

Unsupervised Learning and Clustering

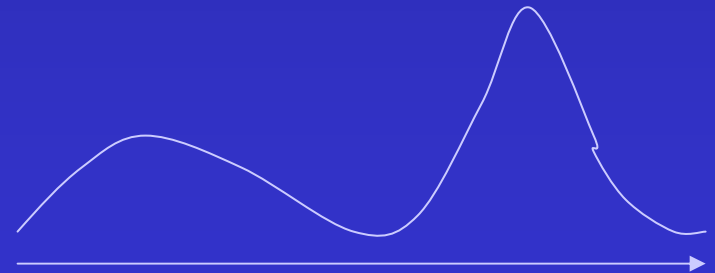
Why consider unlabeled samples?

1. Collecting and labeling large set of samples is costly
 - Getting recorded speech is free, labeling is time consuming
2. Classifier could be designed on small set of labeled samples and tuned on a large unlabeled set
3. Train on large unlabeled set and use supervision on groupings found
4. Characteristics of patterns may change with time
5. Unsupervised methods can be used to find useful features
6. Exploratory data analysis may discover presence of significant subclasses affecting design

Mixture Densities and Identifiability

- Samples come from c classes
- Priors are known $P(\omega_j)$
- Forms of the class-conditional are known
- Values for their parameters are **unknown**
- Probability density function of samples is:

$$p(\mathbf{x}, \theta) = \sum_{j=1}^c p(\mathbf{x} | \omega_j, \theta_j) P(\omega_j)$$



Gradient Ascent for Mixtures

- Mixture density:

$$p(x, \theta) = \sum_{j=1}^c p(x | \omega_j, \theta_j) P(\omega_j)$$

- Likelihood of observed samples:

$$p(D | \theta) = \prod_{k=1}^n p(x_k | \theta)$$

- Log-likelihood:

$$l = \sum_{k=1}^n \ln p(x_k | \theta)$$

- Gradient w.r.t. θ_j :

$$\nabla_{\theta_j} l = \sum_{k=1}^n \frac{1}{p(x_k | \theta)} \nabla_{\theta_j} \left[\sum_{j=1}^c p(x_k | \omega_j, \theta_j) P(\omega_j) \right]$$

- MLE must satisfy:

$$\sum_{k=1}^n P(\omega_i | x_k, \hat{\theta}) \nabla_{\theta_i} \ln p(x_k | \omega_i, \hat{\theta}_i) = 0$$

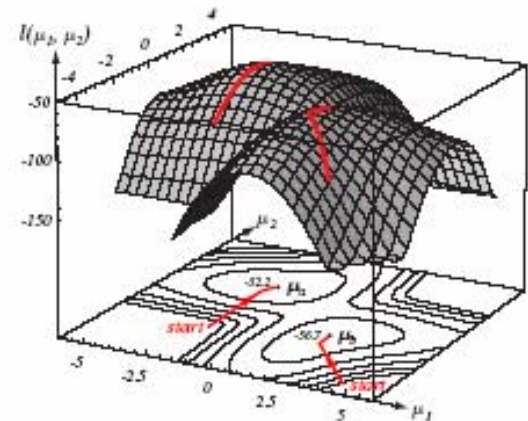
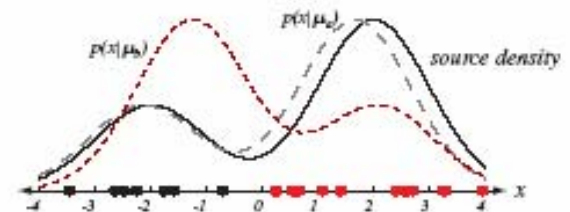
Gaussian Mixture

- Unknown mean vectors, yields

$$\hat{\mu}_i = \frac{\sum_{k=1}^n P(\omega_i | x_k, \hat{\mu}) x_k}{\sum_{k=1}^n P(\omega_i | x_k, \hat{\mu})} \text{ where } \hat{\mu} = (\hat{\mu}_1, \dots, \hat{\mu}_c)^t$$

