

CHAPTER 7:

# CLUSTERING

# Semiparametric Density Estimation

2

- Parametric: Assume a single model for  $p(\mathbf{x} | C_i)$  (Chapters 4 and 5)
- Semiparametric:  $p(\mathbf{x}|C_i)$  is a mixture of densities  
Multiple possible explanations/prototypes:  
Different handwriting styles, accents in speech
- Nonparametric: No model; data speaks for itself (Chapter 8)

# Mixture Densities

3

$$p(\mathbf{x}) = \sum_{i=1}^k p(\mathbf{x} | G_i) P(G_i)$$

where  $G_i$  the components/groups/clusters,

$P(G_i)$  mixture proportions (priors),

$p(\mathbf{x} | G_i)$  component densities

Gaussian mixture where  $p(\mathbf{x}|G_i) \sim N(\boldsymbol{\mu}_i, \Sigma_i)$

parameters  $\Phi = \{P(G_i), \boldsymbol{\mu}_i, \Sigma_i\}_{i=1}^k$

unlabeled sample  $X=\{\mathbf{x}^t\}_t$  (unsupervised learning)

# Classes vs. Clusters

4

□ Supervised:  $X = \{\mathbf{x}^t, \mathbf{r}^t\}_t$

□ Classes  $C_i, i=1, \dots, K$

$$p(\mathbf{x}) = \sum_{i=1}^K p(\mathbf{x} | C_i) P(C_i)$$

where  $p(\mathbf{x} | C_i) \sim N(\boldsymbol{\mu}_i, \Sigma_i)$

□  $\Phi = \{P(C_i), \boldsymbol{\mu}_i, \Sigma_i\}_{i=1}^K$

$$\hat{P}(C_i) = \frac{\sum_t r_i^t}{N} \quad \mathbf{m}_i = \frac{\sum_t r_i^t \mathbf{x}^t}{\sum_t r_i^t}$$

$$\mathbf{S}_i = \frac{\sum_t r_i^t (\mathbf{x}^t - \mathbf{m}_i)(\mathbf{x}^t - \mathbf{m}_i)^T}{\sum_t r_i^t}$$

□ Unsupervised :  $X = \{\mathbf{x}^t\}_t$

□ Clusters  $G_i, i=1, \dots, k$

$$p(\mathbf{x}) = \sum_{i=1}^k p(\mathbf{x} | G_i) P(G_i)$$

where  $p(\mathbf{x} | G_i) \sim N(\boldsymbol{\mu}_i, \Sigma_i)$

□  $\Phi = \{P(G_i), \boldsymbol{\mu}_i, \Sigma_i\}_{i=1}^k$

Labels  $\mathbf{r}_i^t$  ?

