

Clustering algorithms

Konstantinos Koutroumbas

Unit 7

- CFO clustering algorithms: Discussion
- Hierarchical clustering algorithms:
 - the agglomerative case (based on matrix theory)

CFO clustering algorithms: Final remarks (1)

Relating hard, fuzzy and probabilistic clustering

(point representatives, squared Euclidean distance)

A. Generalized Hard Algorithmic Scheme (GHAS) – *k*-means algorithm

$$\text{minimize}_{U, \Theta} J(U, \Theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij} \|\mathbf{x}_i - \boldsymbol{\theta}_j\|^2$$

subject to **(a)** $u_{ij} \in \{0,1\}$, $i = 1, \dots, N, j = 1, \dots, m$, and **(b)** $\sum_{j=1}^m u_{ij} = 1, i = 1, \dots, N$.

The Isodata or *k*-Means or *c*-Means algorithm

- Choose arbitrary **initial estimates** $\boldsymbol{\theta}_j(0)$ for the $\boldsymbol{\theta}_j$'s, $j=1, \dots, m$.

- $t = 0$

- **Repeat**

- For $i = 1$ to N *% Determination of the partition*

- o For $j = 1$ to m

$$u_{ij}(t) = \begin{cases} 1, & \text{if } \|\mathbf{x}_i - \boldsymbol{\theta}_j(t)\|^2 = \min_{q=1, \dots, m} \|\mathbf{x}_i - \boldsymbol{\theta}_q(t)\|^2 \\ 0, & \text{otherwise} \end{cases}$$

- o End {For- j }

- End {For- i }

- $t = t + 1$

- For $j = 1$ to m *% Parameter updating*

- o Set

$$\boldsymbol{\theta}_j(t) = \frac{\sum_{i=1}^N u_{ij}(t-1) \mathbf{x}_i}{\sum_{i=1}^N u_{ij}(t-1)}, j = 1, \dots, m$$

- End {For- j }

- **Until no change** in $\boldsymbol{\theta}_j$'s **occurs** between **two successive iterations**

CFO clustering algorithms: Final remarks (1)

Relating hard, fuzzy and probabilistic clustering

(point representatives, squared Euclidean distance)

B. Generalized Fuzzy Algorithmic Scheme (GFAS) – Fuzzy c-means algorithm

$$\text{minimize}_{U, \Theta} J(U, \Theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij}^q ||\mathbf{x}_i - \boldsymbol{\theta}_j||^2$$

subject to **(a)** $u_{ij} \in (0,1)$, $i = 1, \dots, N, j = 1, \dots, m$, and **(b)** $\sum_{j=1}^m u_{ij} = 1, i = 1, \dots, N$.

- Choose $\boldsymbol{\theta}_j(0)$ as initial estimates for $\boldsymbol{\theta}_j, j=1, \dots, m$.
- $t = 0$
- Repeat
 - For $i = 1$ to N % Determination of u'_{ij} s
 - o For $j = 1$ to m

$$u_{ij}(t) = \frac{1}{\sum_{k=1}^m \left(\frac{d(\mathbf{x}_i, \boldsymbol{\theta}_j(t))}{d(\mathbf{x}_i, \boldsymbol{\theta}_k(t))} \right)^{\frac{1}{q-1}}}$$

- o End {For-j}
 - End {For-i}
 - $t = t + 1$
 - For $j = 1$ to m % Parameter updating
 - o Set

$$\boldsymbol{\theta}_j(t) = \frac{\sum_{i=1}^N u_{ij}^q(t-1) \mathbf{x}_i}{\sum_{i=1}^N u_{ij}^q(t-1)}, j = 1, \dots, m$$

- End {For-j}
- Until a termination criterion is met.

CFO clustering algorithms: Final remarks (1)

Relating hard, fuzzy and probabilistic clustering

(point representatives, squared Euclidean distance)

C. Generalized Probabilistic Algorithmic Scheme (GPrAS) – the normal pdfs case

$$\text{minimize}_{\Theta, P} J(\Theta, P) = - \sum_{i=1}^N \sum_{j=1}^m P(j|\mathbf{x}_i) \ln(p(\mathbf{x}_i|j; \boldsymbol{\theta}_j) P_j)$$

It is **(a)** $P(j|\mathbf{x}_i) \in (0,1)$, $i = 1, \dots, N, j = 1, \dots, m$, and **(b)** $\sum_{j=1}^m P(j|\mathbf{x}_i) = 1$, $i = 1, \dots, N$.

- Choose $\boldsymbol{\mu}_j(0)$, $\Sigma_j(0)$, $P_j(0)$ as **initial estimates** for $\boldsymbol{\mu}_j, \Sigma_j, P_j$, resp., $j = 1, \dots, m$
- $t = 0$
- **Repeat**
 - For $i = 1$ to N % *Expectation step*
 - o For $j = 1$ to m

$$P(j|\mathbf{x}_i; \boldsymbol{\theta}^{(t)}, P^{(t)}) = \frac{p(\mathbf{x}_i|j; \boldsymbol{\theta}_j^{(t)}) P_j^{(t)}}{\sum_{q=1}^m p(\mathbf{x}_i|q; \boldsymbol{\theta}_q^{(t)}) P_q^{(t)}} \equiv \gamma_{ji}^{(t)}$$

- o End {For- j }
- End {For- i }
- $t = t + 1$
- For $j = 1$ to m % *Parameter updating – Maximization step*
 - o Set

$$\boldsymbol{\mu}_j^{(t)} = \frac{\sum_{i=1}^N \gamma_{ji}^{(t-1)} \mathbf{x}_i}{\sum_{i=1}^N \gamma_{ji}^{(t-1)}}, \quad \Sigma_j^{(t)} = \frac{\sum_{i=1}^N \gamma_{ji}^{(t-1)} (\mathbf{x}_i - \boldsymbol{\mu}_j) (\mathbf{x}_i - \boldsymbol{\mu}_j)^T}{\sum_{i=1}^N \gamma_{ji}^{(t-1)}} \quad j = 1, \dots, m$$

$$P_j^{(t)} = \frac{1}{N} \sum_{i=1}^N \gamma_{ji}^{(t-1)}, \quad j = 1, \dots, m$$

- End {For- j }
- **Until** a **termination criterion** is met.

CFO clustering algorithms: Final remarks (1)

Relating hard, fuzzy and probabilistic clustering

(point representatives, squared Euclidean distance)

Consider the **GPrAS cost function**

$$J(\Theta, P) = - \sum_{i=1}^N \sum_{j=1}^m P(j|\mathbf{x}_i) \ln(p(\mathbf{x}_i|j; \boldsymbol{\theta}_j) P_j)$$

with

$$p(\mathbf{x}_i|j; \boldsymbol{\theta}_j) = \frac{1}{(2\pi)^{\frac{l}{2}} |\Sigma_j|^{\frac{1}{2}}} \exp\left(-\frac{(\mathbf{x}_i - \boldsymbol{\mu}_j)^T \Sigma_j^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j)}{2}\right)$$

$$\boldsymbol{\theta}_j = \{\boldsymbol{\mu}_j, \Sigma_j\}$$

It is $J(\boldsymbol{\theta}, P) = - \sum_{i=1}^N \sum_{j=1}^m P(j|\mathbf{x}_i) \ln\left(\frac{1}{(2\pi)^{\frac{l}{2}} |\Sigma_j|^{\frac{1}{2}}} \exp\left(-\frac{(\mathbf{x}_i - \boldsymbol{\mu}_j)^T \Sigma_j^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j)}{2}\right) P_j\right) =$

Term **A** $- \sum_{i=1}^N \sum_{j=1}^m P(j|\mathbf{x}_i) \ln\left(\frac{1}{(2\pi)^{\frac{l}{2}} |\Sigma_j|^{\frac{1}{2}}}\right)$

Term **B** $+ \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^m P(j|\mathbf{x}_i) (\mathbf{x}_i - \boldsymbol{\mu}_j)^T \Sigma_j^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j)$

Term **C** $- \sum_{i=1}^N \sum_{j=1}^m P(j|\mathbf{x}_i) \ln P_j$

CFO clustering algorithms: Final remarks (1)

Relating hard, fuzzy and probabilistic clustering

(point representatives, squared Euclidean distance)

Assumption 1: $\Sigma_j = \Sigma = \text{constant}$, $j = 1, \dots, m$. Then

$$\begin{aligned} \text{Term } \mathbf{A} &= - \sum_{i=1}^N \sum_{j=1}^m P(j|\mathbf{x}_i) \ln \left(\frac{1}{(2\pi)^{\frac{l}{2}} |\Sigma|^{\frac{1}{2}}} \right) \\ &= - \ln \left(\frac{1}{(2\pi)^{\frac{l}{2}} |\Sigma|^{\frac{1}{2}}} \right) \sum_{i=1}^N \sum_{j=1}^m P(j|\mathbf{x}_i) = - \ln \left(\frac{1}{(2\pi)^{\frac{l}{2}} |\Sigma|^{\frac{1}{2}}} \right) \sum_{i=1}^N 1 \\ &= -N \ln \left(\frac{1}{(2\pi)^{\frac{l}{2}} |\Sigma|^{\frac{1}{2}}} \right) = \text{constant} \end{aligned}$$

Assumption 2: $P_j = \frac{1}{m}$, $j = 1, \dots, m$. Then

Term **C**

$$= - \sum_{i=1}^N \sum_{j=1}^m P(j|\mathbf{x}_i) \ln \frac{1}{m} = - \ln \frac{1}{m} \sum_{i=1}^N \sum_{j=1}^m P(j|\mathbf{x}_i) = -N \ln \frac{1}{m} = \text{constant}$$

CFO clustering algorithms: Final remarks (1)

Relating hard, fuzzy and probabilistic clustering

(point representatives, squared Euclidean distance)

Based on the previous two results, it follows that

$$\text{minimize} \left(- \sum_{i=1}^N \sum_{j=1}^m P(j|\mathbf{x}_i) \ln(p(\mathbf{x}_i|j; \boldsymbol{\theta}_j) P_j) \right)$$



$$\text{minimize} \left(\sum_{i=1}^N \sum_{j=1}^m P(j|\mathbf{x}_i) (\mathbf{x}_i - \boldsymbol{\mu}_j)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_j) \right)$$

$$\boldsymbol{\Sigma}_j = \boldsymbol{\Sigma}$$

Assumption 3(a): Approximate $P(j|\mathbf{x}_i)$ as

$$P(j|\mathbf{x}_i) = \begin{cases} 1, & P(j|\mathbf{x}_i) = \max_{s=1, \dots, m} P(s|\mathbf{x}_i) \quad (\equiv u_{ij}) \\ 0, & \text{otherwise} \end{cases}$$

In this case, $GPrAS \Leftrightarrow k - \text{means}$ (for $\boldsymbol{\Sigma} = \sigma^2 I$)

Assumption 3(b): Approximate $P(j|\mathbf{x}_i)$ as

$$P(j|\mathbf{x}_i) = \frac{1}{\sum_{k=1}^m \left(\frac{d(\mathbf{x}_i, \boldsymbol{\theta}_j(t))}{d(\mathbf{x}_i, \boldsymbol{\theta}_k(t))} \right)^{\frac{1}{q-1}}}$$

WARNING: Valid ONLY from a mathematical formulation point of view. NOT from a conceptual point of view.

In this case, $GPrAS \Leftrightarrow \text{fuzzy } c - \text{means}$ (for $\boldsymbol{\Sigma} = \sigma^2 I$)

CFO clustering algorithms: Final remarks (1)

Relating hard, fuzzy and probabilistic clustering

(point representatives, squared Euclidean distance)

Remarks:

The **hard**, **fuzzy** and **probabilistic CFO** clustering algorithms (with point representatives and squared Euclidean distance) :

- are **partition algorithms**.
- they **share** the “**sum-to-one**” constraint.
- they can be related to each other (due to the “sum-to-one” constraint).

The **possibilistic** CFO clustering algorithms (point representatives and squared Euclidean distance) :

- are **mode seeking algorithms**
- no “**sum-to-one**” constraint is associated with them
- they can not be related to the hard, fuzzy and probabilistic CFO clustering algorithms (due to the absence of the sum-to-one constraint).

CFO clustering algorithms: Final remarks (2)

The role of q in the fuzzy clustering

Consider the minimization problem for fuzzy clustering

$$\text{minimize}_{U, \Theta} J(U, \Theta) = \sum_{i=1}^N \sum_{j=1}^m u_{ij}^q d_{ij}$$

$$d_{ij} = d(x_i, \theta_j)$$

subject to **(a)** $u_{ij} \in (0,1)$, $i = 1, \dots, N, j = 1, \dots, m$, and **(b)** $\sum_{j=1}^m u_{ij} = 1, i = 1, \dots, N$.

Expanding $J(U, \Theta)$, we have

$$J(U, \Theta) = \begin{array}{cccc} u_{11}^q d_{11} + & u_{12}^q d_{12} + & \dots & u_{1m}^q d_{1m} \\ u_{21}^q d_{21} + & u_{22}^q d_{22} + & \dots & u_{2m}^q d_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ u_{N1}^q d_{N1} + & u_{N2}^q d_{N2} + & \dots & u_{Nm}^q d_{Nm} \end{array}$$

Assumption: d_{ij} 's are fixed.