- Leading to an iterative scheme for improving estimates

$$\hat{\mu}_i(j+1) = \frac{\sum_{k=1}^n P(\omega_i | x_k, \hat{\mu}(j)) x_k}{\sum_{k=1}^n P(\omega_i | x_k, \hat{\mu}(j))}$$



k-means clustering

- Gaussian case with all parameters unknown leads to a formulation:

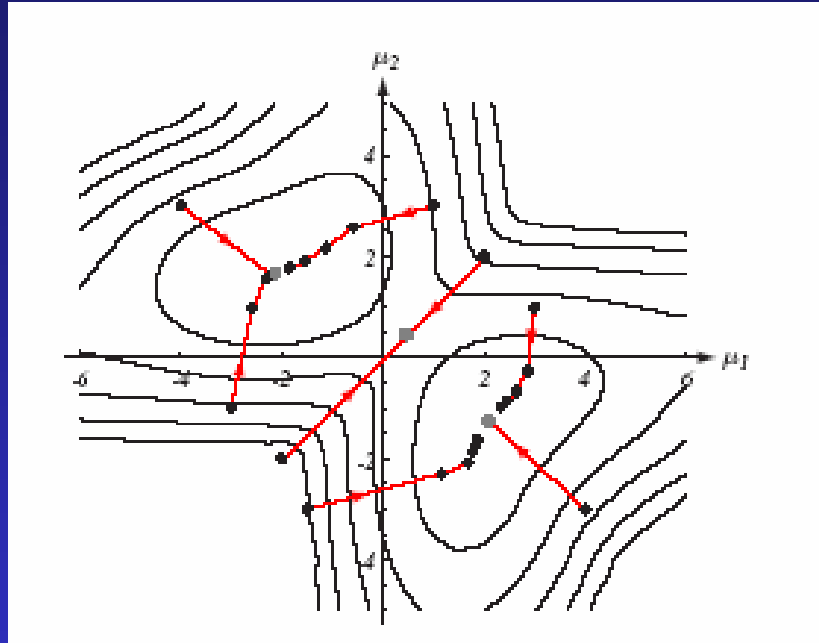
begin initialize $n, c, \mu_1, \mu_2, \dots, \mu_c$

- do classify n samples according to nearest μ_i
- recompute μ_i

until no change in μ_i

end

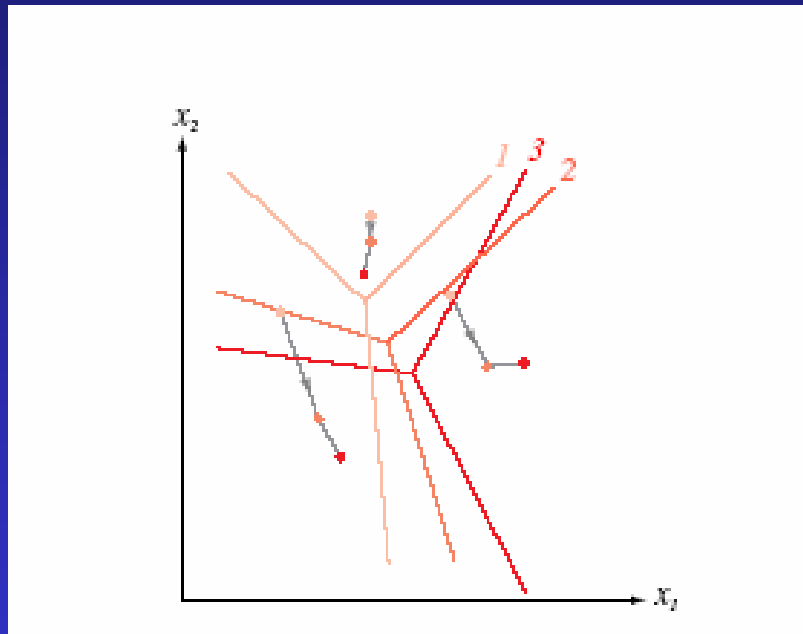
k-means clustering with one feature



One-dimensional
example

Six starting points lead local maxima whereas two for both of which $\mu_1(0) = \mu_2(0)$ lead to a saddle point

