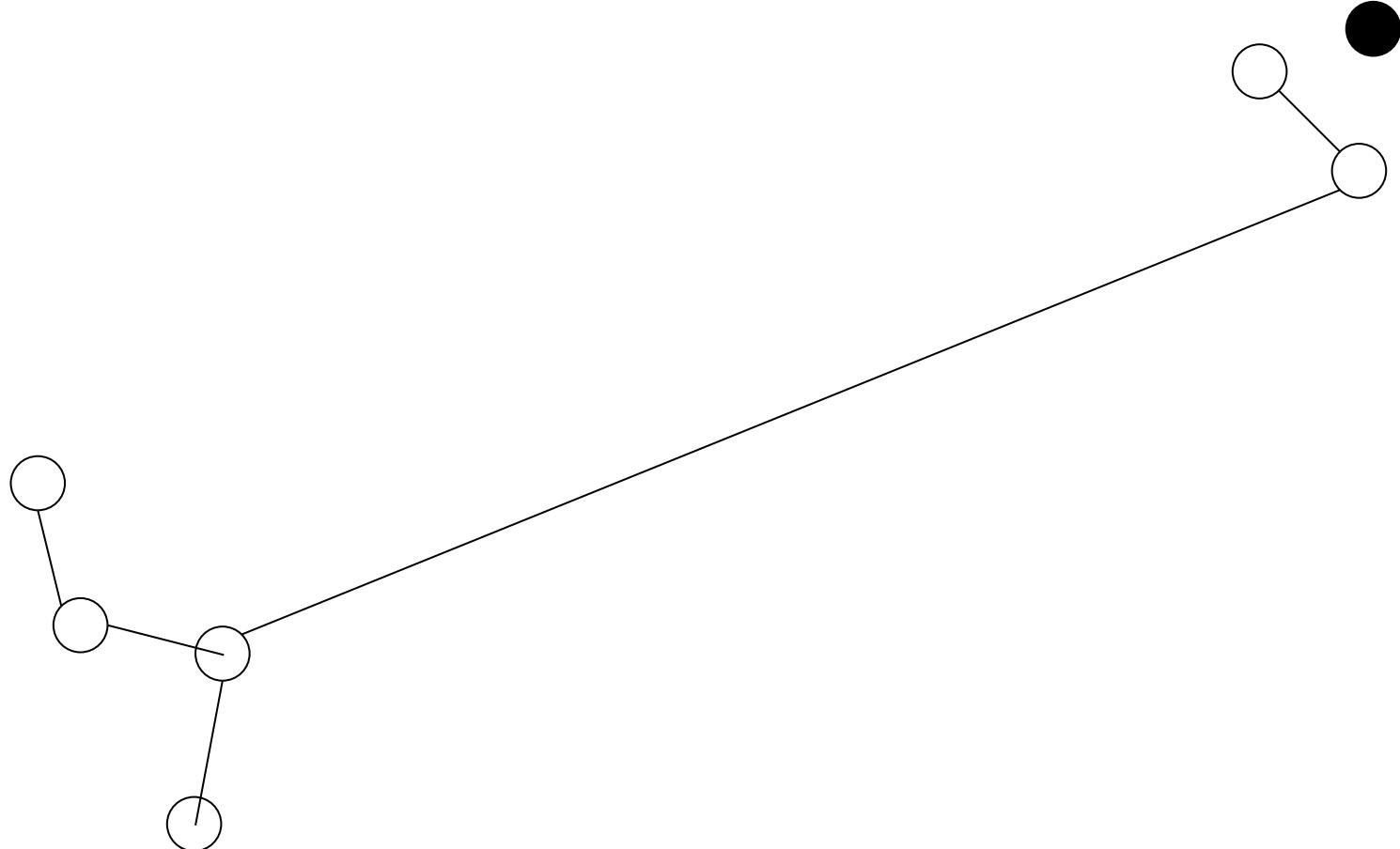


VAT ordering example

MELBOURNE

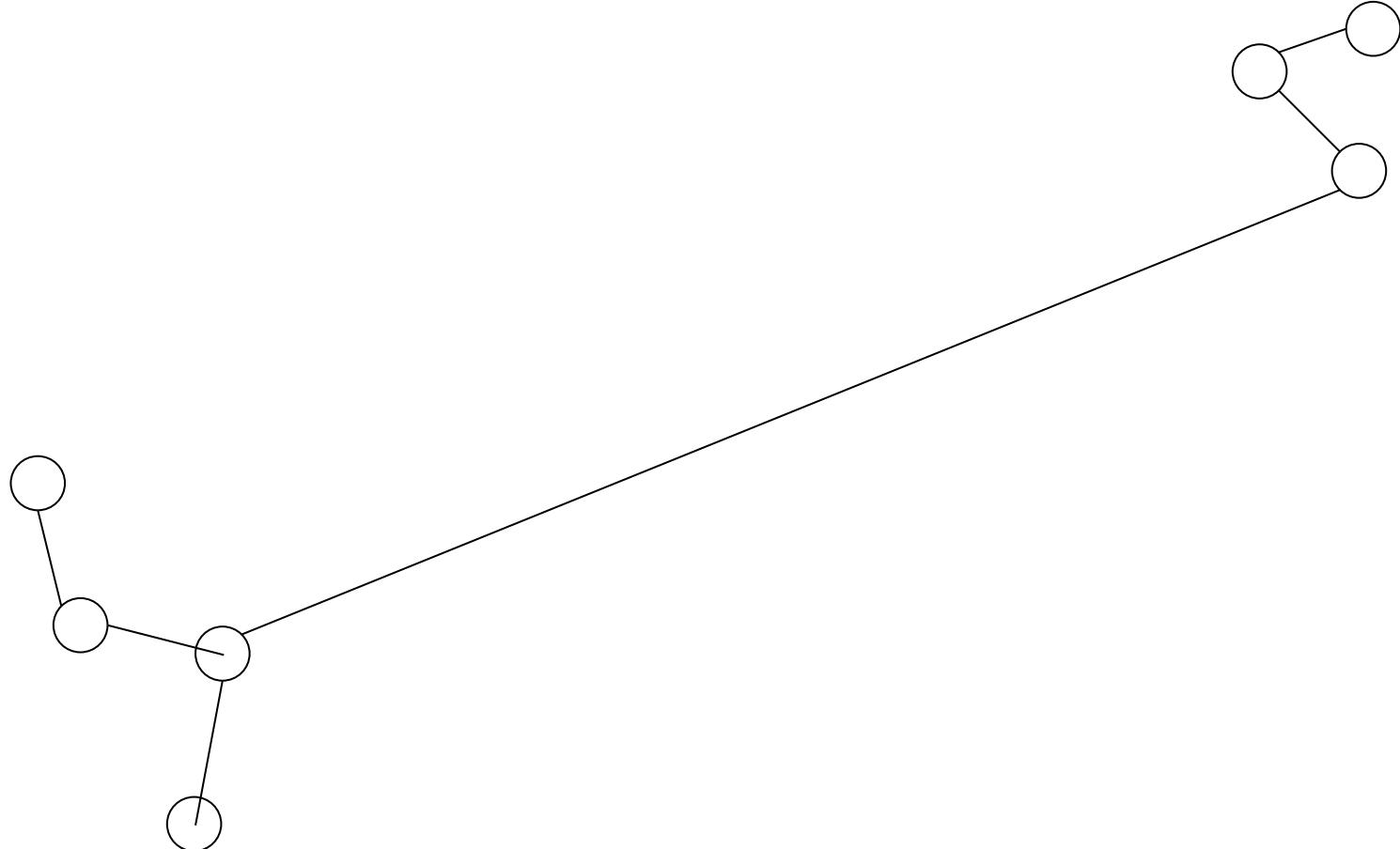






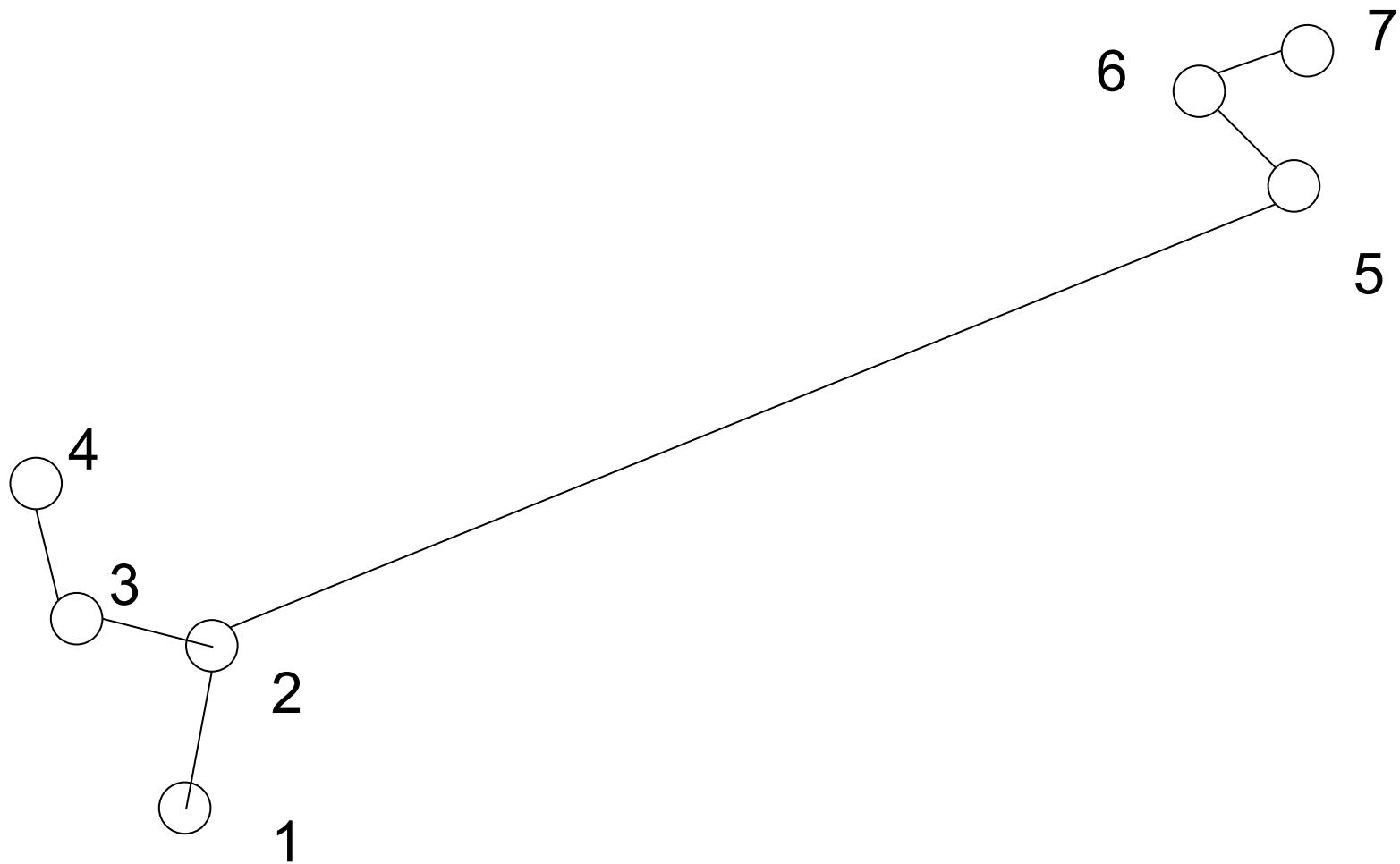


VAT ordering example





VAT ordering example





Given an $N \times N$ dissimilarity matrix \mathbf{D}

Let $K = \{1, \dots, N\}$, $I = J = \{\}$

Pick the two least similar objects o_a and o_b from \mathbf{D}