Then, due to the sum-to-one constraint, $J(U, \Theta)$ is **minimized** if each of the summation in the rows of the above expansion is minimized.

Let s_i : $d_{is_i} = \min_{j=1, \dots, m} d_{ij}, i = 1, \dots, N$

Then,

$$u_{i1}^q d_{i1} + \dots + u_{im}^q d_{im} \geq \left(\sum_{j=1}^m u_{ij}^q \right) d_{is_i}$$

CFO clustering algorithms: Final remarks (2)

The role of q in the fuzzy clustering

$$A_i = u_{i1}^q d_{i1} + \dots + u_{im}^q d_{im} \geq \left(\sum_{j=1}^m u_{ij}^q \right) d_{is_i}$$

For $q = 1$, it is $\sum_{j=1}^m u_{ij} = 1$. Thus

$$A_i = u_{i1} d_{i1} + \dots + u_{im} d_{im} \geq d_{is_i}$$

Clearly, the **equality holds** for $u_{is_i} = 1$ and $u_{ij} = 0$, for $j = 1, \dots, m, j \neq s_i$

In other words the minimum possible value of A_i is achieved for the hard cluster solution. Thus, **no fuzzy clustering** (where more than one u_{ij} 's are positive) **minimizes** the A_i .

For $q > 1$, in the **hard clustering** case, the minimum possible value of A_i is still d_{is_i} .

For $q > 1$, in the **fuzzy clustering** case, it is $\sum_{j=1}^m u_{ij}^q < 1$. Thus

$$d_{is_i} > \left(\sum_{j=1}^m u_{ij}^q \right) d_{is_i}$$

Thus, in this cases, there are choices for u_{ij} 's with more than one of them being positive (fuzzy case) that achieve lower value for A_i than the best hard clustering.

The **larger** the value of q , the **more fuzzy clusterings achieve** for A_i value $< d_{is_i}$. ¹⁰

CFO clustering algorithms: Final remarks (2)

The role of q in the fuzzy clustering

Example: $X = \{x_1, x_2, x_3, x_4\}$

$$x_1 = [0,0]^T, x_2 = [2,0]^T, x_3 = [0,3]^T, x_4 = [2,3]^T$$

$$\theta_1 = [1,0]^T, \theta_2 = [1,3]^T \text{ (fixed)}$$

$$q = 1 \text{ (hard case): Best solution } U_{hard} = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix}, J_{hard} = 4$$

$q = 2$ (fuzzy case): Focus on x_1 :

Question: Is it possible to have $u_{11}^2 \cdot 1 + u_{12}^2 \cdot \sqrt{10} < 1$? **(A)**

Since $u_{12} = 1 - u_{11}$, **(A)** becomes

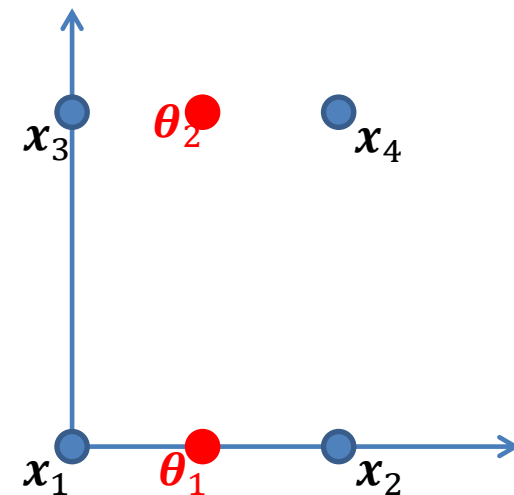
$$u_{11}^2 \cdot 1 + (1 - u_{11})^2 \cdot \sqrt{10} < 1 \Leftrightarrow$$

$$(\sqrt{10} + 1)u_{11}^2 - 2\sqrt{10}u_{11} + \sqrt{10} - 1 < 0 \Leftrightarrow$$

$$u_{11} \in (0.52, 1) \Rightarrow u_{12} \in (0, 0.48)$$

For example, if $u_{11} = 0.7$ ($u_{12} = 0.3$), it is

$$u_{11}^2 \cdot 1 + u_{12}^2 \cdot \sqrt{10} = 0.7^2 \cdot 1 + 0.3^2 \cdot \sqrt{10} = 0.77 < 1$$



$d(x_i, \theta_j)$	$\theta_1 = (1,0)$	$\theta_2 = (1,3)$
$x_1 = (0,0)$	$d_{11} = 1$	$d_{12} = \sqrt{10}$
$x_2 = (2,0)$	$d_{21} = 1$	$d_{22} = \sqrt{10}$
$x_3 = (0,3)$	$d_{31} = \sqrt{10}$	$d_{32} = 1$
$x_4 = (2,3)$	$d_{41} = \sqrt{10}$	$d_{42} = 1$

CFO clustering algorithms: Final remarks (3)

The role of q in the possibilistic clustering

Consider the minimization problem for possibilistic clustering

$$\text{minimize}_{U, \Theta} J(\mathbf{u}_j, \boldsymbol{\theta}_j) = \sum_{i=1}^N u_{ij}^q d_{ij} + \eta_j \sum_{i=1}^N (1 - u_{ij})^q$$

subject to $u_{ij} \in (0,1)$, $i = 1, \dots, N, j = 1, \dots, m$.

For $q = 1$, $J(\mathbf{u}_j, \boldsymbol{\theta}_j)$ is written as

$$J(\mathbf{u}_j, \boldsymbol{\theta}_j) = \sum_{i=1}^N [u_{ij}(d_{ij} - \eta_j) + \eta_j]$$

Thus, minimizing $J(\mathbf{u}_j, \boldsymbol{\theta}_j)$ is equivalent to minimizing

$$\sum_{i=1}^N u_{ij}(d_{ij} - \eta_j)$$

For fixed $\boldsymbol{\theta}_j (\Rightarrow \text{fixed } d(x_i, \boldsymbol{\theta}_j) \equiv d_{ij})$, the latter achieves its **minimum** (negative) value by selecting $u_{ij} = 1$, for $d_{ij} < \eta_j$ and $u_{ij} = 0$, for $d_{ij} > \eta_j$.

However, in the above situation, all points having distance less than η_j from $\boldsymbol{\theta}_j$, they all have the same weight in the determination of $\boldsymbol{\theta}_j$, while all the other points have no influence in the determination of $\boldsymbol{\theta}_j$.

CFO clustering algorithms: Final remarks (3)

The role of q in the possibilistic clustering

Consider the minimization problem for possibilistic clustering

$$\text{minimize}_{U, \Theta} J(\mathbf{u}_j, \boldsymbol{\theta}_j) = \sum_{i=1}^N u_{ij}^q d_{ij} + \eta_j \sum_{i=1}^N (1 - u_{ij})^q$$

subject to $u_{ij} \in (0,1)$, $i = 1, \dots, N, j = 1, \dots, m$.

➤ For $q > 1$, (for fixed $\boldsymbol{\theta}_j (\Rightarrow \text{fixed } d(x_i, \boldsymbol{\theta}_j) \equiv d_{ij})$) it is

$$u_{ij} = \frac{1}{1 + \left(\frac{d_{ij}}{\eta_j}\right)^{\frac{1}{q-1}}}$$

Thus, points for which $d_{ij} > \eta_j$ have $u_{ij} > 0$.

➤ Furthermore, as $q \rightarrow \infty$, (for fixed $\boldsymbol{\theta}_j (\Rightarrow \text{fixed } d(x_i, \boldsymbol{\theta}_j) \equiv d_{ij})$) it is

$$u_{ij} \rightarrow \frac{1}{2}$$

Thus, **all points** have the **same degree of compatibility** with **all clusters**.

CFO clustering algorithms: Final remarks (4)

The role of q in the parameters updating in fuzzy and possibilistic clustering

Let $0 < u_1 < u_2 < 1$.

We define $\Delta u = u_2 - u_1$ and $\Delta u^q = u_2^q - u_1^q$ ($q > 1$).

For $q = 2$, it is $\Delta u^2 = u_2^2 - u_1^2 = \Delta u(u_2 + u_1)$.

If $u_2 + u_1 > 1$ (the respective points are “close” to the representative under consideration), then $\Delta u^2 > \Delta u$

Therefore, the function $f(u) = u^q$, **enhances** the **difference** between two values of u .

If $u_2 + u_1 < 1$ (the respective points are “away” from the representative under consideration), then $\Delta u^2 < \Delta u$

Therefore, the function $f(u) = u^q$, **diminishes** the **difference** between two values of u .

CFO clustering algorithms: Final remarks (4)

The role of q in the parameters updating in fuzzy and possibilistic clustering

Consider the updating equation for the point representative case and the squared Euclidean distance case (**fuzzy** and **1st possibilistic** clust. algorithms)

$$\theta_j(t) = \frac{\sum_{i=1}^N u_{ij}^q(t-1) \mathbf{x}_i}{\sum_{i=1}^N u_{ij}^q(t-1)}, j = 1, \dots, m$$

For $q > 1$, and since $u_{ij} \in (0,1)$, the previous observation indicates that the \mathbf{x}_i 's with **high** (**low**) u_{ij} , will have **more** (**much less**) significant contribution to the estimation of $\theta_j(t)$, compared with the $q = 1$ case.

Example: Let $\mathbf{x}_1 = [0, 0]^T$ and $\mathbf{x}_2 = [10, 10]^T$, and $u_{1j} = 0.1$, $u_{2j} = 0.9$. Then

$$\theta_j = \frac{u_{1j} \mathbf{x}_1 + u_{2j} \mathbf{x}_2}{u_{1j} + u_{2j}} = \begin{bmatrix} 9 \\ 9 \end{bmatrix} \quad (q = 1)$$

and

$$\theta_j = \frac{u_{1j}^q \mathbf{x}_1 + u_{2j}^q \mathbf{x}_2}{u_{1j}^q + u_{2j}^q} = \begin{bmatrix} 9.9 \\ 9.9 \end{bmatrix} \quad (q = 2)$$

Optimization theory – Basic concepts

Let $J(\mathbf{w})$ be a continuous function of \mathbf{w} .

Problem (P1): Determine the **position \mathbf{w}^*** where the function $J(\mathbf{w})$ achieves its **minimum** value.

A simple method for solving **(P1)** is that of **gradient descent**.

-Initialize $\mathbf{w} = \mathbf{w}(0)$

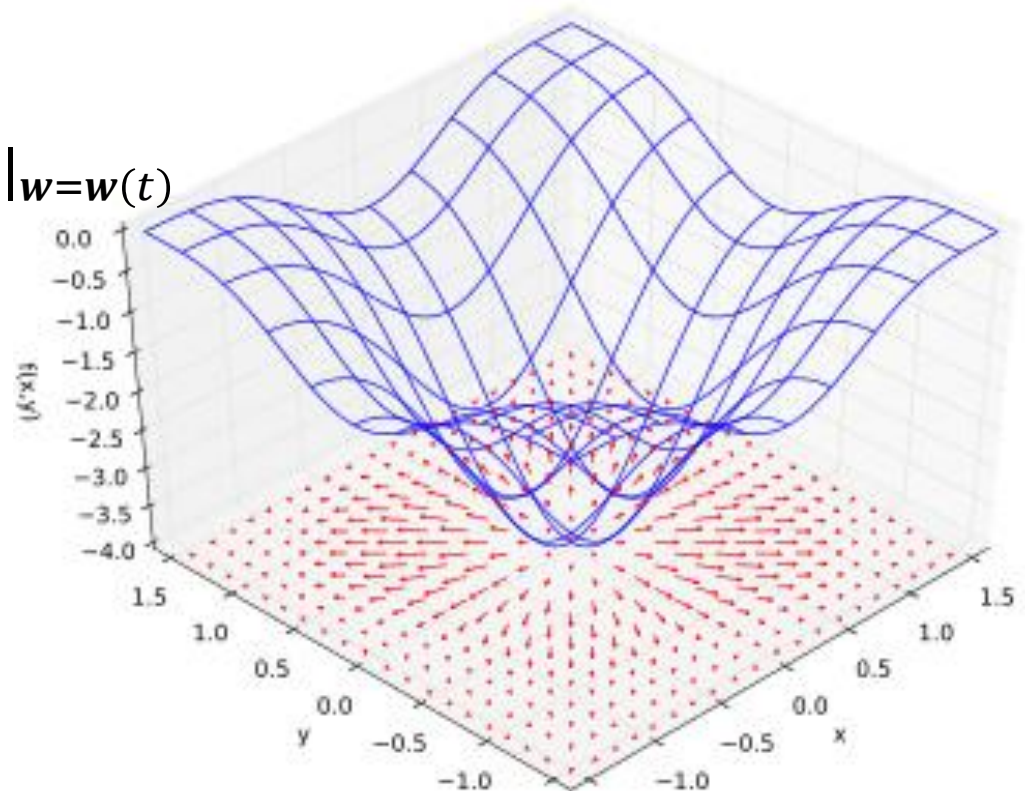
- $t = 0$

-Repeat

$$- \mathbf{w}(t + 1) = \mathbf{w}(t) - \mu \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} \Big|_{\mathbf{w}=\mathbf{w}(t)}$$