k -means clustering with two features



Two-dimensional
example

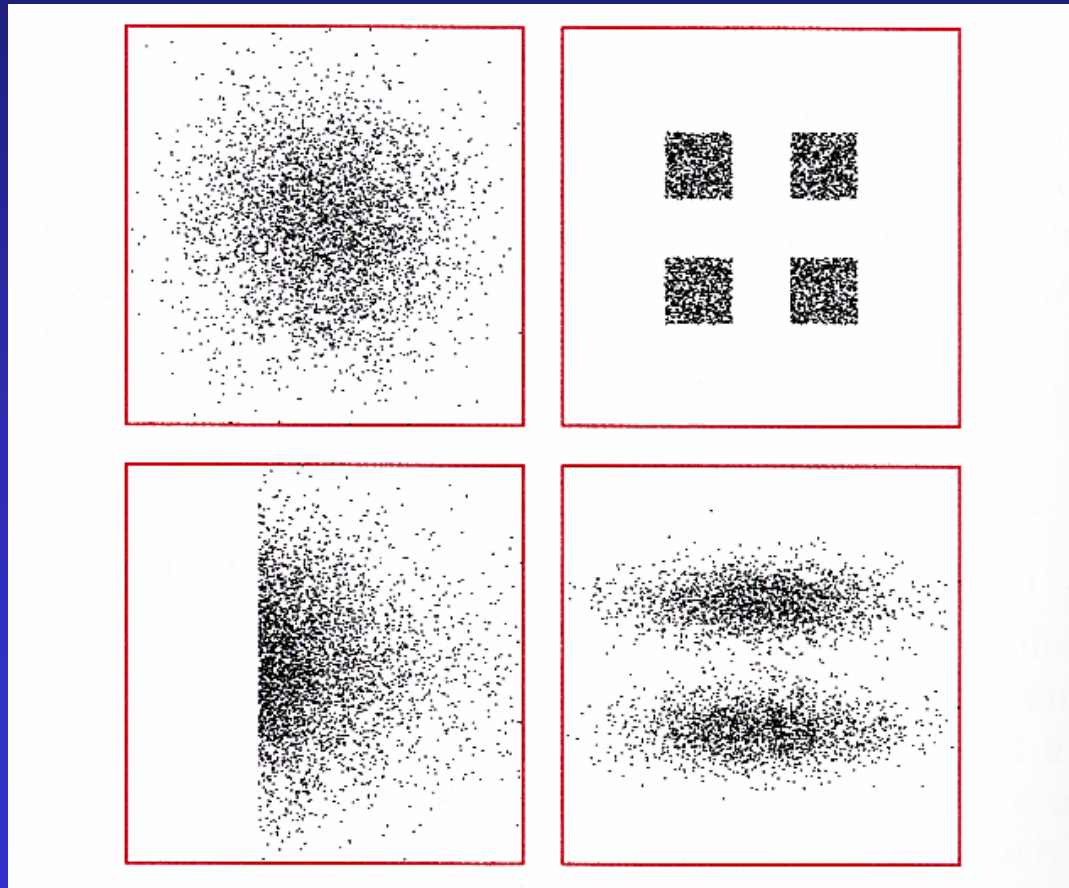
There are three means and there are three steps in the iteration. Voronoi tessellations based on means are shown

Data Description and Clustering

Data Description

- Learning the structure of multidimensional patterns from a set of unlabelled samples
- Form clouds of points in d -dimensional space
- If data were from a single normal distribution, mean and covariance metric would suffice as a description