# $k$ -Means Clustering

5

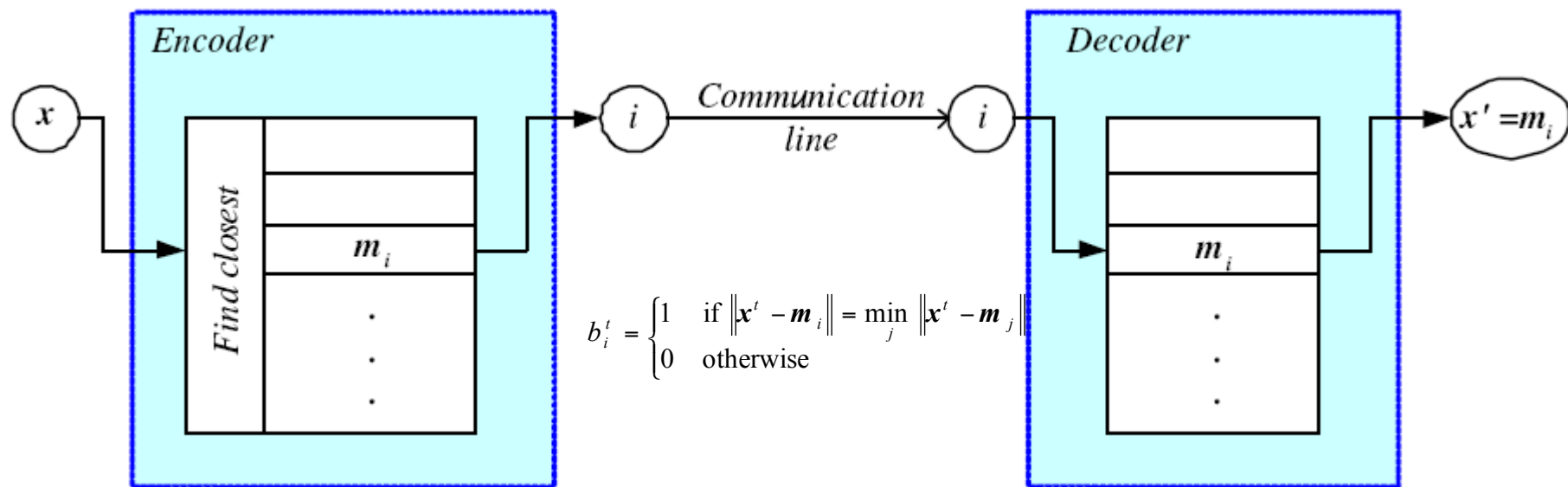
- Find  $k$  reference vectors (prototypes/codebook vectors/codewords) which best represent data
- Reference vectors,  $\mathbf{m}_j, j = 1, \dots, k$
- Use nearest (most similar) reference:

$$\|\mathbf{x}^t - \mathbf{m}_i\| = \min_j \|\mathbf{x}^t - \mathbf{m}_j\|$$

- Reconstruction error  $E(\{\mathbf{m}_i\}_{i=1}^k | \mathbf{X}) = \sum_t \sum_i b_i^t \|\mathbf{x}^t - \mathbf{m}_i\|$   
$$b_i^t = \begin{cases} 1 & \text{if } \|\mathbf{x}^t - \mathbf{m}_i\| = \min_j \|\mathbf{x}^t - \mathbf{m}_j\| \\ 0 & \text{otherwise} \end{cases}$$

# Encoding/Decoding

6



# $k$ -means Clustering

7

Initialize  $\mathbf{m}_i, i = 1, \dots, k$ , for example, to  $k$  random  $\mathbf{x}^t$

