*Lecture Slides for*

INTRODUCTION TO

*Machine Learning*

ETHEM ALPAYDIN
© The MIT Press, 2004

*alpaydin@boun.edu.tr*
*http://www.cmpe.boun.edu.tr/~ethem/i2ml*

CHAPTER 7:

*Clustering*

# *Semiparametric Density Estimation*

- **Parametric:** Assume a single model for $p(x \mid C_i)$ (Chapter 4 and 5)
- **Semiparametric:** $p(x \mid C_i)$ 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*

$$p(\boldsymbol{x}) = \sum_{i=1}^{k} p(\boldsymbol{x} \mid \mathcal{G}_i) P(\mathcal{G}_i)$$

where $\mathcal{G}_i$ the components/groups/clusters,
$P(\mathcal{G}_i)$ mixture proportions (priors),
$p(\boldsymbol{x} \mid \mathcal{G}_i)$ component densities

Gaussian mixture where $p(\boldsymbol{x} \mid \mathcal{G}_i) \sim \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$
parameters $\Phi = \{P(\mathcal{G}_i), \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i\}_{i=1}^{k}$
unlabeled sample $\mathcal{X} = \{\boldsymbol{x}^t\}_t$ (unsupervised learning)

# *Classes vs. Clusters*

- Supervised: $\mathcal{X} = \{\, \boldsymbol{x}^t, \boldsymbol{r}^t \}_t$
- Classes $C_i\ i=1,\dots,K$

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

where $p(\,\boldsymbol{x} \mid C_i) \sim \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$

- $\Phi = \{P(C_i), \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i\}^K_{i=1}$

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

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

- Unsupervised : $\mathcal{X} = \{\, \boldsymbol{x}^t \}_t$
- Clusters $\mathcal{G}_i\ i=1,\dots,k$

$$p(\boldsymbol{x}) = \sum_{i=1}^{k} p(\boldsymbol{x} \mid \mathcal{G}_i) P(\mathcal{G}_i)$$

where $p(\,\boldsymbol{x} \mid \mathcal{G}_i) \sim \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$

- $\Phi = \{P(\mathcal{G}_i), \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i\}^k_{i=1}$

Labels, $\boldsymbol{r}^t_i$ ?

# *k-Means Clustering*

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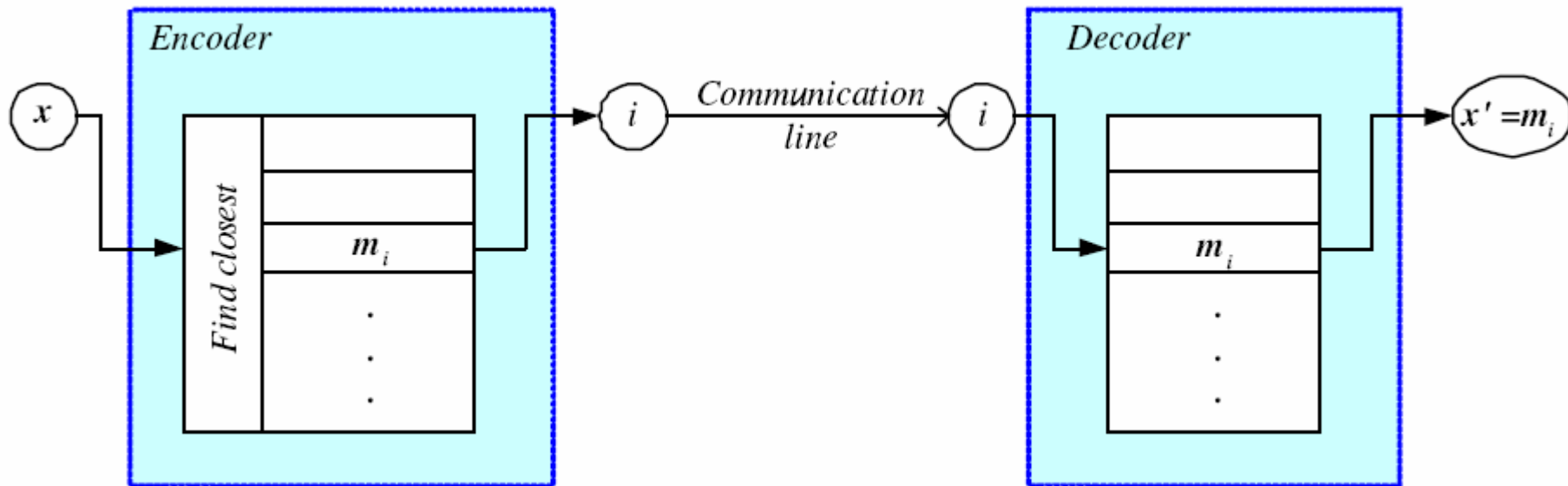- Use nearest (most similar) reference:

$$\left\|\boldsymbol{x}^t - \boldsymbol{m}_i\right\| = \min_j \left\|\boldsymbol{x}^t - \boldsymbol{m}_j\right\|$$