Data sets having identical statistics upto second order, i.e., same μ and Σ

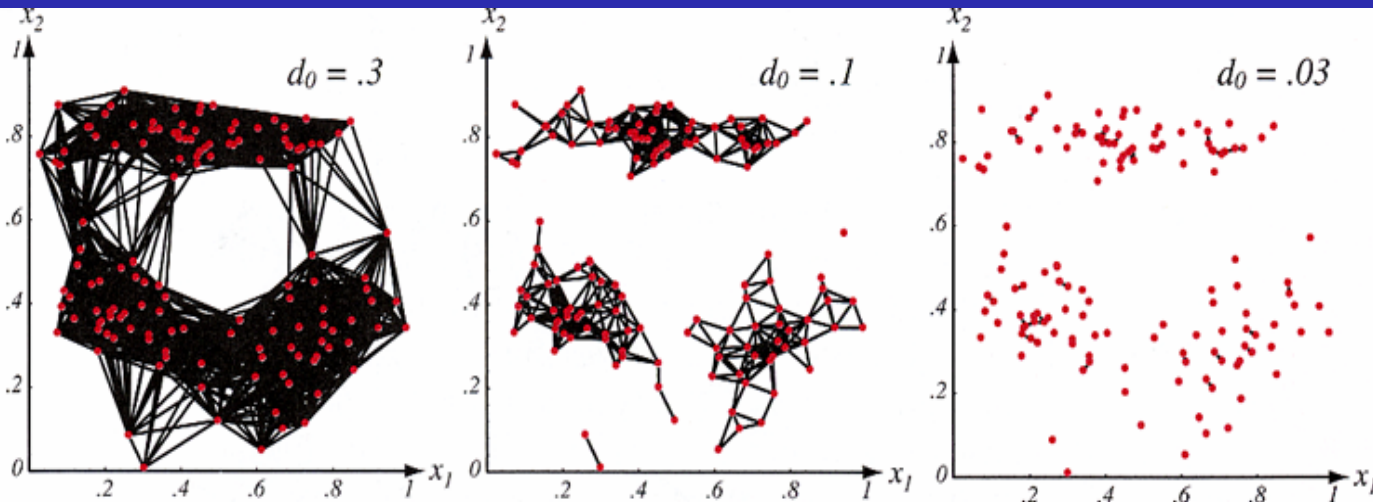


Mixture of c normal distributions approach

- Estimating parameters is non-trivial
- Assumption of particular parametric forms
 - can lead to poor or meaningless results
- Alternatively use nonparametric approach:
 - peaks or modes can indicate clusters
- If goal is to find sub-classes use clustering procedures

Similarity Measures

- Two Issues
 1. How to measure similarity between samples?
 2. How to evaluate partitioning?
- If distance is a good measure of dissimilarity
distance between samples in same cluster must be smaller
than distance between samples in different clusters
- Two samples belong to the same cluster if distance
between them is less than a threshold d_0



Distance
threshold
affects number
and size of
clusters

Similarity Measures for Clustering

- Minkowski Metric

$$d(x, x') = \left(\sum_{k=1}^d |x_k - x'_k|^q \right)^{1/q}$$

$q = 2$ is Euclidean,
 $q = 1$ is Manhattan or city block metric

- Metric based on data itself:
 - Mahalanobis distance
- Angle between vectors as similarity

$$s(x, x') = \frac{x^t x'}{\|x\| \|x'\|}$$

Cosine of angle between vectors is invariant to rotation and dilation but not translation and general linear transformations

Binary Feature Similarity Measures

$$s(x, x') = \frac{x^t x'}{\|x\| \|x'\|}$$

Numerator

no of attributes possessed by both x and x'

Denominator

$(x^t x x'^t x')^{1/2}$ is geometric mean of no of attributes possessed by x and x'

$$s(x, x') = \frac{x^t x'}{d}$$

Fraction of attributes shared

$$s(x, x') = \frac{x^t x'}{x^t x + x'^t x' - x^t x'}$$

Tanimoto coefficient: Ratio of number of shared attributes to number possessed by x or x'

Issues in Choice of Similarity Function

- Tanimoto coefficient used in Information Retrieval and Taxonomy
- Fundamental issues in Measurement Theory
- Combining features is tricky: inches versus meters
- Nominal, ordinal, interval and ratio scales