Repeat

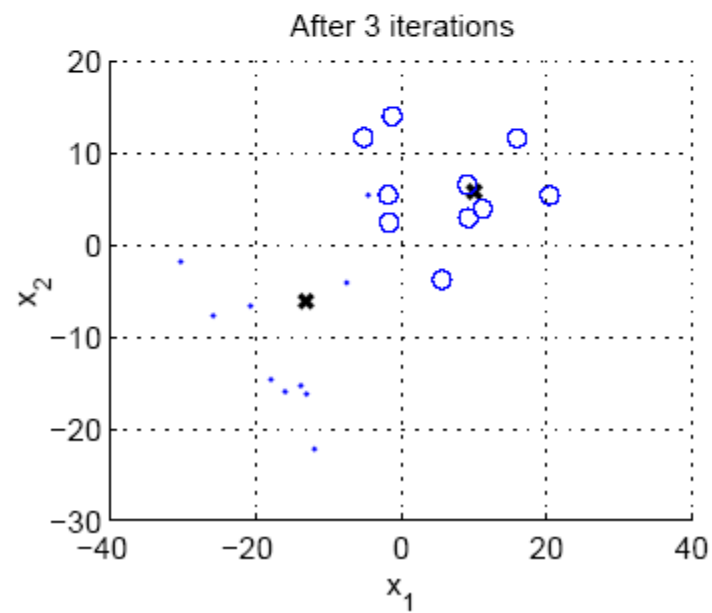
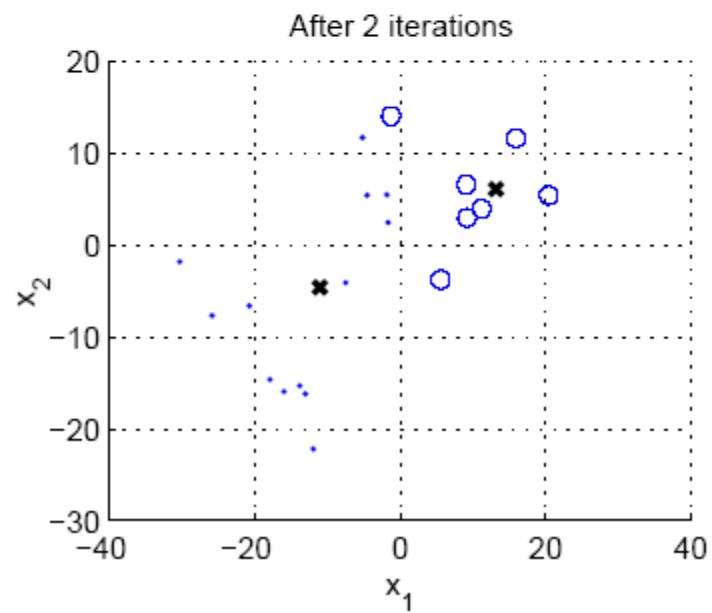
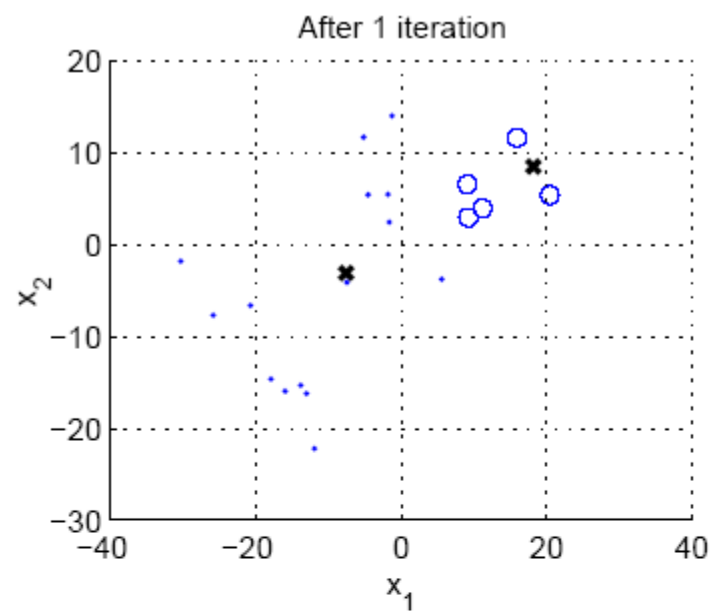
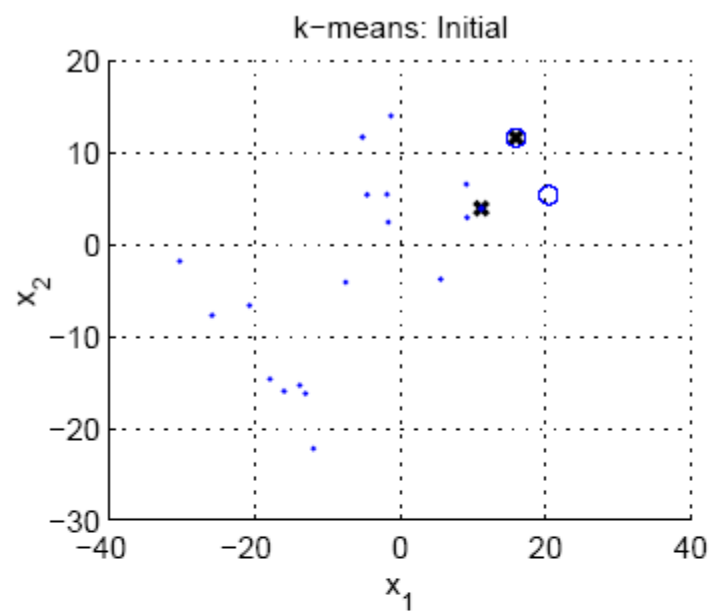
For all  $\mathbf{x}^t \in \mathcal{X}$