- Reconstruction error

$$E\left(\{\boldsymbol{m}_i\}_{i=1}^k \mid \mathcal{X}\right) = \sum_t \sum_i b_i^t \left\|\boldsymbol{x}^t - \boldsymbol{m}_i\right\|$$

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

# Encoding/Decoding



$$b_i^t = \begin{cases} 1 & \text{if } \left\| x^t - m_i \right\| = \min_j \left\| x^t - m_j \right\| \\ 0 & \text{otherwise} \end{cases}$$

# k-means Clustering

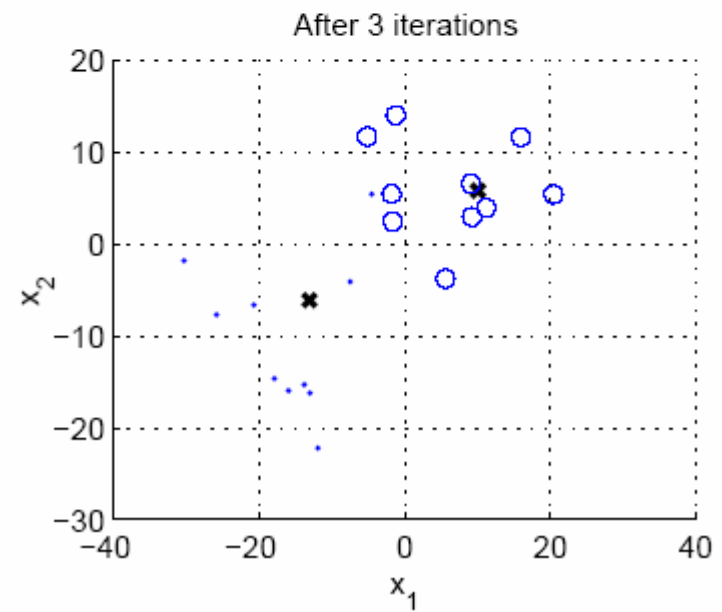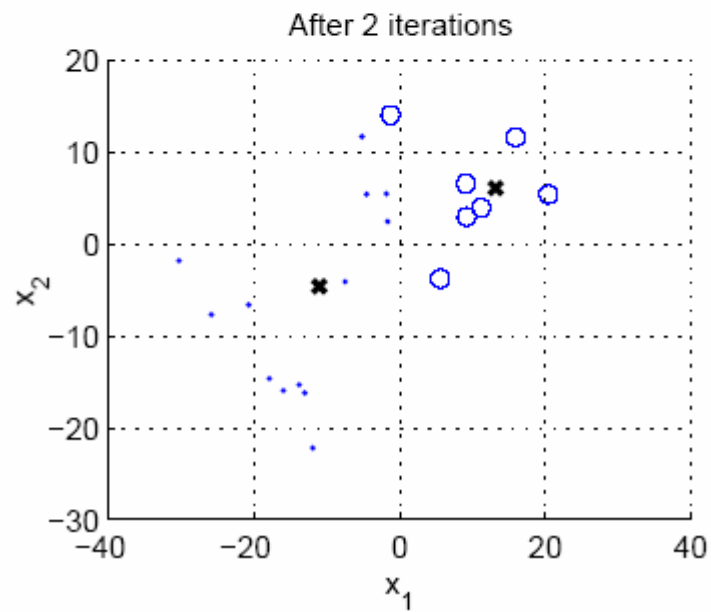Initialize $\boldsymbol{m}_i, i = 1, \ldots, k$, for example, to $k$ random $\boldsymbol{x}^t$
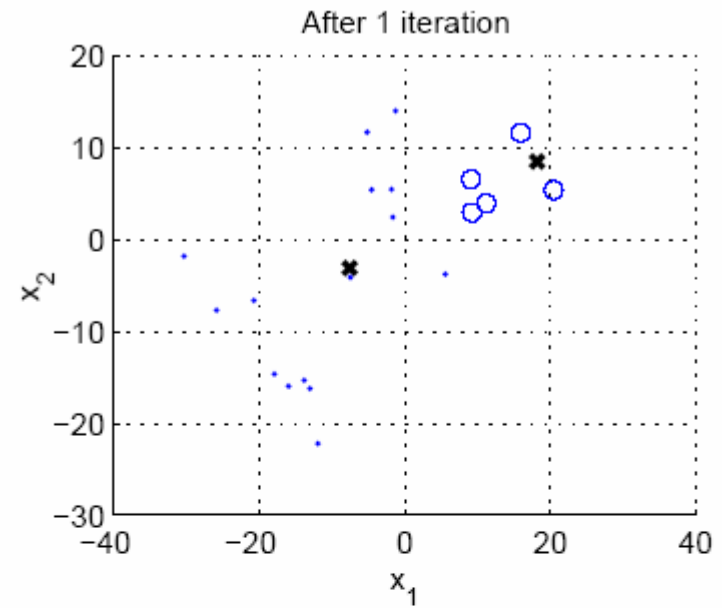
Repeat