Criterion Functions for Clustering

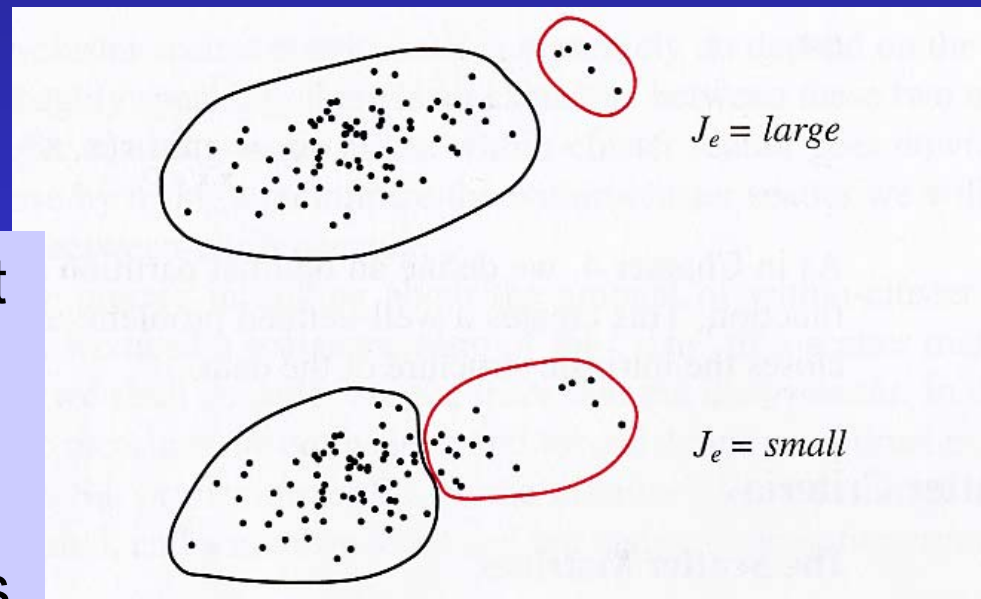
Sum of squared errors criterion

Mean of samples in D_i

$$m_i = \frac{1}{2} \sum_{x \in D_i} x$$

$$J_e = \sum_{i=1}^c \sum_{x \in D_i} \|x - m_i\|^2$$

Criterion is not best
when two clusters
are of unequal size
Suitable when they
are compact clouds



Related Minimum Variance Criteria

$$J_e = \frac{1}{2} \sum_{i=1}^c n_i \bar{s}_i$$

$$\text{where } \bar{s}_i = \frac{1}{n_i} \sum_{x \in D_i} \sum_{x' \in D_i} \underbrace{\|x - x'\|^2}_{\text{similarity function}}$$