$$b_i^t \leftarrow \begin{cases} 1 & \text{if } \|\mathbf{x}^t - \mathbf{m}_i\| = \min_j \|\mathbf{x}^t - \mathbf{m}_j\| \\ 0 & \text{otherwise} \end{cases}$$

For all  $\mathbf{m}_i, i = 1, \dots, k$

$$\mathbf{m}_i \leftarrow \sum_t b_i^t \mathbf{x}^t / \sum_t b_i^t$$

Until  $\mathbf{m}_i$  converge





# Expectation-Maximization (EM)

9

- Log likelihood with a mixture model

$$\begin{aligned} L(\Phi | X) &= \log \prod_t p(\mathbf{x}^t | \Phi) \\ &= \sum_t \log \sum_{i=1}^k p(\mathbf{x}^t | G_i) P(G_i) \end{aligned}$$

- Assume hidden variables  $z$ , which when known, make optimization much simpler
- Complete likelihood,  $L_c(\Phi | X, Z)$ , in terms of  $x$  and  $z$
- Incomplete likelihood,  $L(\Phi | X)$ , in terms of  $x$

# E- and M-steps

10

Iterate the two steps:

1. E-step: Estimate  $z$  given  $X$  and current  $\Phi$
2. M-step: Find new  $\Phi'$  given  $z$ ,  $X$ , and old  $\Phi$ .

$$\text{E - step : } Q(\Phi | \Phi^l) = E[\mathcal{L}_C(\Phi | X, Z) | X, \Phi^l]$$

$$\text{M - step : } \Phi^{l+1} = \arg \max_{\Phi} Q(\Phi | \Phi^l)$$

$Q$  is the expectation of the complete likelihood given  $X$  and current parameters  $\Phi^l$