- $t = t + 1$

-Until convergence



Optimization theory – Basic concepts

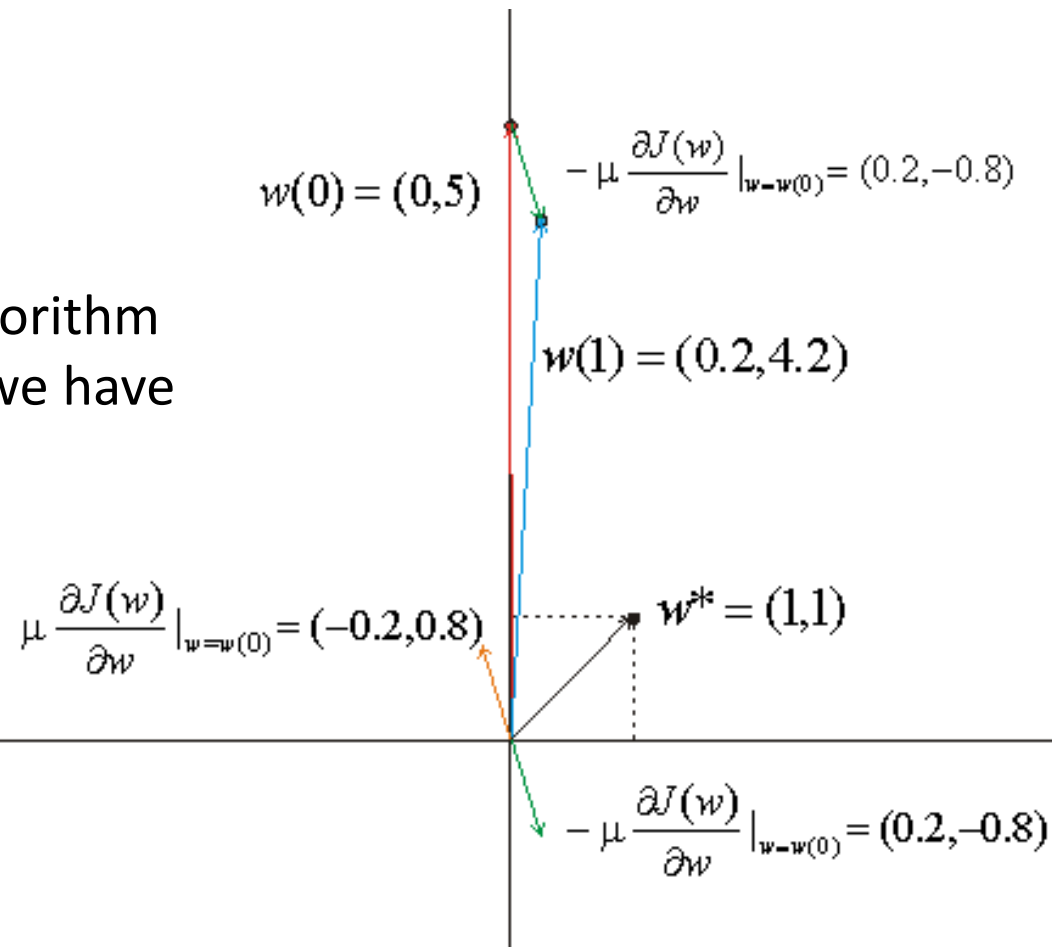
-An example: Let $\mathbf{w} = [w_1, w_2]^T$ and $J(\mathbf{w}) = (w_1 - 1)^2 + (w_2 - 1)^2$. Clearly, the minimum value of $J(\mathbf{w})$ is met at $\mathbf{w}^* = [1, 1]^T$.

-It is
$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \begin{bmatrix} 2w_1 - 2 \\ 2w_2 - 2 \end{bmatrix}$$

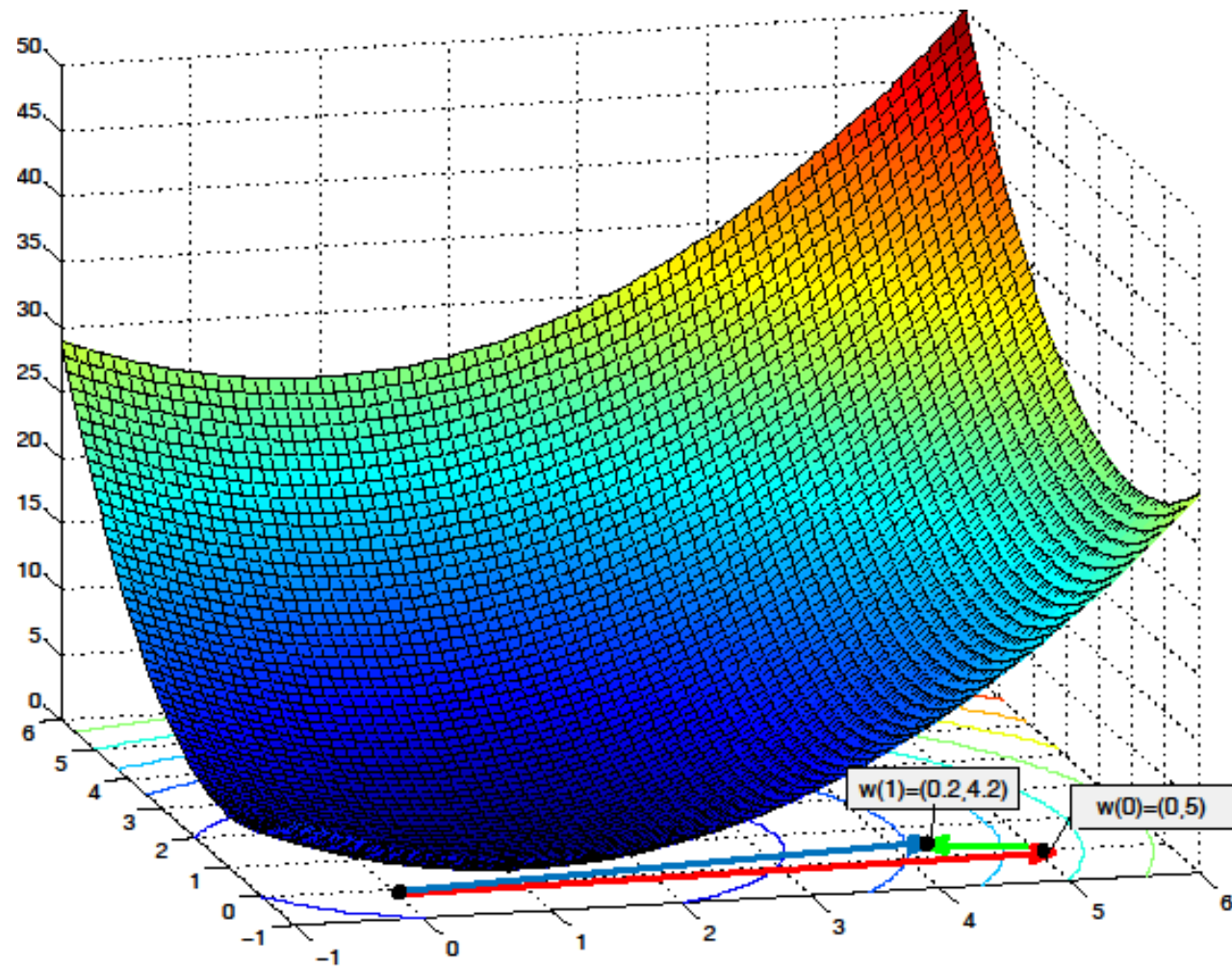
-Applying the gradient descent algorithm for $\mathbf{w}(0) = [0, 5]^T$, and $\mu = 0.1$, we have

$$\mathbf{w}(1) = \begin{bmatrix} 0 \\ 5 \end{bmatrix} - 0.1 \begin{bmatrix} -2 \\ 8 \end{bmatrix} = \begin{bmatrix} 0.2 \\ 4.2 \end{bmatrix}$$

-Thus, $\mathbf{w}(1)$ comes closer to \mathbf{w}^* .



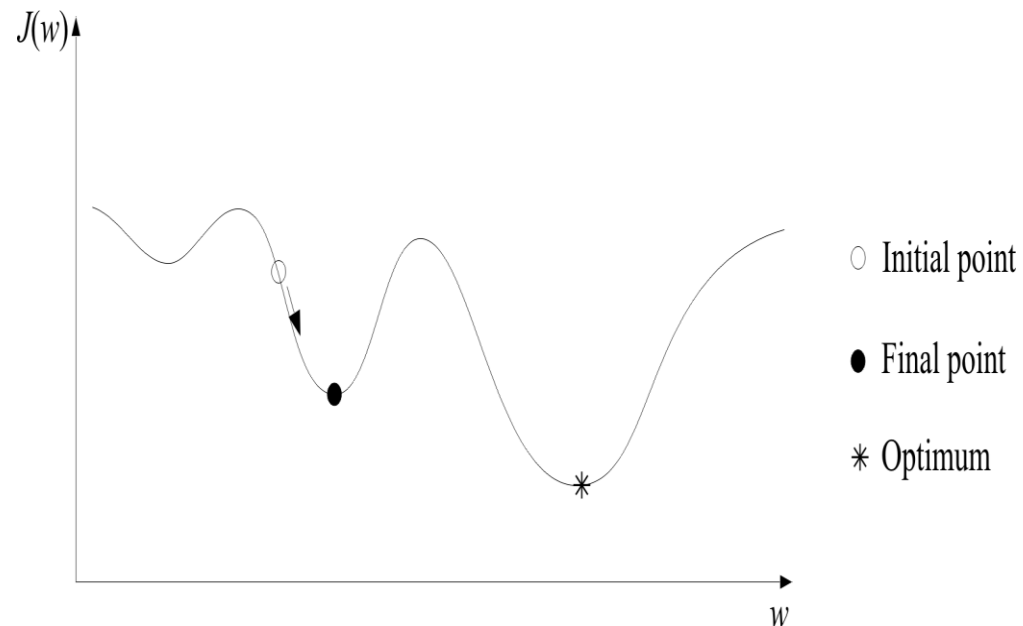
Optimization theory – Basic concepts



Optimization theory – Basic concepts

Remarks for gradient descent:

- The value of μ should be chosen **not too large**, in order to avoid oscillations around the minimum and **not too small** in order to avoid unnecessary delays in the convergence
- If $J(\mathbf{w})$ has **more than one local minima**, the gradient descent will converge (in general) to the one that is closest to $\mathbf{w}(0)$.
- If the algorithm is trapped to a **local minimum** that correspond to a poor solution, the only way to **escape** from it is to **re-initialize** the algorithm from another initial position.
- It can be proved that, under certain conditions, the algorithm **converges asymptotically** to a **local minimum** of $J(\mathbf{w})$.



Optimization theory – Basic concepts

Let $J(\mathbf{w})$ be a continuous function of \mathbf{w} .

Problem (P2): Determine the **position \mathbf{w}^*** where the function $J(\mathbf{w})$ achieves its **minimum** value, under the constraint that \mathbf{w} satisfies some **equality constraints**.

For **linear equality constraints**, the problem is stated as follows

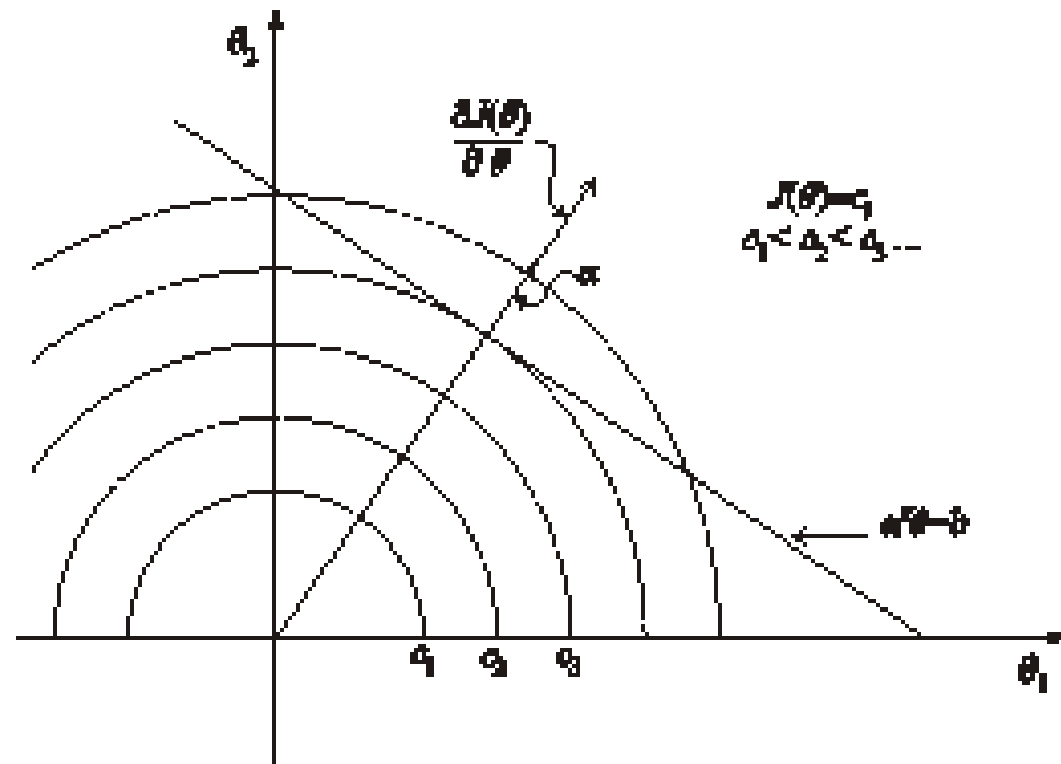
- Minimize $J(\mathbf{w})$
- Subject to the constraints $A\mathbf{w} = \mathbf{b}$, where A an $m \times l$ matrix and \mathbf{b} an m -dim. Vector.