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

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

    For all $\boldsymbol{m}_i, i = 1, \ldots, k$

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

Until $\boldsymbol{m}_i$ converge

*Lecture Notes for E ALPAYDIN 2004 Introduction to Machine Learning © The MIT Press (V1.1)*

# *Expectation-Maximization (EM)*

- Log likelihood with a mixture model

$$\mathcal{L}(\Phi \mid \mathcal{X}) = \log \prod_t p(\mathbf{x}^t \mid \Phi)$$

$$= \sum_t \log \sum_{i=1}^{k} p(\mathbf{x}^t \mid \mathcal{G}_i) P(\mathcal{G}_i)$$

- Assume hidden variables $z$, which when known, make optimization much simpler

- Complete likelihood, $\mathcal{L}_c(\Phi \mid \mathcal{X}, \mathcal{Z})$, in terms of $\mathbf{x}$ and $\mathbf{z}$

- Incomplete likelihood, $\mathcal{L}(\Phi \mid \mathcal{X})$, in terms of $\mathbf{x}$

# E- and M-steps

- Iterate the two steps
1. E-step: Estimate $z$ given $\mathcal{X}$ and current $\Phi$
2. M-step: Find new $\Phi'$ given $z$, $\mathcal{X}$, and old $\Phi$.

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

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

An increase in $\mathcal{Q}$ increases incomplete likelihood

$$\mathcal{L}\left(\Phi^{l+1} \mid \mathcal{X}\right) \geq \mathcal{L}\left(\Phi^l \mid \mathcal{X}\right)$$

# *EM in Gaussian Mixtures*

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

- E-step:

$$E\!\left[z_i^t | \mathcal{X}, \Phi^l\right] = \frac{p\!\left(\boldsymbol{x}^t \mid \mathcal{G}_i, \Phi^l\right) P\!\left(\mathcal{G}_i\right)}{\sum_j p\!\left(\boldsymbol{x}^t \mid \mathcal{G}_j, \Phi^l\right) P\!\left(\mathcal{G}_j\right)}$$