An increase in  $Q$  increases incomplete likelihood

$$L(\Phi^{l+1} | X) \geq L(\Phi^l | X)$$

# EM in Gaussian Mixtures

11

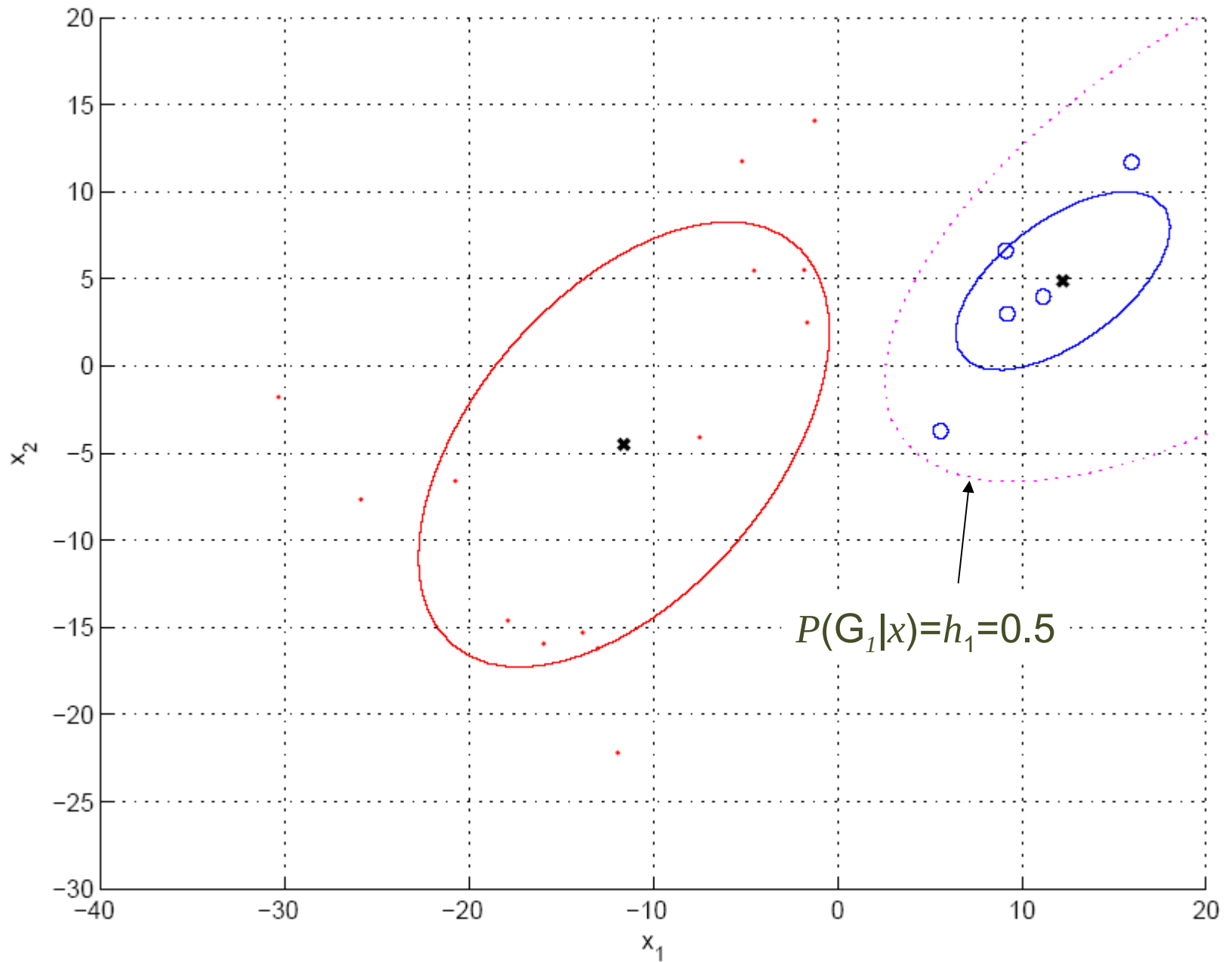
- $z_i^t = 1$  if  $\mathbf{x}^t$  belongs to  $G_i$ , 0 otherwise (labels  $\mathbf{r}_i^t$  of supervised learning); assume  $p(\mathbf{x} | G_i) \sim \mathcal{N}(\boldsymbol{\mu}_i, \Sigma_i)$

- E-step: 
$$E[z_i^t | \mathbf{X}, \Phi^l] = \frac{p(\mathbf{x}^t | G_i, \Phi^l) P(G_i)}{\sum_j p(\mathbf{x}^t | G_j, \Phi^l) P(G_j)}$$
$$= P(G_i | \mathbf{x}^t, \Phi^l) \equiv h_i^t$$

- M-step: 
$$P(G_i) = \frac{\sum_t h_i^t}{N} \quad \mathbf{m}_i^{l+1} = \frac{\sum_t h_i^t \mathbf{x}^t}{\sum_t h_i^t}$$
$$\mathbf{S}_i^{l+1} = \frac{\sum_t h_i^t (\mathbf{x}^t - \mathbf{m}_i^{l+1})(\mathbf{x}^t - \mathbf{m}_i^{l+1})^T}{\sum_t h_i^t}$$

*Use estimated labels in place of unknown labels*

EM solution



# Mixtures of Latent Variable Models

13

Regularize clusters, to avoid overfitting when attribute dimensionality is much larger than sample size