Solution: Lagrange multipliers

Minimize

- $L(\mathbf{w}) = J(\mathbf{w}) + \boldsymbol{\lambda}^T(A\mathbf{w} - \mathbf{b})$

- $\boldsymbol{\lambda}$ is an m -dim vector that is estimated through the constraints $A\mathbf{w} = \mathbf{b}$



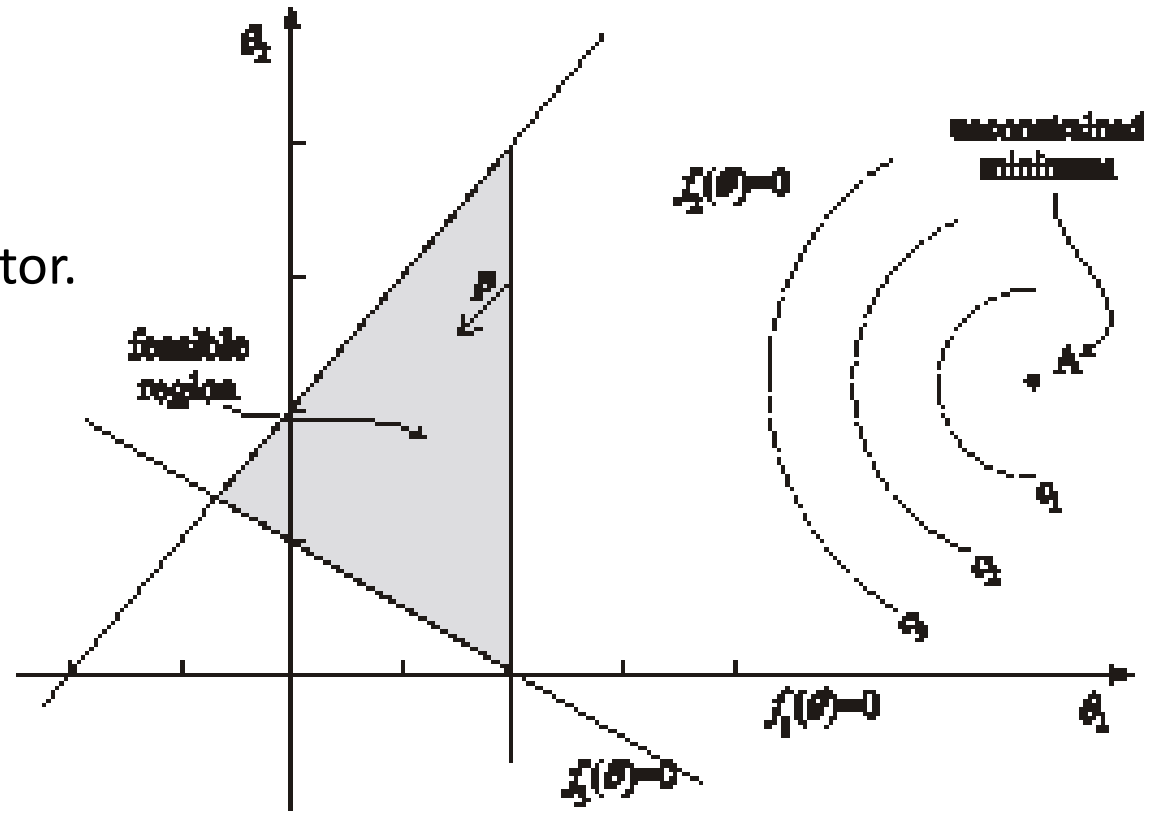
Optimization theory – Basic concepts

Let $J(\mathbf{w})$ be a continuous function of \mathbf{w} .

Problem (P3): Determine the **position \mathbf{w}^*** where the function $J(\mathbf{w})$ achieves its **minimum** value, under the constraint that \mathbf{w} satisfies some **inequality constraints**.

For **linear inequality constraints**, the problem is stated as follows

- Minimize $J(\mathbf{w})$
- Subject to the constraints $A\mathbf{w} \geq \mathbf{b}$, where A an $m \times l$ matrix and \mathbf{b} an m -dim. Vector.



Hierarchical Clustering Algorithms

- ✓ They produce a **hierarchy** of (**hard**) clusterings instead of a **single** clustering.
- ✓ They find applications in:
 - Social sciences
 - Biological taxonomy
 - Modern biology
 - Medicine
 - Archaeology
 - Computer science and engineering

Hierarchical Clustering Algorithms

Let $X = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, $\mathbf{x}_i = [x_{i1}, \dots, x_{il}]^T$.

Recall that:

- In hard clustering each vector belongs **exclusively** to a single cluster.
- An **m -(hard) clustering** of X , \mathfrak{R} , is a partition of X into m sets (clusters) C_1, \dots, C_m , so that:

- $C_j \neq \emptyset, j = 1, \dots, m$
- $\bigcup_{j=1}^m C_j = X$
- $C_i \cap C_j = \emptyset, i \neq j, i, j = 1, 2, \dots, m$

By the definition: $\mathfrak{R} = \{C_j, j = 1, \dots, m\}$

Hierarchical Clustering Algorithms

➤ **Definition:** A clustering \mathcal{R}_1 consisting of k clusters is said to be **nested** in the clustering \mathcal{R}_2 consisting of r ($< k$) clusters, if **each cluster in \mathcal{R}_1 is a subset of a cluster in \mathcal{R}_2 .**

We write $\mathcal{R}_1 \angle \mathcal{R}_2$

Example: Let $\mathcal{R}_1 = \{\{\mathbf{x}_1, \mathbf{x}_3\}, \{\mathbf{x}_4\}, \{\mathbf{x}_2, \mathbf{x}_5\}\}$, $\mathcal{R}_2 = \{\{\mathbf{x}_1, \mathbf{x}_3, \mathbf{x}_4\}, \{\mathbf{x}_2, \mathbf{x}_5\}\}$,

$\mathcal{R}_3 = \{\{\mathbf{x}_1, \mathbf{x}_4\}, \{\mathbf{x}_3\}, \{\mathbf{x}_2, \mathbf{x}_5\}\}$, $\mathcal{R}_4 = \{\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_4\}, \{\mathbf{x}_3, \mathbf{x}_5\}\}$.

It is $\mathcal{R}_1 \angle \mathcal{R}_2$, **but not** $\mathcal{R}_1 \angle \mathcal{R}_3$, $\mathcal{R}_1 \angle \mathcal{R}_4$, $\mathcal{R}_1 \angle \mathcal{R}_1$.

Hierarchical Clustering Algorithms

Remarks:

- Hierarchical clustering algorithms produce a **hierarchy of nested clusterings**.
- They involve **N steps** at the most.
- At each step t , the clustering \mathfrak{R}_t is produced by \mathfrak{R}_{t-1} .

➤ Main strategies:

Agglomerative hierarchical clustering algorithms	Divisive hierarchical clustering algorithms
$\mathfrak{R}_0 = \{\{\mathbf{x}_1\}, \dots, \{\mathbf{x}_N\}\}$	$\mathfrak{R}_0 = \{\{\mathbf{x}_1, \dots, \mathbf{x}_N\}\}$
\dots	\dots
$\mathfrak{R}_{N-1} = \{\{\mathbf{x}_1, \dots, \mathbf{x}_N\}\}$	$\mathfrak{R}_{N-1} = \{\{\mathbf{x}_1\}, \dots, \{\mathbf{x}_N\}\}$
$\mathfrak{R}_0 \angle \dots \angle \mathfrak{R}_{N-1}$	$\mathfrak{R}_{N-1} \angle \dots \angle \mathfrak{R}_0$

Agglomerative Clustering Algorithms

Let $g(C_i, C_j)$ a **proximity function** between two clusters C_i and C_j of X .

Generalized Agglomerative Scheme (GAS)

➤ Initialization

- **Choose** $\mathcal{R}_0 = \{\{x_1\}, \dots, \{x_N\}\}$
- $t = 0$

➤ Repeat

- $t = t + 1$
- **Choose** (C_i, C_j) in \mathcal{R}_{t-1} such that

$$g(C_i, C_j) = \begin{cases} \min_{r,s} g(C_r, C_s), & \text{if } g \text{ is a disim. function} \\ \max_{r,s} g(C_r, C_s), & \text{if } g \text{ is a sim. function} \end{cases}$$

- Define $C_q = C_i \cup C_j$ and produce $\mathcal{R}_t = (\mathcal{R}_{t-1} - \{C_i, C_j\}) \cup \{C_q\}$

➤ Until all vectors lie in a single cluster.

Agglomerative Clustering Algorithms

Remarks:

- If two vectors come together into a single cluster at level t of the hierarchy, they will remain in the same cluster for all subsequent clusterings. As a consequence, there **is no way** to recover a “**poor**” clustering that may have occurred in an earlier level of hierarchy.
- Number of operations: $O(N^3)$

Agglomerative Clustering Algorithms

Definitions of some useful quantities:

Let $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$, with $\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{il}]^T$.

- **Pattern matrix** ($D(X)$): An $N \times l$ matrix whose i -th row is \mathbf{x}_i (transposed).
- **Proximity (similarity or dissimilarity) matrix** ($P(X)$): An $N \times N$ matrix whose (i, j) element equals the proximity $\wp(\mathbf{x}_i, \mathbf{x}_j)$ (similarity $s(\mathbf{x}_i, \mathbf{x}_j)$, dissimilarity $d(\mathbf{x}_i, \mathbf{x}_j)$).

Example 1: Let $X = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5\}$, with

$$\mathbf{x}_1 = [1, 1]^T, \mathbf{x}_2 = [2, 1]^T, \mathbf{x}_3 = [5, 4]^T, \mathbf{x}_4 = [6, 5]^T, \mathbf{x}_5 = [6.5, 6]^T$$

Pattern matrix

Euclidean distance

Tanimoto distance

$$D(X) = \begin{bmatrix} 1 & 1 \\ 2 & 1 \\ 5 & 4 \\ 6 & 5 \\ 6.5 & 6 \end{bmatrix} \quad P(X) = \begin{bmatrix} 0 & 1 & 5 & 6.4 & 7.4 \\ 1 & 0 & 4.2 & 5.7 & 6.7 \\ 5 & 4.2 & 0 & 1.4 & 2.5 \\ 6.4 & 5.7 & 1.4 & 0 & 1.1 \\ 7.4 & 6.7 & 2.5 & 1.1 & 0 \end{bmatrix} \quad P'(X) = \begin{bmatrix} 1 & 0.75 & 0.26 & 0.21 & 0.18 \\ 0.75 & 1 & 0.44 & 0.35 & 0.20 \\ 0.26 & 0.44 & 1 & 0.96 & 0.90 \\ 0.21 & 0.35 & 0.96 & 1 & 0.98 \\ 0.18 & 0.20 & 0.90 & 0.98 & 1 \end{bmatrix}$$

Agglomerative Clustering Algorithms

Definitions of some useful quantities:

➤ **Threshold dendrogram** (or **dendrogram**): It is an effective way of **representing the sequence of clusterings**, which are produced by an agglomerative algorithm.