Can be replaced by other similarity function $s(x, x')$

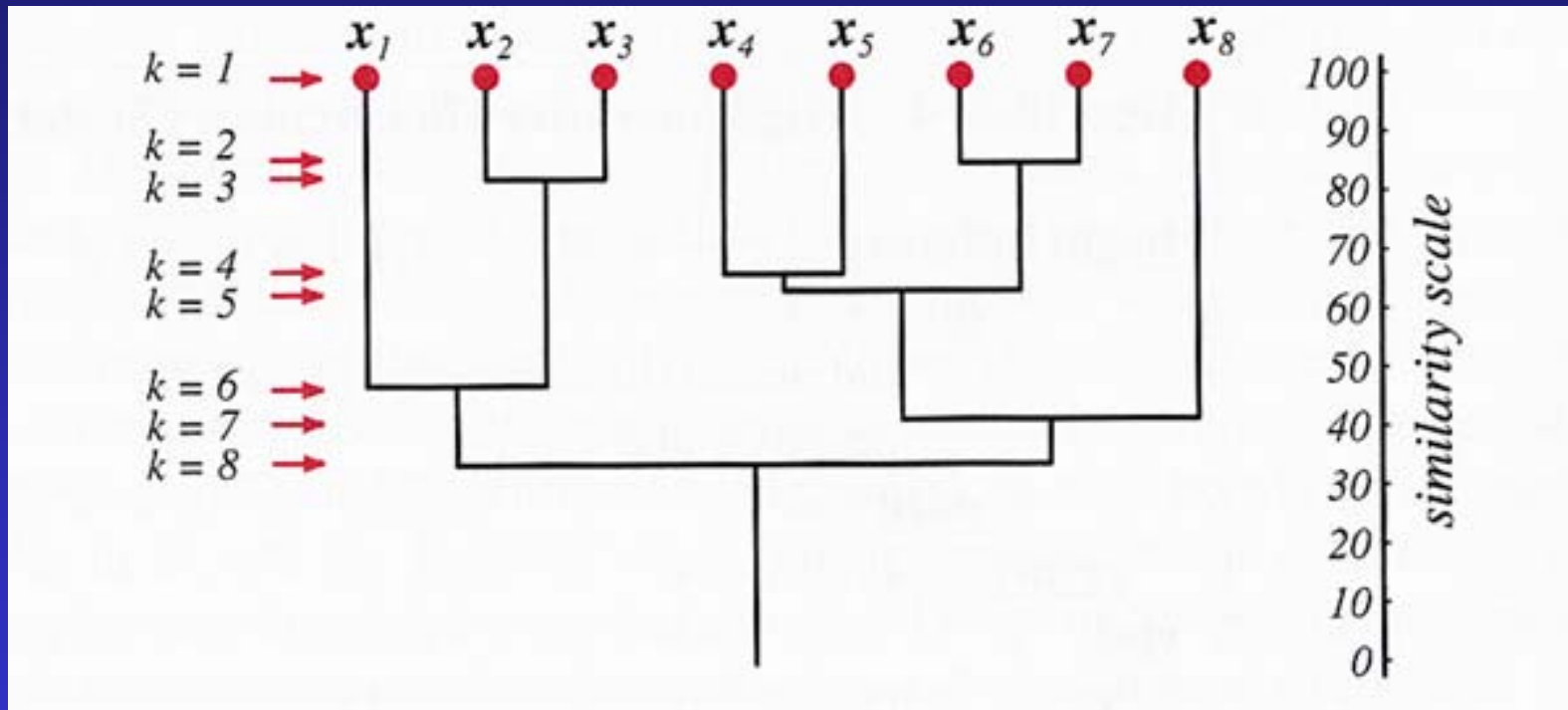
Optimal partition extremizes the criterion function

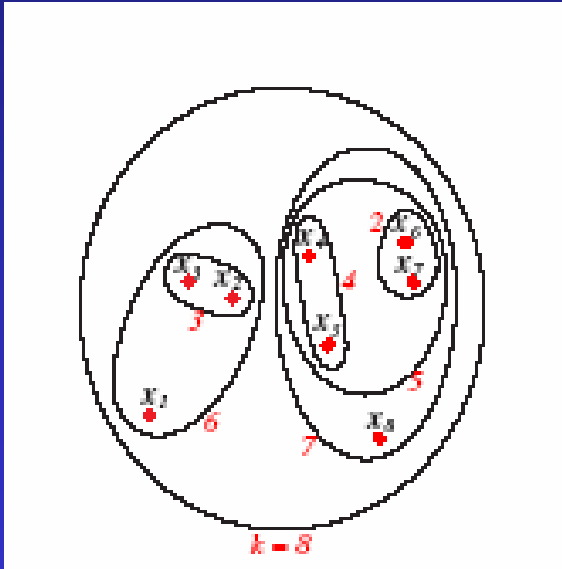
Scatter Criteria

- Derived from Scatter Matrices
- Trace criterion
- Determinant Criterion
- Invariant Criteria