- Assume shared/diagonal covariance matrices

but these assumptions may not be appropriate for a given dataset,

OR

- Use PCA/FA in the clusters to decrease dimensionality:

Mixtures of PCA/FA

$$p(\mathbf{x}_t | G_i) = \mathcal{N}(\mathbf{x}_t | \mathbf{m}_i, \mathbf{V}_i \mathbf{V}_i^T + \boldsymbol{\Psi}_i)$$

where  $\mathbf{V}_i$  and  $\boldsymbol{\Psi}_i$  are factor loadings and variances of cluster  $G_i$

Can use EM to learn  $\mathbf{V}_i$  (Ghahramani and Hinton, 1997; Tipping and Bishop, 1999)

# After Clustering

14

- Dimensionality reduction methods find correlations between features and group features
  - Clustering methods find similarities between instances and group instances
  - Allows knowledge extraction through
    - number of clusters,
    - prior probabilities,
    - cluster parameters, i.e., center, range of features.
- Example: CRM, customer segmentation

# Clustering as Preprocessing

15

- Estimated group labels  $h_j$  (soft) or  $b_j$  (hard) may be seen as the dimensions of a new  $k$  dimensional space, where we can then learn our discriminant or regressor.
- **Local** representation (only one  $b_j$  is 1, all others are 0; only few  $h_j$  are nonzero) vs **Distributed** representation (After PCA; all  $z_j$  are nonzero)

# Mixture of Mixtures

16

- In classification, the input comes from a mixture of classes (supervised).
- If each class is also a mixture, e.g., of Gaussians, (unsupervised), we have a mixture of mixtures:

$$p(\mathbf{x} | \mathcal{C}_i) = \sum_{j=1}^{k_i} p(\mathbf{x} | G_{ij}) P(G_{ij})$$

$$p(\mathbf{x}) = \sum_{i=1}^K p(\mathbf{x} | \mathcal{C}_i) P(\mathcal{C}_i)$$



# Spectral Clustering

17

- Cluster using predefined pairwise similarities  $B_{rs}$  instead of using Euclidean or Mahalanobis distance
- Can be used even if instances not vectorially represented
- Steps:
  - I. Use Laplacian Eigenmaps (chapter 6) to map to a new  $\mathbf{z}$  space using  $B_{rs}$
  - II. Use  $k$ -means in this new  $\mathbf{z}$  space for clustering

# Hierarchical Clustering

18

- Cluster based on similarities/distances
- Distance measure between instances  $\mathbf{x}^r$  and  $\mathbf{x}^s$

Minkowski ( $L_p$ ) (Euclidean for  $p = 2$ )

$$d_m(\mathbf{x}^r, \mathbf{x}^s) = \left[ \sum_{j=1}^d \left( x_j^r - x_j^s \right)^p \right]^{1/p}$$

City-block distance

$$d_{cb}(\mathbf{x}^r, \mathbf{x}^s) = \sum_{j=1}^d |x_j^r - x_j^s|$$

# Agglomerative Clustering

19

- Start with  $N$  groups each with one instance and merge two closest groups at each iteration
- Distance between two groups  $G_i$  and  $G_j$ :