Example 1 (cont.): If $d_{min}^{SS}(C_i, C_j)$ is employed as the distance measure **between two sets** and the **Euclidean** one as the distance measure **between two vectors**, the following series of clusterings are produced:

$$\{\{x_1\}, \{x_2\}, \{x_3\}, \{x_4\}, \{x_5\}\}$$

$$\{\{x_1, x_2\}, \{x_3\}, \{x_4\}, \{x_5\}\}$$

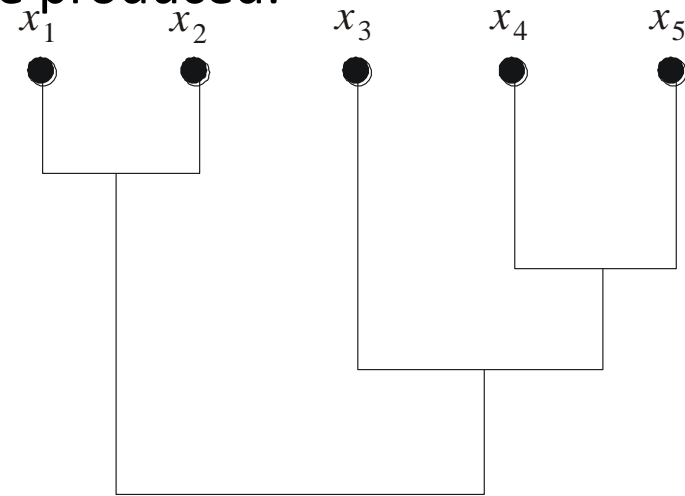
$$\{\{x_1, x_2\}, \{x_3\}, \{x_4, x_5\}\}$$

$$\{\{x_1, x_2\}, \{x_3, x_4, x_5\}\}$$

$$\{\{x_1, x_2, x_3, x_4, x_5\}\}$$

$$D(X) = \begin{bmatrix} 1 & 1 \\ 2 & 1 \\ 5 & 4 \\ 6 & 5 \\ 6.5 & 6 \end{bmatrix}$$

$$P(X) = \begin{bmatrix} 0 & 1 & 5 & 6.4 & 7.4 \\ 1 & 0 & 4.2 & 5.7 & 6.7 \\ 5 & 4.2 & 0 & 1.4 & 2.5 \\ 6.4 & 5.7 & 1.4 & 0 & 1.1 \\ 7.4 & 6.7 & 2.5 & 1.1 & 0 \end{bmatrix}$$



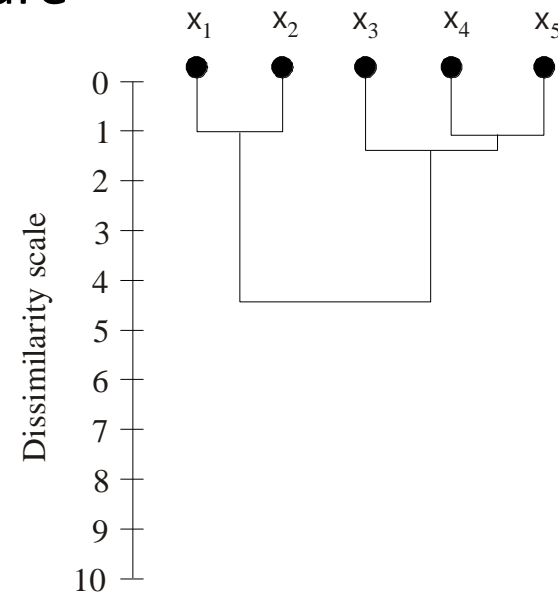
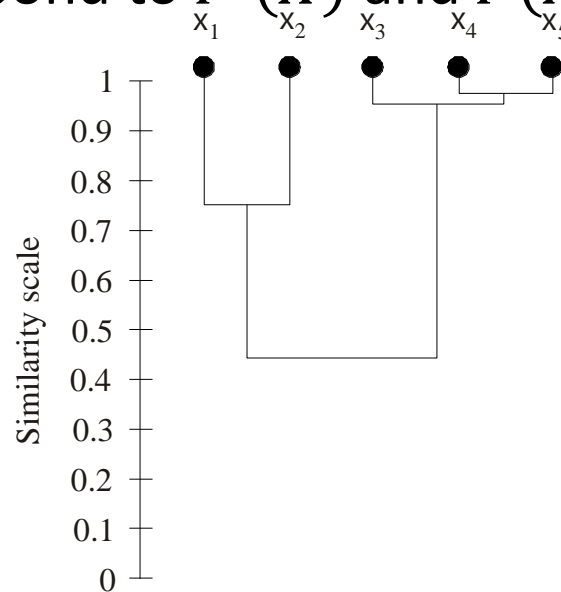
Agglomerative Clustering Algorithms

Definitions of some useful quantities:

➤ **Proximity (dissimilarity or similarity) dendrogram**: A dendrogram that takes into account the **level of proximity** (dissimilarity or similarity) where two clusters are **merged for the first time**.

Example 1 (cont.): In terms of the previous example, the proximity dendrograms that correspond to $P'(X)$ and $P(X)$ are

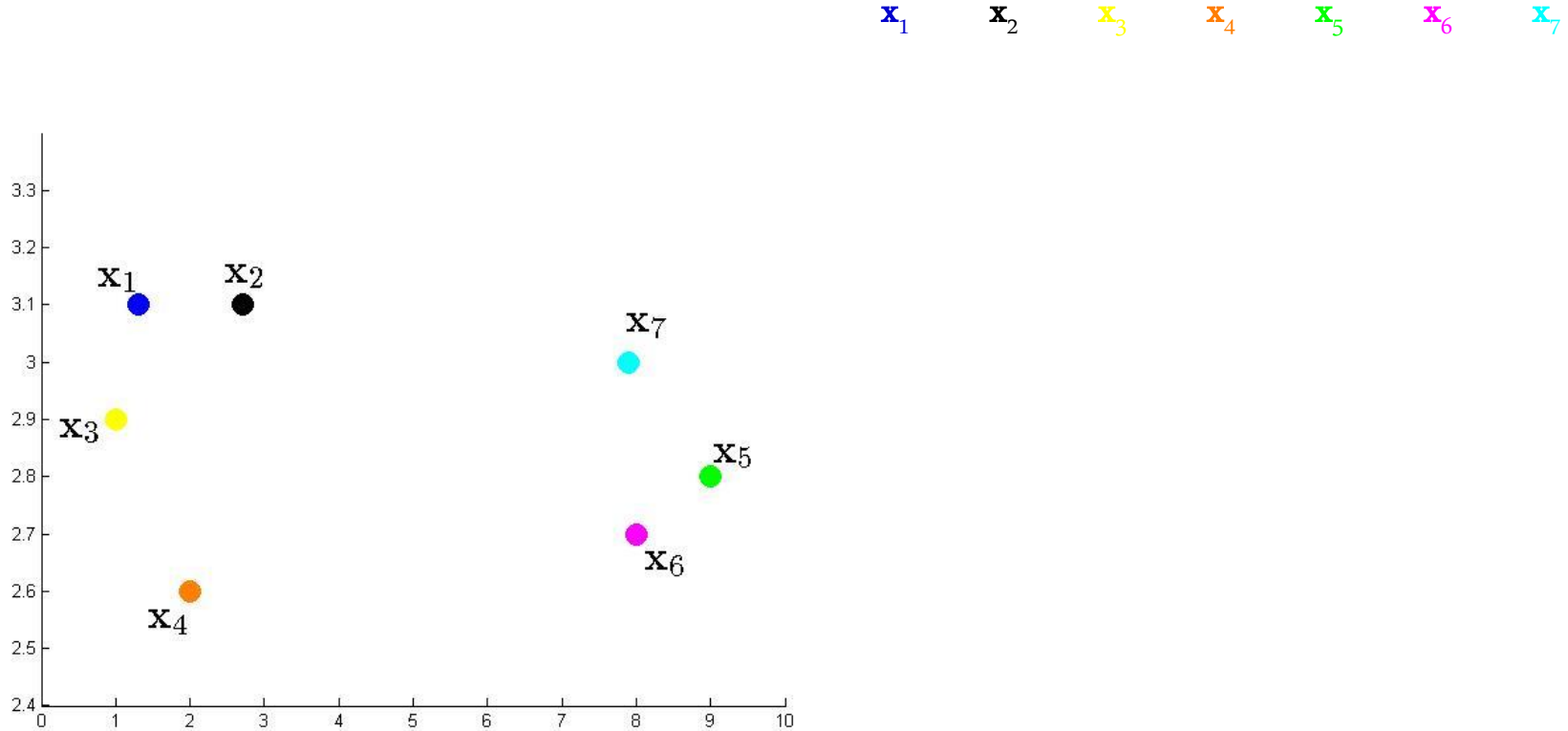
$$P(X) = \begin{bmatrix} 0 & 1 & 5 & 6.4 & 7.4 \\ 1 & 0 & 4.2 & 5.7 & 6.7 \\ 5 & 4.2 & 0 & 1.4 & 2.5 \\ 6.4 & 5.7 & 1.4 & 0 & 1.1 \\ 7.4 & 6.7 & 2.5 & 1.1 & 0 \end{bmatrix}$$



Remark: One can readily observe the **level in which a cluster is formed** and the **level in which it is absorbed** in a larger cluster (**indication of the natural clustering**).

Agglomerative Clustering Algorithms

Example:

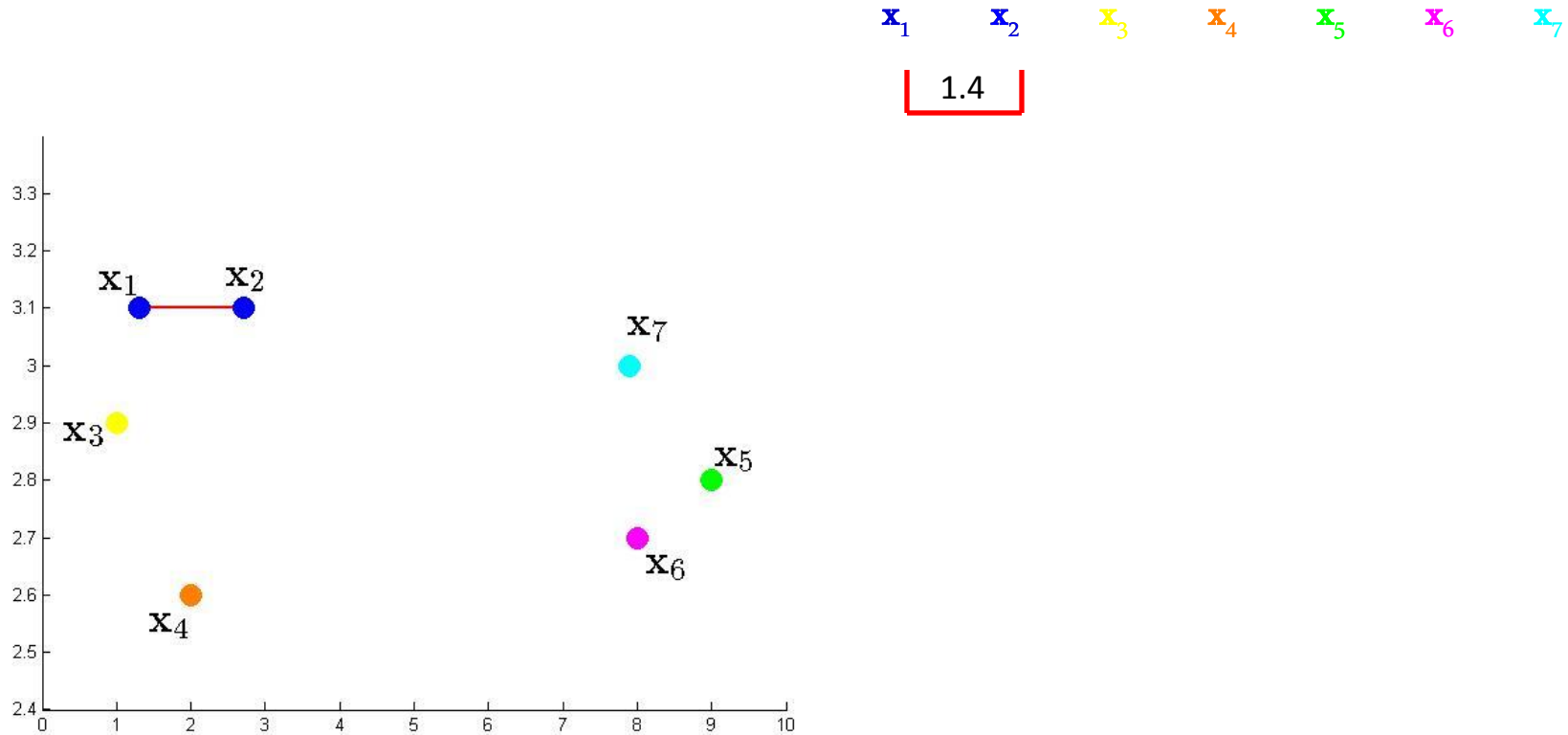


Agglomerative philosophy:

- In the **initial clustering** all **data vectors** **belong** to **different clusters**.
- At each step a **new clustering** is defined by **merging** the **two most similar clusters** to one.
- At the **final clustering** all **vectors** **belong** to the **same cluster**.