$$= P\!\left(\mathcal{G}_i \mid \boldsymbol{x}^t, \Phi^l\right) \equiv h_i^t$$
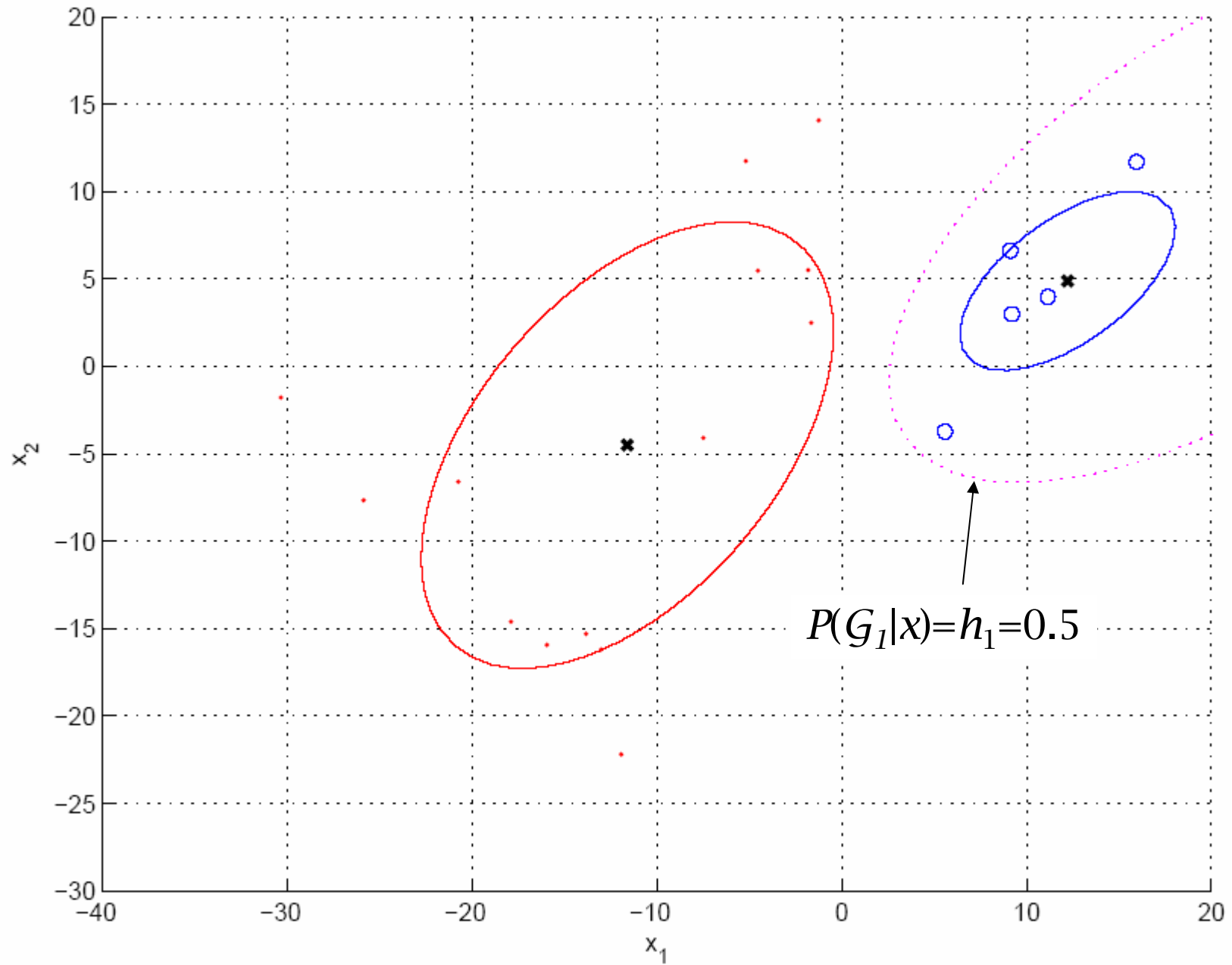
- M-step:

$$P(\mathcal{G}_i) = \frac{\sum_t h_i^t}{N} \qquad \boldsymbol{m}_i^{l+1} = \frac{\sum_t h_i^t \boldsymbol{x}^t}{\sum_t h_i^t}$$

*Use estimated labels in place of unknown labels*

$$\boldsymbol{S}_i^{l+1} = \frac{\sum_t h_i^t \left(\boldsymbol{x}^t - \boldsymbol{m}_i^{l+1}\right)\left(\boldsymbol{x}^t - \boldsymbol{m}_i^{l+1}\right)^T}{\sum_t h_i^t}$$

EM solution

$P(\mathcal{G}_1|x)=h_1=0.5$

*Lecture Notes for E ALPAYDIN 2004 Introduction to Machine Learning © The MIT Press (V1.1)*

# *Mixtures of Latent Variable Models*

- Regularize clusters
1. Assume shared/diagonal covariance matrices
2. Use PCA/FA to decrease dimensionality: Mixtures of PCA/FA

$$p(\boldsymbol{x}_t \mid \mathcal{G}_i) = \mathcal{N}\left(\boldsymbol{m}_i, \mathbf{V}_i \mathbf{V}_i^T + \psi_i\right)$$

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

# *After Clustering*

- 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*

- 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*

- 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:

$$p(\boldsymbol{x} \mid C_i) = \sum_{j=1}^{k_i} p(\boldsymbol{x} \mid \mathcal{G}_{ij}) P(\mathcal{G}_{ij})$$