$P(1) = a; I = \{a\}; J = K - \{a\}$

For $r = 2, \dots, N$

Select (i, j) : pair of most similar objects o_i and o_j from \mathbf{D}

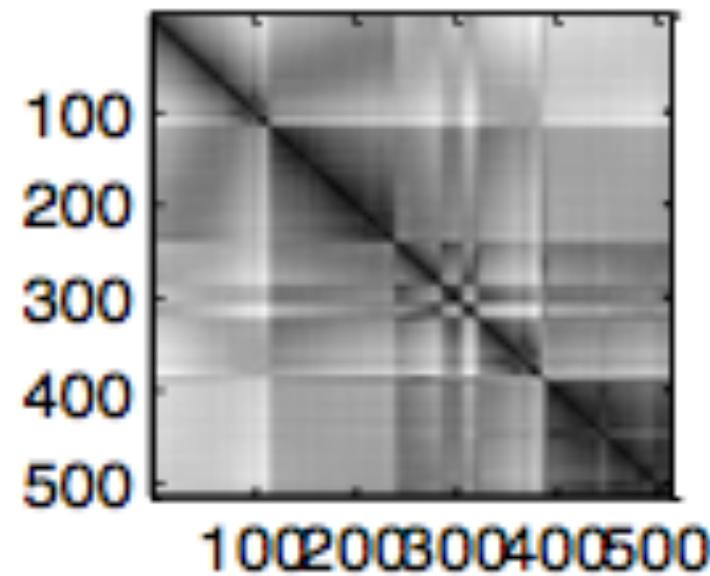
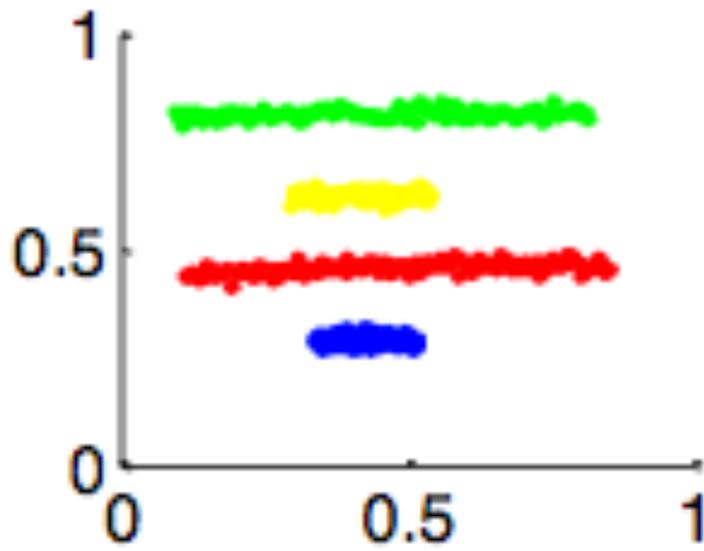
Such that $i \in I, j \in J$

$P(r) = j; I = I \cup \{j\}; J = J - \{j\};$

Obtain reordered dissimilarity matrix \mathbf{D}^* from permutation P



- VAT algorithm won't be effective in every situation
 - For complex shaped datasets (either significant overlap or irregular geometries between different clusters), the quality of the VAT image may significantly degrade.



- You will practice in workshop



- An application of VAT, role discovery in company data
<http://www.youtube.com/watch?v=l3tkUpGTTmQ&authuser=0>
- In this video, they represent each employee by a vector describing their level of access to various organisational entities

Employee	Entity 1	Entity 2	...	Entity N
James	Yes	No	Yes
Bob	No	No	Yes

Can also represent each entity by a vector describing which employees have access to it

Entity	James	Bob	...	Kate
1	Yes	No		No

- Material partly adapted from
 - “Data Mining Concepts and Techniques”, Han et al, 2nd edition 2006.