Agglomerative Clustering Algorithms

Example:

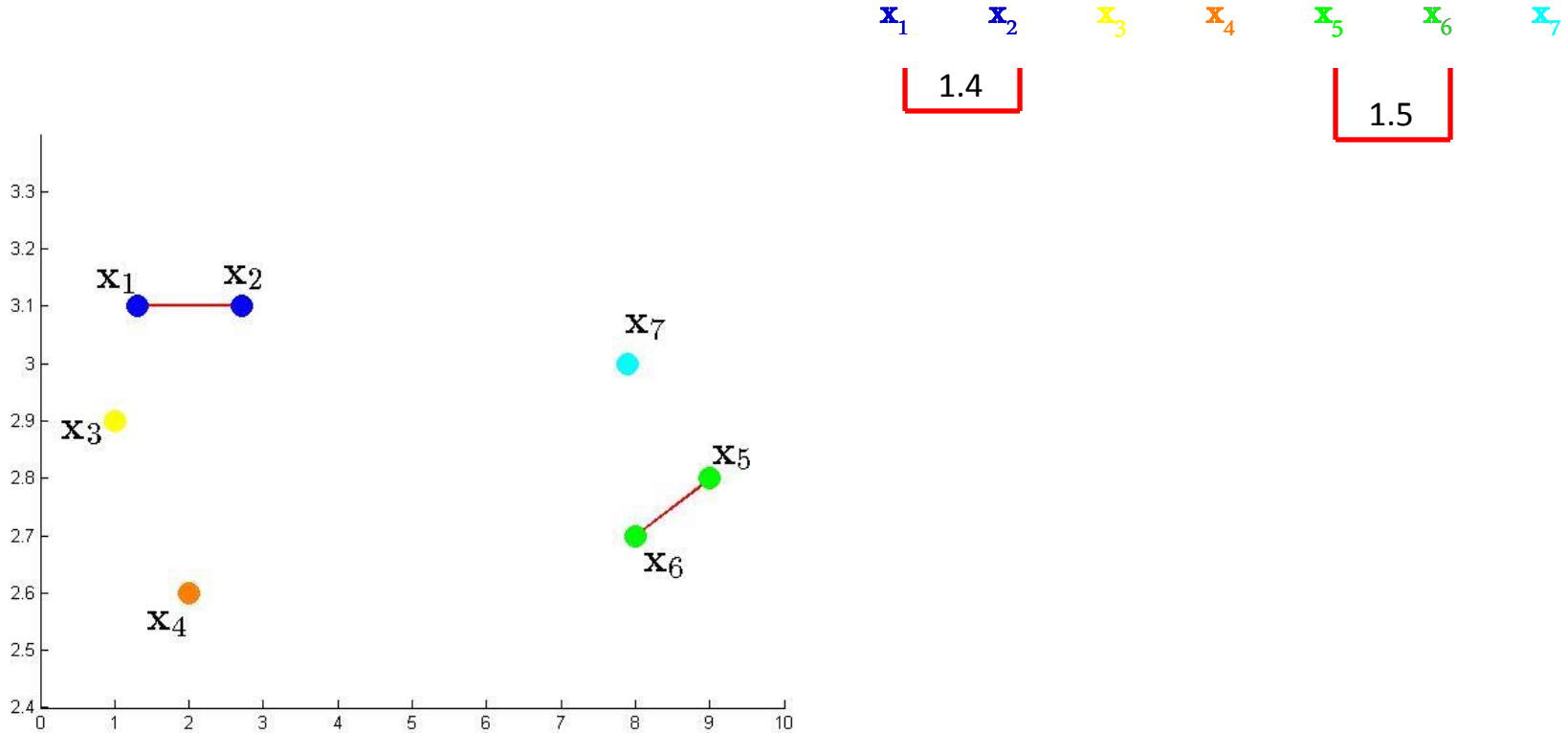


Agglomerative philosophy:

- In the **initial clustering** all **data vectors** **belong** to **different clusters**.
- At each step **a new clustering** is defined by **merging** the **two most similar clusters** to one.
- At the **final clustering** all **vectors** **belong** to the **same cluster**.

Agglomerative Clustering Algorithms

Example:

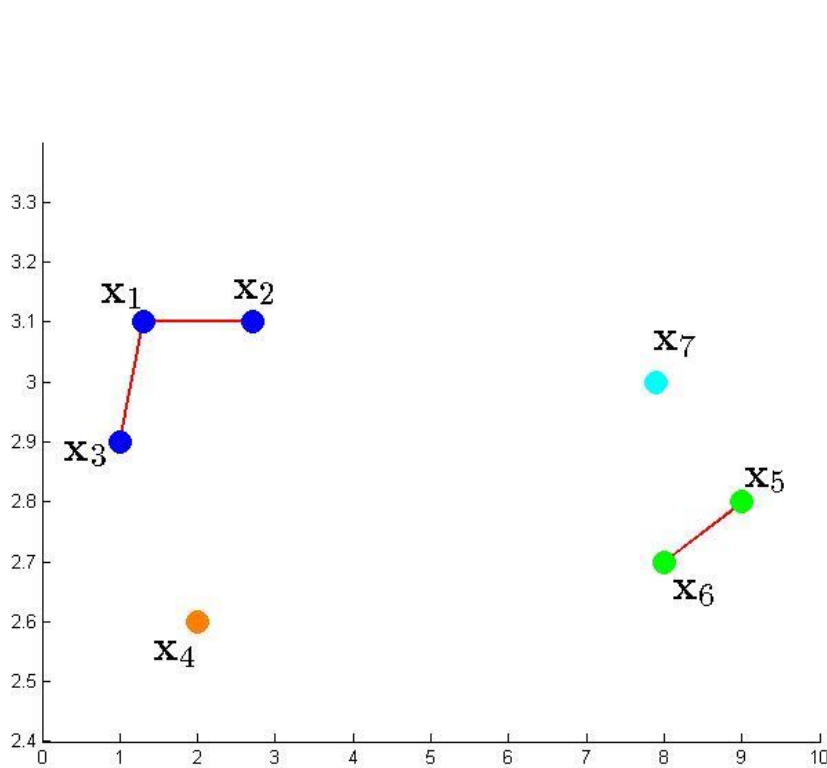


Agglomerative philosophy:

- In the **initial clustering** all **data vectors** **belong** to **different clusters**.
- At each step a **new clustering** is defined by **merging** the **two most similar clusters** to one.
- At the **final clustering** all **vectors** **belong** to the **same cluster**.

Agglomerative Clustering Algorithms

Example:

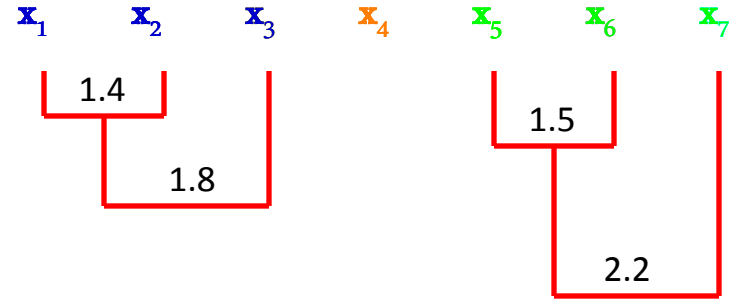
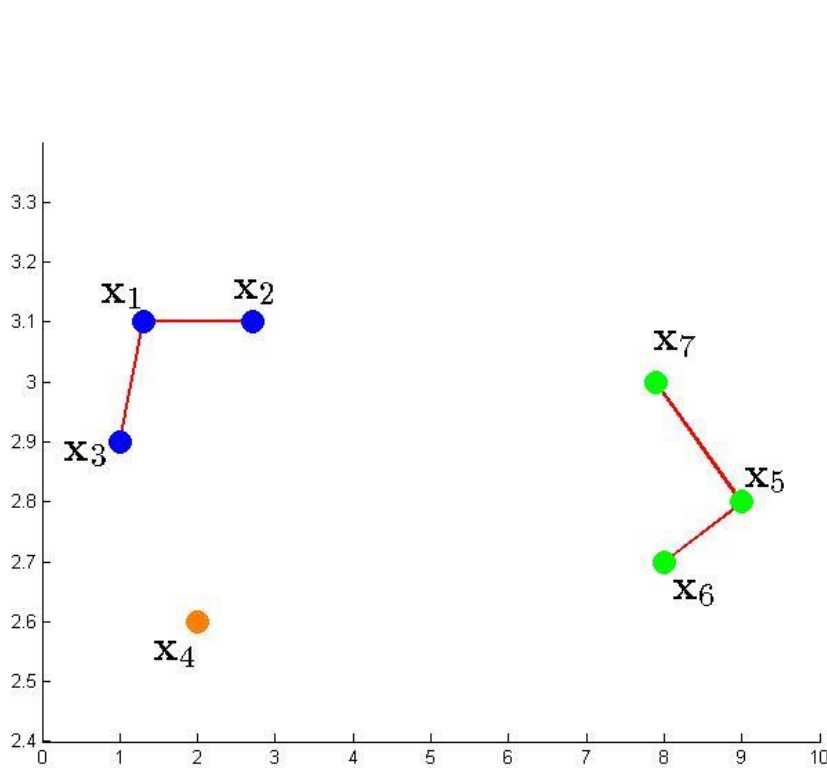


Agglomerative philosophy:

- In the **initial clustering** all **data vectors** **belong** to **different clusters**.
- At each step a **new clustering** is defined by **merging** the **two most similar clusters** to one.
- At the **final clustering** all **vectors** **belong** to the **same cluster**.

Agglomerative Clustering Algorithms

Example:

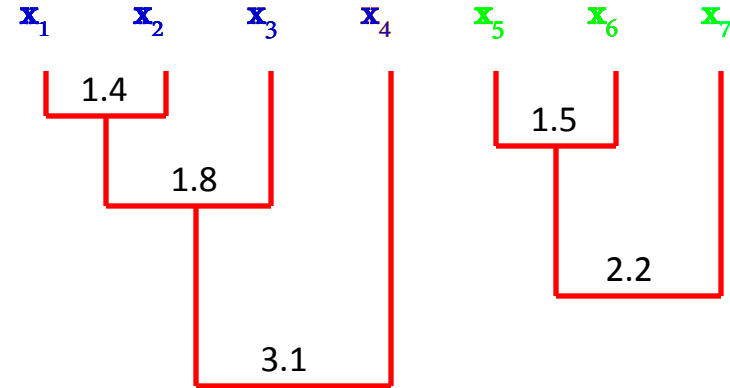
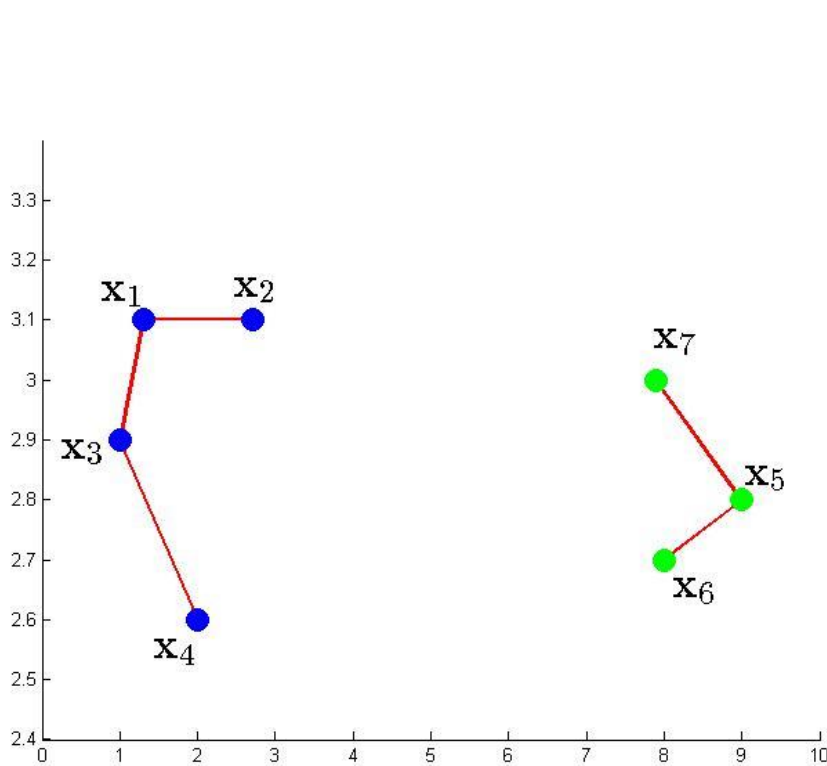


Agglomerative philosophy:

- In the **initial clustering** all **data vectors** **belong** to **different clusters**.
- At each step a **new clustering** is defined by **merging** the **two most similar clusters** to one.
- At the **final clustering** all **vectors** **belong** to the **same cluster**.