Hierarchical Clustering

Dendrogram



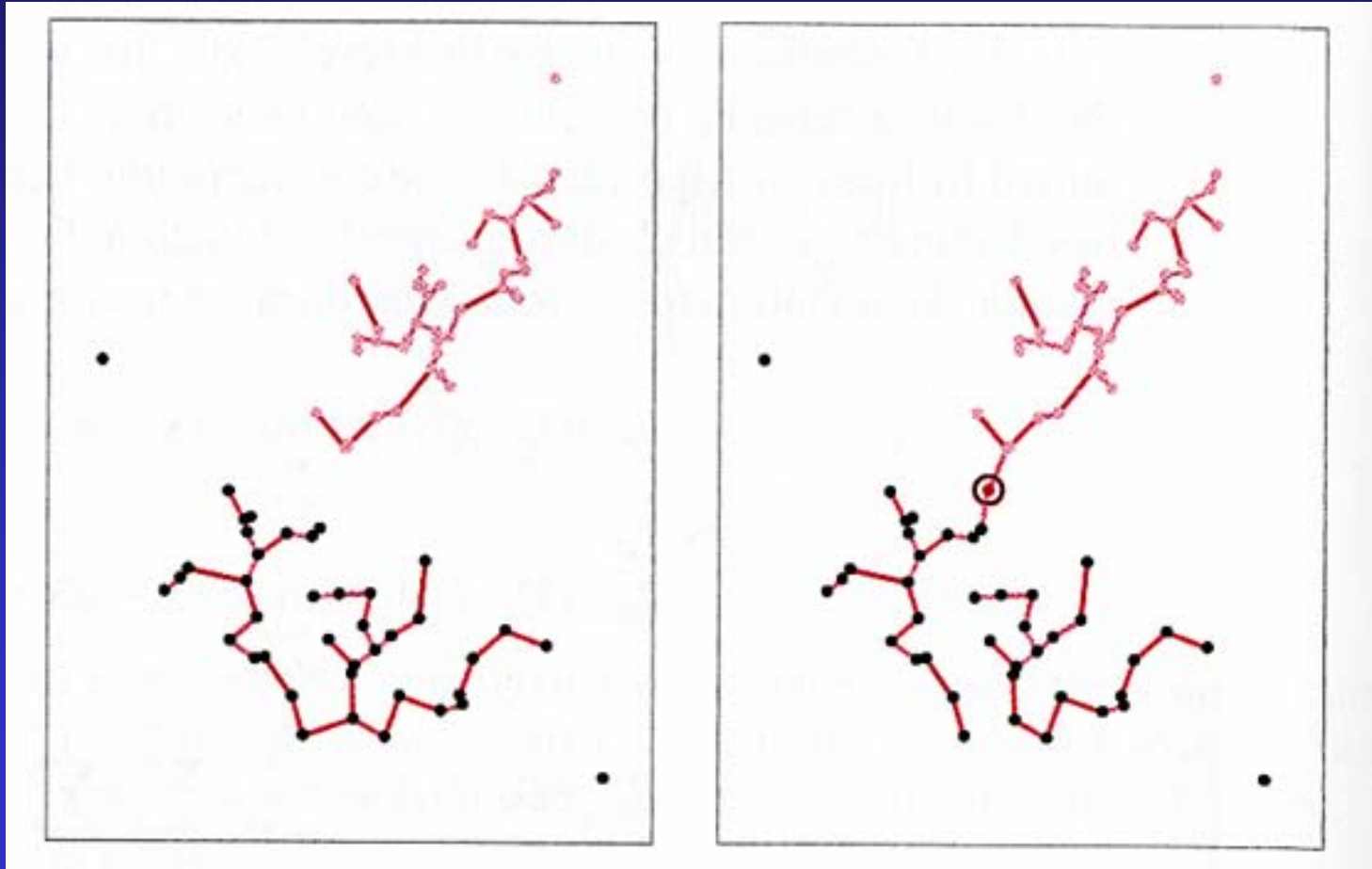


Agglomerative Algorithm

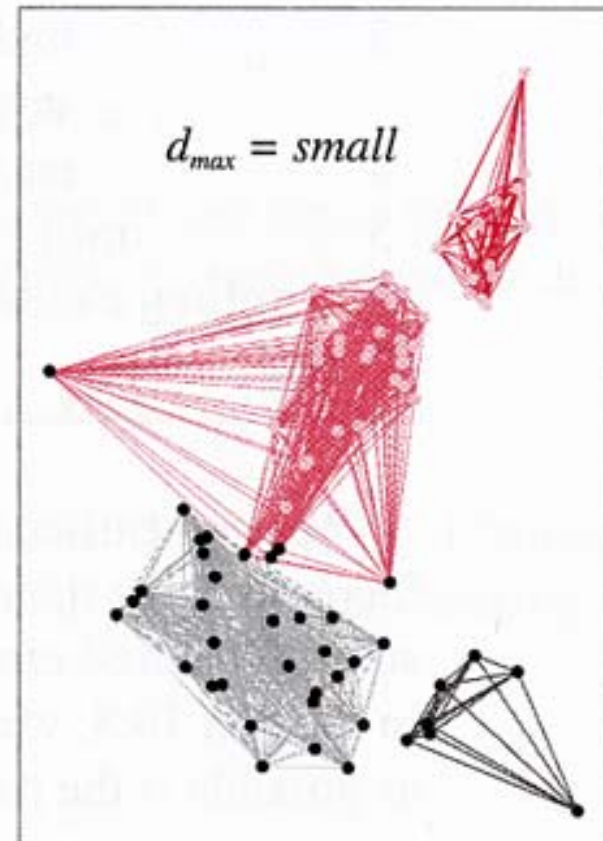
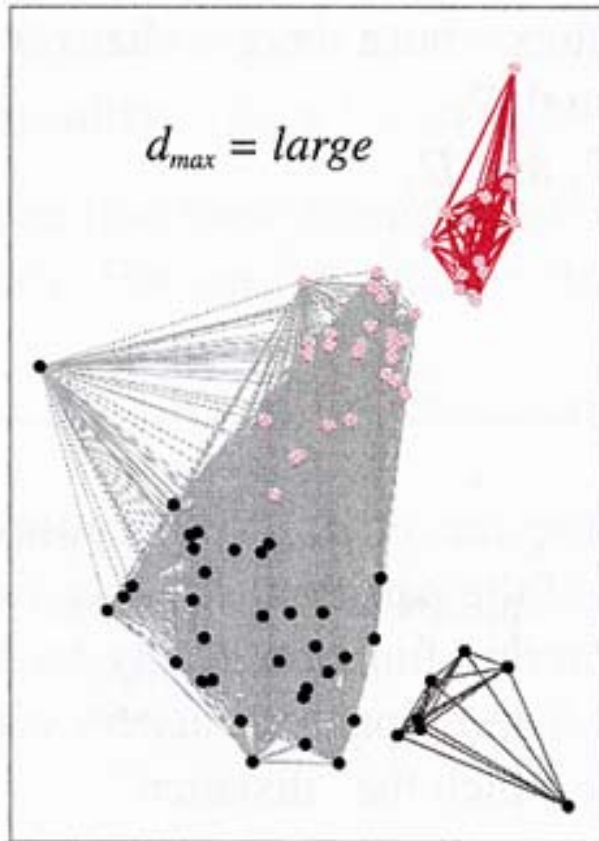
■ **Algorithm 4. (Agglomerative Hierarchical Clustering)**

```
1 begin initialize  $c, \hat{c} \leftarrow n, \mathcal{D}_i \leftarrow \{\mathbf{x}_i\}, i = 1, \dots, n$   
2           do  $\hat{c} \leftarrow \hat{c} - 1$   
3           find nearest clusters, say,  $\mathcal{D}_i$  and  $\mathcal{D}_j$   
4           merge  $\mathcal{D}_i$  and  $\mathcal{D}_j$   
5           until  $c = \hat{c}$   
6 return  $c$  clusters  
7 end
```

Nearest Neighbor Algorithm



Farthest Neighbor Algorithm



How to determine nearest clusters

$$d_{min}(\mathcal{D}_i, \mathcal{D}_j) = \min_{\substack{\mathbf{x} \in \mathcal{D}_i \\ \mathbf{x}' \in \mathcal{D}_j}} \|\mathbf{x} - \mathbf{x}'\|$$

$$d_{max}(\mathcal{D}_i, \mathcal{D}_j) = \max_{\substack{\mathbf{x} \in \mathcal{D}_i \\ \mathbf{x}' \in \mathcal{D}_j}} \|\mathbf{x} - \mathbf{x}'\|$$

$$d_{avg}(\mathcal{D}_i, \mathcal{D}_j) = \frac{1}{n_i n_j} \sum_{\mathbf{x} \in \mathcal{D}_i} \sum_{\mathbf{x}' \in \mathcal{D}_j} \|\mathbf{x} - \mathbf{x}'\|$$

$$d_{mean}(\mathcal{D}_i, \mathcal{D}_j) = \|\mathbf{m}_i - \mathbf{m}_j\|.$$