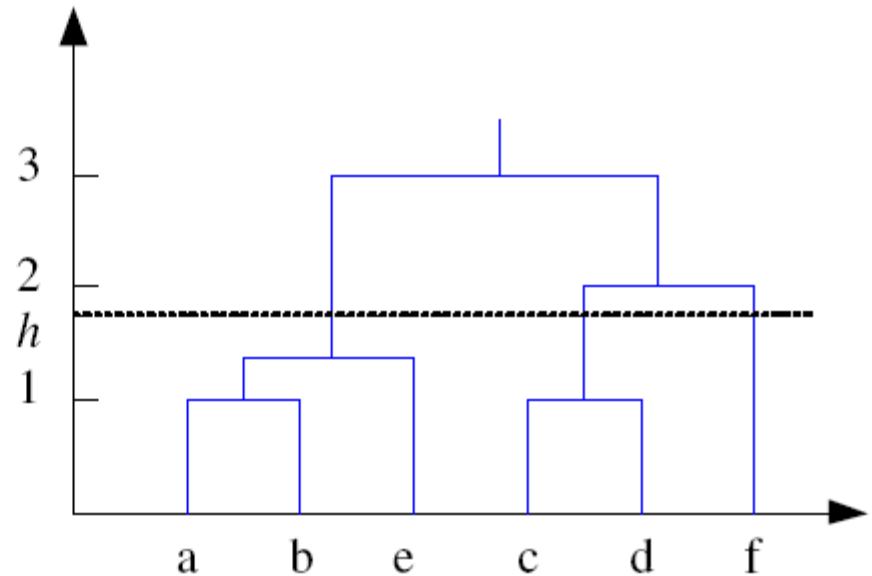
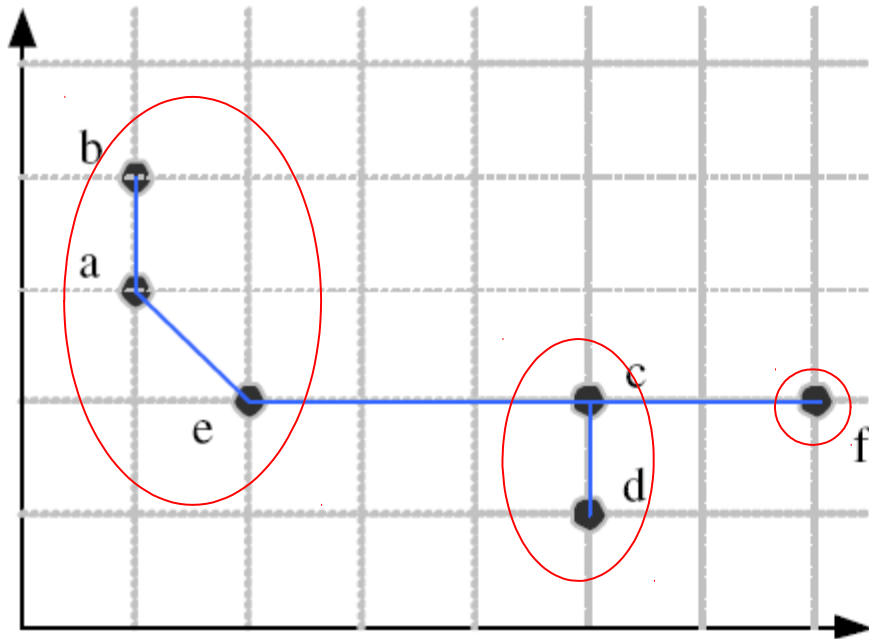
- Single-link: 
$$d(G_i, G_j) = \min_{\mathbf{x}^r \in G_i, \mathbf{x}^s \in G_j} d(\mathbf{x}^r, \mathbf{x}^s)$$

- Complete-link: 
$$d(G_i, G_j) = \max_{\mathbf{x}^r \in G_i, \mathbf{x}^s \in G_j} d(\mathbf{x}^r, \mathbf{x}^s)$$

- Average-link, centroid 
$$d(G_i, G_j) = \text{ave}_{\mathbf{x}^r \in G_i, \mathbf{x}^s \in G_j} d(\mathbf{x}^r, \mathbf{x}^s)$$

# Example: Single-Link Clustering

20



*Dendrogram*

# Choosing $k$

21

- Defined by the application, e.g., image quantization
- Plot data (after PCA) and check for clusters
- Incremental (leader-cluster) algorithm: Add one at a time until “elbow” (reconstruction error/log likelihood/intergroup distances)
- Manually check for meaning