Agglomerative Clustering Algorithms

Example:

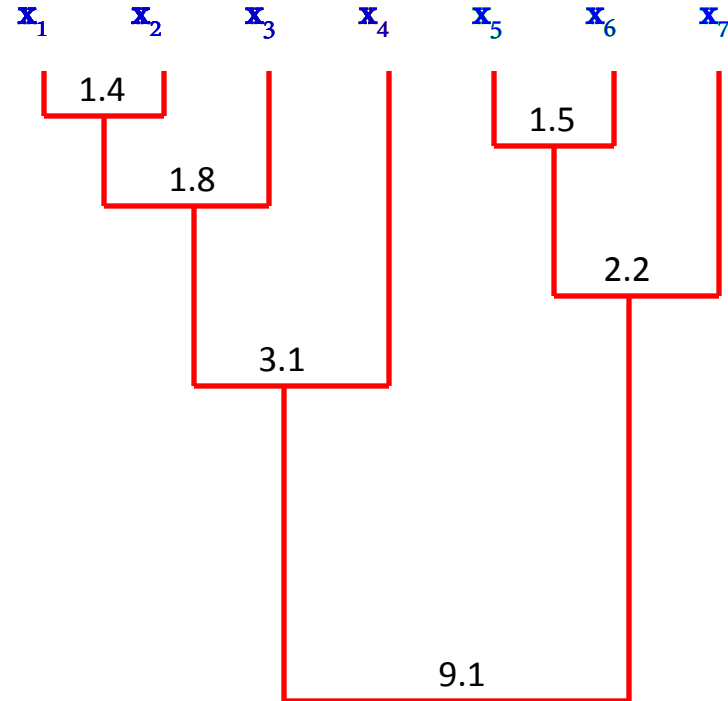
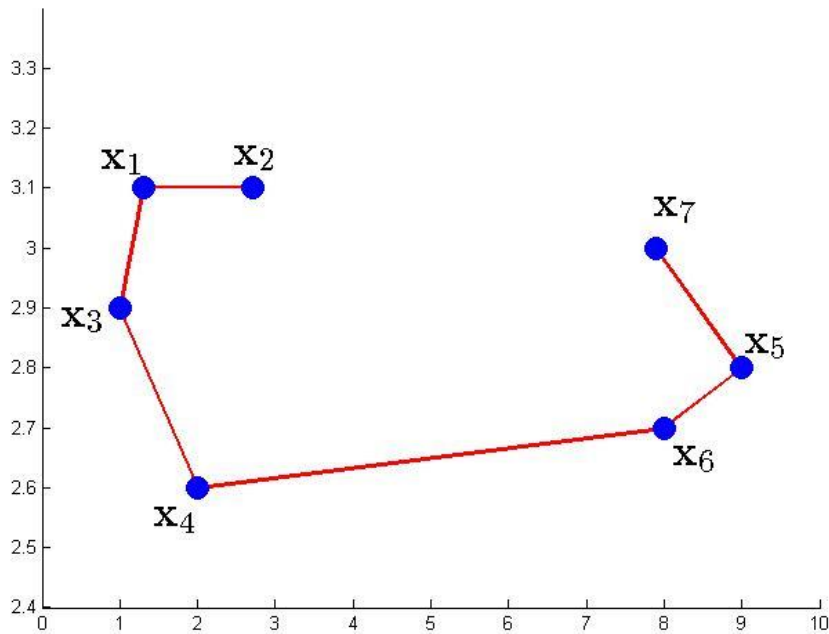


Agglomerative philosophy:

- In the **initial clustering** all **data vectors** **belong** to **different clusters**.
- At each step a **new clustering** is defined by **merging** the **two most similar clusters** to one.
- At the **final clustering** all **vectors** **belong** to the **same cluster**.

Agglomerative Clustering Algorithms

Example:

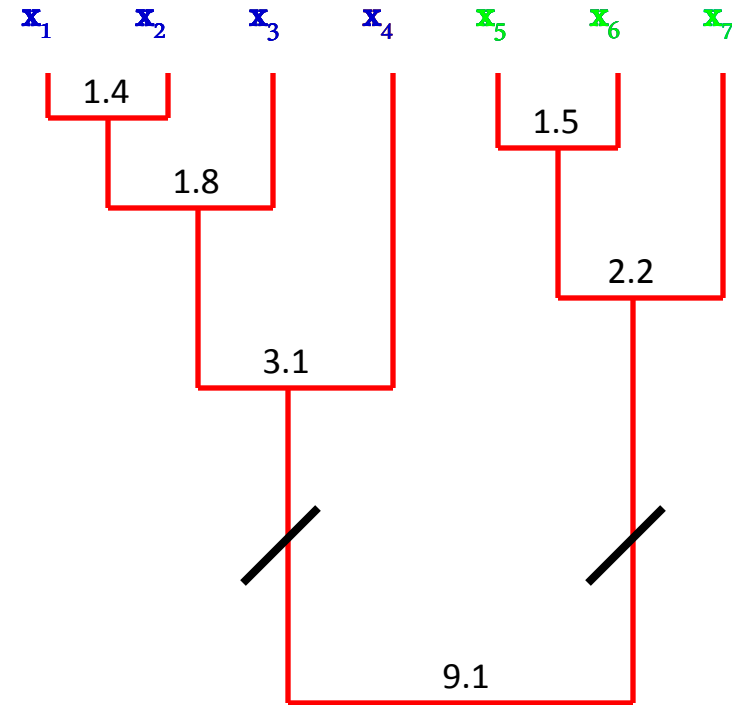
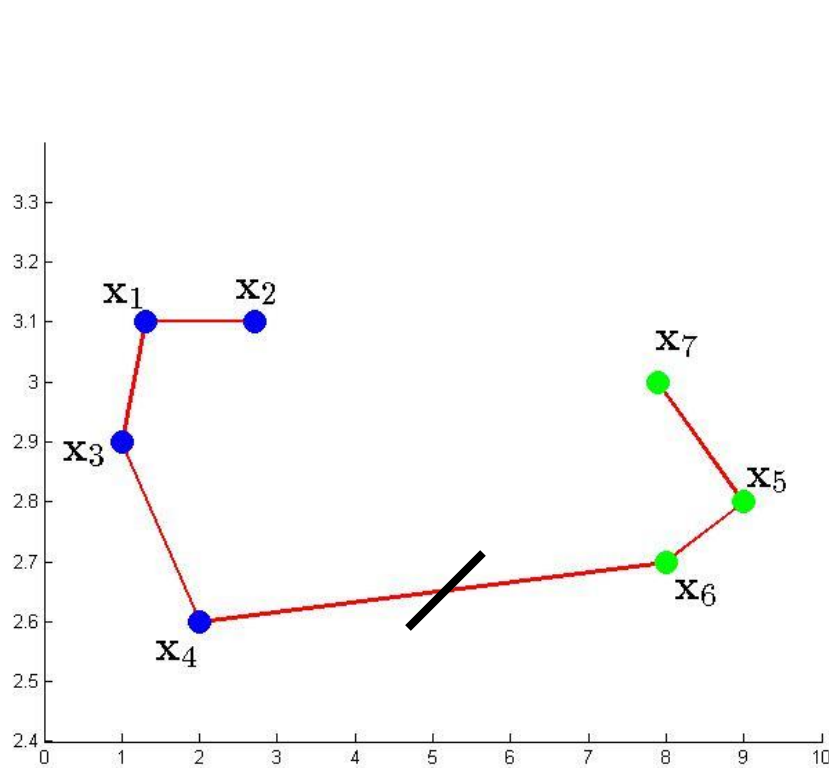


Agglomerative philosophy:

- In the **initial clustering** all **data vectors** belong to **different clusters**.
- At each step a **new clustering** is defined by **merging** the **two most similar clusters** to one.
- At the **final clustering** all **vectors** belong to the **same cluster**.

Agglomerative Clustering Algorithms

Example:



Agglomerative philosophy:

- In the **initial clustering** all **data vectors** belong to **different clusters**.
- At each step a **new clustering** is defined by **merging** the **two most similar clusters** to one.
- At the **final clustering** all **vectors** belong to the **same cluster**.

Agglomerative Clustering Algorithms

According to the mathematical tools used for their expression, **agglomerative algorithms** are divided into:

- Algorithms based on **matrix theory**.
- Algorithms based on **graph theory**.

NOTE: In the sequel we consider only **dissimilarity measures**.

➤ Algorithms based on matrix theory.

- They take as input the $N \times N$ dissimilarity matrix $P_0 = P(X)$.
- At each **level** t where two clusters C_i and C_j are **merged** to C_q , the dissimilarity matrix P_t is extracted from P_{t-1} by:
 - **Deleting** the two **rows** and **columns** of P_t that **correspond** to C_i and C_j .
 - **Adding** a **new row** and a **new column** that contain the **distances** of **newly formed** $C_q = C_i \cup C_j$ from each of the **remaining clusters** C_s , via a relation of the form

$$d(C_q, C_s) = f(d(C_i, C_s), d(C_j, C_s), d(C_i, C_j))$$

Agglomerative matrix theory based Clustering Algorithms

- A number of distance functions comply with the following update equation

$$d(C_q, C_s) = a_i d(C_i, C_s) + a_j (d(C_j, C_s) + b d(C_i, C_j) + c |d(C_i, C_s) - d(C_j, C_s)|) \quad (1)$$

Algorithms that follow the above equation are:

- **Single link (SL) algorithm** ($a_i = 1/2, a_j = 1/2, b = 0, c = -1/2$). In this case

$$d(C_q, C_s) = \min\{d(C_i, C_s), d(C_j, C_s)\} \quad (2)$$

- **Complete link (CL) algorithm** ($a_i = 1/2, a_j = 1/2, b = 0, c = 1/2$). In this case

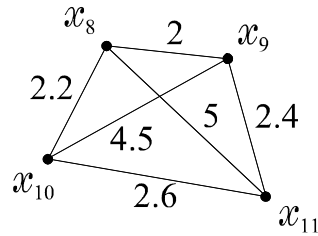
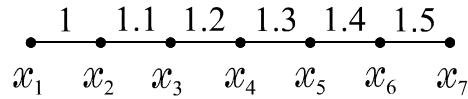
$$d(C_q, C_s) = \max\{d(C_i, C_s), d(C_j, C_s)\}$$

Remarks:

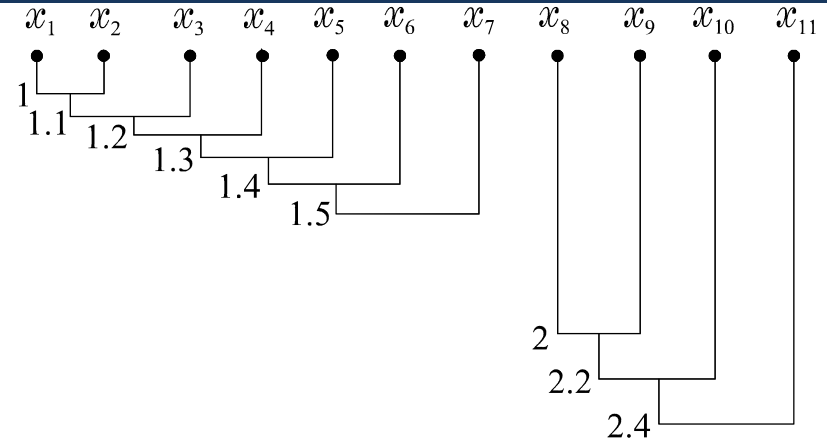
- **Single link** forms clusters at **low dissimilarities** while **complete link** forms clusters at **high dissimilarities**.
- **Single link** tends to form **elongated clusters** (*chaining effect*) while **complete link** tends to form **compact clusters**.
- The **rest algorithms** are **compromises** between these two extremes.

Agglomerative matrix theory based Clustering Algorithms

Example:



(a)

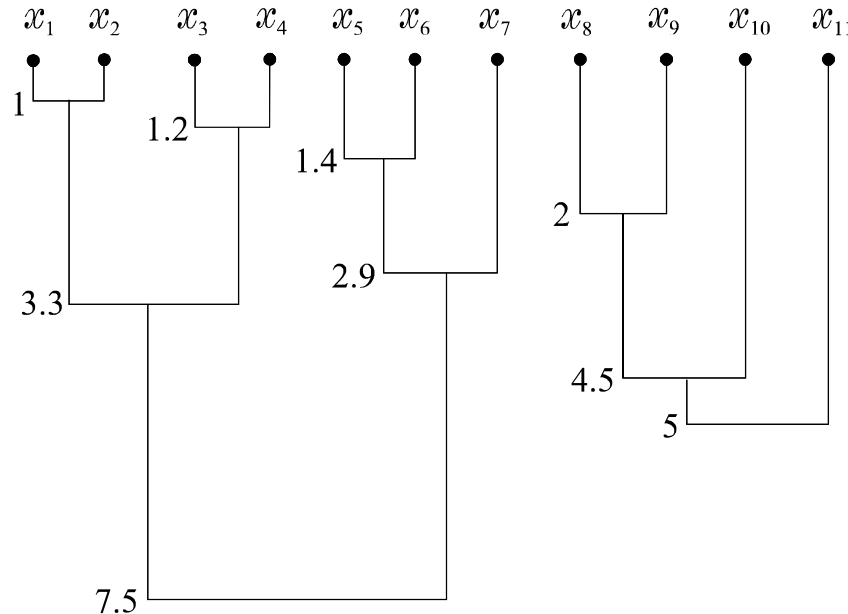


(b)

(a) The data set X .

(b) The **single link** algorithm **dissimilarity dendrogram**.

(c) The **complete link** algorithm **dissimilarity dendrogram**.