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

# *Hierarchical Clustering*

- Cluster based on similarities/distances

- Distance measure between instances $\boldsymbol{x}^r$ and $\boldsymbol{x}^s$
  Minkowski ($L_p$) (Euclidean for $p = 2$)

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

City-block distance

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

# *Agglomerative Clustering*

- Start with *N* groups each with one instance and merge two closest groups at each iteration
- Distance between two groups $G_i$ and $G_j$:
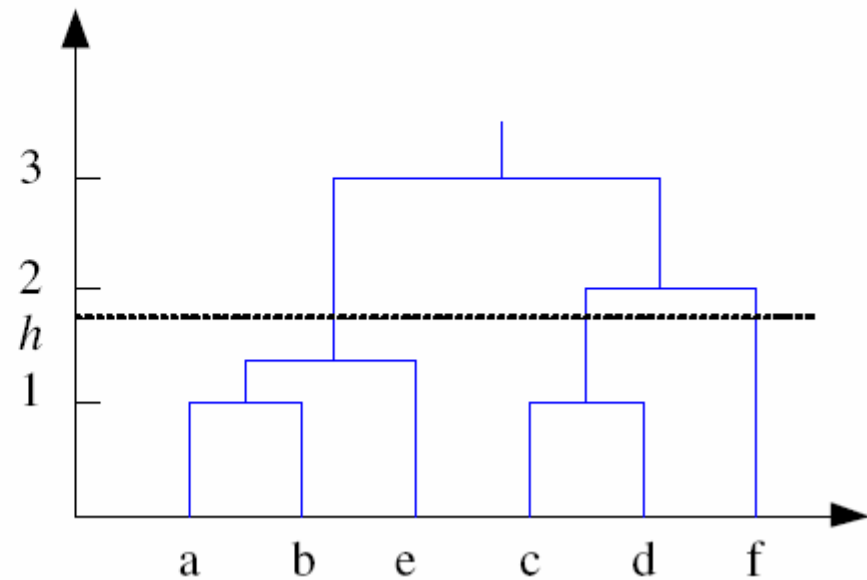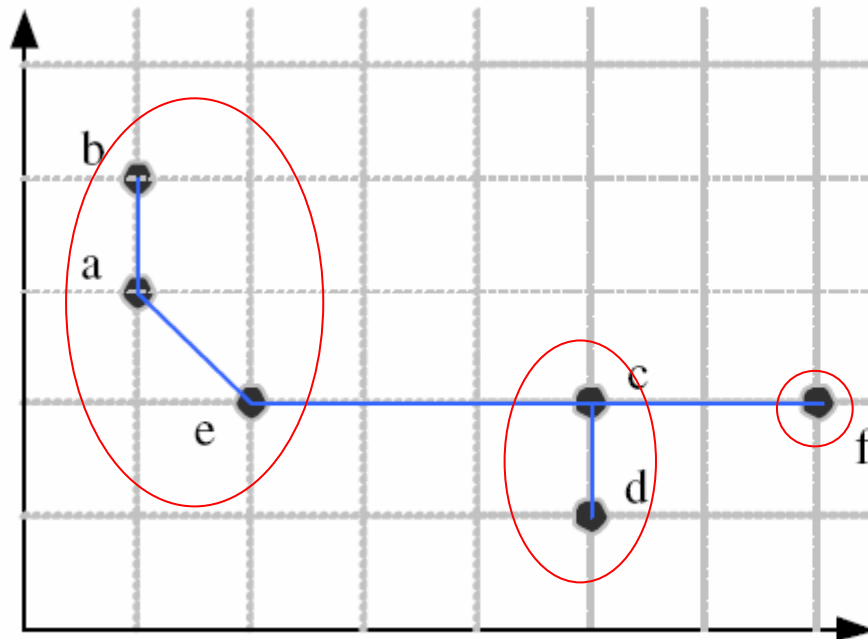    - Single-link:
    
    $$d(G_i, G_j) = \min_{x^r \in G_i, x^s \in G_j} d(x^r, x^s)$$
    
    - Complete-link:
    
    $$d(G_i, G_j) = \max_{x^r \in G_i, x^s \in G_j} d(x^r, x^s)$$
    
    - Average-link, centroid

# *Example: Single-Link Clustering*



*Dendrogram*

# *Choosing k*

- 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)
- Manual check for meaning