(c)

Agglomerative matrix theory based Clustering Algorithms

- **Weighted Pair Group Method Average (WPGMA)** ($a_i = 1/2, a_j = 1/2, b = 0, c = 0$). In this case:

$$d(C_q, C_s) = a_i d(C_i, C_s) + a_j (d(C_j, C_s) + b d(C_i, C_j) + c |d(C_i, C_s) - d(C_j, C_s)|)$$

$$d(C_q, C_s) = \frac{1}{2} (d(C_i, C_s) + d(C_j, C_s))$$

- **Unweighted Pair Group Method Average (UPGMA)** ($a_i = n_i/(n_i + n_j), a_j = n_j/(n_i + n_j), b = 0, c = 0$, where n_i is the **cardinality** of C_i). In this case:

$$d(C_q, C_s) = \frac{n_i}{n_i + n_j} d(C_i, C_s) + \frac{n_j}{n_i + n_j} d(C_j, C_s)$$

- **Unweighted Pair Group Method Centroid (UPGMC)** ($a_i = n_i/(n_i + n_j), a_j = n_j/(n_i + n_j), b = -n_i n_j/(n_i + n_j)^2, c = 0$). In this case:

$$d_{qs} = \frac{n_i}{n_i + n_j} d_{is} + \frac{n_j}{n_i + n_j} d_{js} - \frac{n_i n_j}{(n_i + n_j)^2} d_{ij}$$

For the **UPGMC**, if d_{ij} is defined as the **squared Euclidean distance** between the **means** of C_i and C_j , then it holds that $d_{qs} = \|\mathbf{m}_q - \mathbf{m}_s\|^2$, where \mathbf{m}_q is the **mean** of C_q .

Agglomerative matrix theory based Clustering Algorithms

- **Weighted Pair Group Method Centroid (WPGMC)** ($a_i = 1/2, a_j = 1/2, b = -1/4, c = 0$). In this case

$$d_{qs} = \frac{1}{2}d_{is} + \frac{1}{2}d_{js} - \frac{1}{4}d_{ij}$$

$$d(C_q, C_s) = a_i d(C_i, C_s) + a_j (d(C_j, C_s) + b d(C_i, C_j) + c |d(C_i, C_s) - d(C_j, C_s)|)$$

For **WPGMC** there are cases where $d_{qs} \leq \max\{d_{is}, d_{js}\}$ (**crossover**)

- **Ward or minimum variance algorithm.** Here the distanced d'_{ij} between C_i and C_j is defined as

$$d'_{ij} = \frac{n_i n_j}{n_i + n_j} \|\mathbf{m}_i - \mathbf{m}_j\|^2 \quad (3)$$

d'_{qs} can be expressed in terms of $d'_{is}, d'_{js}, d'_{ij}$ as

$$d'_{qs} = \frac{n_i + n_s}{n_i + n_j + n_s} d'_{is} + \frac{n_j + n_s}{n_i + n_j + n_s} d'_{js} - \frac{n_s}{n_i + n_j + n_s} d'_{ij}$$

Remark: Ward's algorithm forms \mathfrak{R}_{t+1} by merging the two clusters that lead to the **smallest possible increase of the total variance, i.e.,**

$$E_t = \sum_{r=1}^{N-t} \sum_{x \in C_r} \|\mathbf{x} - \mathbf{m}_r\|^2$$

Agglomerative matrix theory based Clustering Algorithms

Example 3: Consider the following dissimilarity matrix (Euclidean distance)

$$P_0 = \begin{bmatrix} 0 & 1 & 2 & 26 & 37 \\ 1 & 0 & 3 & 25 & 36 \\ 2 & 3 & 0 & 16 & 25 \\ 26 & 25 & 16 & 0 & 1.5 \\ 37 & 36 & 25 & 1.5 & 0 \end{bmatrix}$$

$$\mathcal{R}_0 = \{ \{x_1\}, \{x_2\}, \{x_3\}, \{x_4\}, \{x_5\} \},$$

$$\mathcal{R}_1 = \{ \{x_1, x_2\}, \{x_3\}, \{x_4\}, \{x_5\} \},$$

$$\mathcal{R}_2 = \{ \{x_1, x_2\}, \{x_3\}, \{x_4, x_5\} \},$$

$$\mathcal{R}_3 = \{ \{x_1, x_2, x_3\}, \{x_4, x_5\} \},$$

$$\mathcal{R}_4 = \{ \{x_1, x_2, x_3, x_4, x_5\} \}$$

All the algorithms produce the same sequence of clusterings shown above, yet at **different** proximity levels:

	<i>SL</i>	<i>CL</i>	<i>WPGMA</i>	<i>UPGMA</i>	<i>WPGMC</i>	<i>UPGMC</i>	<i>Ward</i>
\mathcal{R}_0	0	0	0	0	0	0	0
\mathcal{R}_1	1	1	1	1	1	1	0.5
\mathcal{R}_2	1.5	1.5	1.5	1.5	1.5	1.5	0.75
\mathcal{R}_3	2	3	2.5	2.5	2.25	2.25	1.5
\mathcal{R}_4	16	37	25.75	27.5	24.69	26.46	31.75

Agglomerative matrix theory based Clustering Algorithms

Example 3 (in detail): (a) The **single-link** case

$$(C_q = C_i \cup C_j, d(C_q, C_s) = \min(d(C_i, C_s), d(C_j, C_s)))$$

$$d(\{x_1, x_2\}, \{x_3\}) = \min(d(\{x_1\}, \{x_3\}), d(\{x_2\}, \{x_3\})) = \min(2, 3) = 2$$

$$d(\{x_1, x_2\}, \{x_4\}) = \min(26, 25) = 25$$

$$d(\{x_1, x_2\}, \{x_5\}) = \min(37, 36) = 36$$

P_0 :

	$\{x_1\}$	$\{x_2\}$	$\{x_3\}$	$\{x_4\}$	$\{x_5\}$		$\{x_1\}$	$\{x_2\}$	$\{x_3\}$	$\{x_4\}$	$\{x_5\}$
$\{x_1\}$	0	1	2	26	37	$\{x_1\}$	0	1	2	26	37
$\{x_2\}$	1	0	3	25	36	$\{x_2\}$	1	0	3	25	36
$\{x_3\}$	2	3	0	16	25	$\{x_3\}$	2	3	0	16	25
$\{x_4\}$	26	25	16	0	1.5	$\{x_4\}$	26	25	16	0	1.5
$\{x_5\}$	37	36	25	1.5	0	$\{x_5\}$	37	36	25	1.5	0

P_1 :

	$\{x_1, x_2\}$	$\{x_3\}$	$\{x_4\}$	$\{x_5\}$		$\{x_1, x_2\}$	$\{x_3\}$	$\{x_4\}$	$\{x_5\}$
$\{x_1, x_2\}$	0	2	25	36	$\{x_1, x_2\}$	0	2	25	36
$\{x_3\}$	2	0	16	25	$\{x_3\}$	2	0	16	25
$\{x_4\}$	25	16	0	1.5	$\{x_4\}$	25	16	0	1.5
$\{x_5\}$	36	25	1.5	0	$\{x_5\}$	36	25	1.5	0

$$d(\{x_1, x_2\}, \{x_4, x_5\}) = \min(25, 36) = 25$$

$$d(\{x_3\}, \{x_4, x_5\}) = \min(16, 25) = 16$$

Agglomerative matrix theory based Clustering Algorithms

Example 3 (in detail): (a) The **single-link** case

$$(C_q = C_i \cup C_j, d(C_q, C_s) = \min(d(C_i, C_s), d(C_j, C_s)))$$

P_2 :

	$\{x_1, x_2\}$	$\{x_3\}$	$\{x_4, x_5\}$
$\{x_1, x_2\}$	0	2	25
$\{x_3\}$	2	0	16
$\{x_4, x_5\}$	25	16	0

→

	$\{x_1, x_2\}$	$\{x_3\}$	$\{x_4, x_5\}$
$\{x_1, x_2\}$	0	2	25
$\{x_3\}$	2	0	16
$\{x_4, x_5\}$	25	16	0

$$d(\{x_1, x_2, x_3\}, \{x_4, x_5\}) = \min(25, 16) = 16$$

P_3 :

	$\{x_1, x_2, x_3\}$	$\{x_4, x_5\}$
$\{x_1, x_2, x_3\}$	0	16
$\{x_4, x_5\}$	16	0

→

	$\{x_1, x_2, x_3\}$	$\{x_4, x_5\}$
$\{x_1, x_2, x_3\}$	0	16
$\{x_4, x_5\}$	16	0

P_4 :

	$\{x_1, x_2, x_3, x_4, x_5\}$
$\{x_1, x_2, x_3, x_4, x_5\}$	0

$$\mathcal{R}_0 = \{\{x_1\}, \{x_2\}, \{x_3\}, \{x_4\}, \{x_5\}\}, (0)$$

$$\mathcal{R}_1 = \{\{x_1, x_2\}, \{x_3\}, \{x_4\}, \{x_5\}\}, (1)$$

$$\mathcal{R}_2 = \{\{x_1, x_2\}, \{x_3\}, \{x_4, x_5\}\}, (1.5)$$

$$\mathcal{R}_3 = \{\{x_1, x_2, x_3\}, \{x_4, x_5\}\}, (2)$$

$$\mathcal{R}_4 = \{\{x_1, x_2, x_3, x_4, x_5\}\}, (16)$$

Agglomerative matrix theory based Clustering Algorithms

Example 3 (in detail): (b) The **complete-link** case

$$(C_q = C_i \cup C_j, d(C_q, C_s) = \max(d(C_i, C_s), d(C_j, C_s)))$$

$$d(\{x_1, x_2\}, \{x_3\}) = \max(d(\{x_1\}, \{x_3\}), d(\{x_2\}, \{x_3\})) = \max(2, 3) = 3$$

$$d(\{x_1, x_2\}, \{x_4\}) = \max(26, 25) = 26$$

$$d(\{x_1, x_2\}, \{x_5\}) = \max(37, 36) = 37$$

P_0 :

	$\{x_1\}$	$\{x_2\}$	$\{x_3\}$	$\{x_4\}$	$\{x_5\}$		$\{x_1\}$	$\{x_2\}$	$\{x_3\}$	$\{x_4\}$	$\{x_5\}$
$\{x_1\}$	0	1	2	26	37	$\{x_1\}$	0	1	2	26	37
$\{x_2\}$	1	0	3	25	36	$\{x_2\}$	1	0	3	25	36
$\{x_3\}$	2	3	0	16	25	$\{x_3\}$	2	3	0	16	25
$\{x_4\}$	26	25	16	0	1.5	$\{x_4\}$	26	25	16	0	1.5
$\{x_5\}$	37	36	25	1.5	0	$\{x_5\}$	37	36	25	1.5	0

P_1 :

	$\{x_1, x_2\}$	$\{x_3\}$	$\{x_4\}$	$\{x_5\}$		$\{x_1, x_2\}$	$\{x_3\}$	$\{x_4\}$	$\{x_5\}$
$\{x_1, x_2\}$	0	3	26	37	$\{x_1, x_2\}$	0	3	26	37
$\{x_3\}$	3	0	16	25	$\{x_3\}$	3	0	16	25
$\{x_4\}$	26	16	0	1.5	$\{x_4\}$	26	16	0	1.5
$\{x_5\}$	37	25	1.5	0	$\{x_5\}$	37	25	1.5	0

$$d(\{x_1, x_2\}, \{x_4, x_5\}) = \max(26, 37) = 37$$

$$d(\{x_3\}, \{x_4, x_5\}) = \max(16, 25) = 25$$

Agglomerative matrix theory based Clustering Algorithms

Example 3 (in detail): (b) The **complete-link** case

$$(C_q = C_i \cup C_j, d(C_q, C_s) = \max(d(C_i, C_s), d(C_j, C_s)))$$

P_2 :

	$\{x_1, x_2\}$	$\{x_3\}$	$\{x_4, x_5\}$
$\{x_1, x_2\}$	0	3	37
$\{x_3\}$	3	0	25
$\{x_4, x_5\}$	37	25	0

→

	$\{x_1, x_2\}$	$\{x_3\}$	$\{x_4, x_5\}$
$\{x_1, x_2\}$	0	3	37
$\{x_3\}$	3	0	25
$\{x_4, x_5\}$	37	25	0

$$d(\{x_1, x_2, x_3\}, \{x_4, x_5\}) = \max(37, 25) = 37$$

P_3 :

	$\{x_1, x_2, x_3\}$	$\{x_4, x_5\}$
$\{x_1, x_2, x_3\}$	0	37
$\{x_4, x_5\}$	37	0

→

	$\{x_1, x_2, x_3\}$	$\{x_4, x_5\}$
$\{x_1, x_2, x_3\}$	0	37
$\{x_4, x_5\}$	37	0

P_4 :

	$\{x_1, x_2, x_3, x_4, x_5\}$
$\{x_1, x_2, x_3, x_4, x_5\}$	0

$$\mathcal{R}_0 = \{\{\underline{x}_1\}, \{\underline{x}_2\}, \{\underline{x}_3\}, \{\underline{x}_4\}, \{\underline{x}_5\}\}, (0)$$

$$\mathcal{R}_1 = \{\{\underline{x}_1, \underline{x}_2\}, \{\underline{x}_3\}, \{\underline{x}_4\}, \{\underline{x}_5\}\}, (1)$$

$$\mathcal{R}_2 = \{\{\underline{x}_1, \underline{x}_2\}, \{\underline{x}_3\}, \{\underline{x}_4, \underline{x}_5\}\}, (1.5)$$

$$\mathcal{R}_3 = \{\{\underline{x}_1, \underline{x}_2, \underline{x}_3\}, \{\underline{x}_4, \underline{x}_5\}\}, (3)$$

$$\mathcal{R}_4 = \{\{\underline{x}_1, \underline{x}_2, \underline{x}_3, \underline{x}_4, \underline{x}_5\}\}, (37)$$