

Clustering algorithms

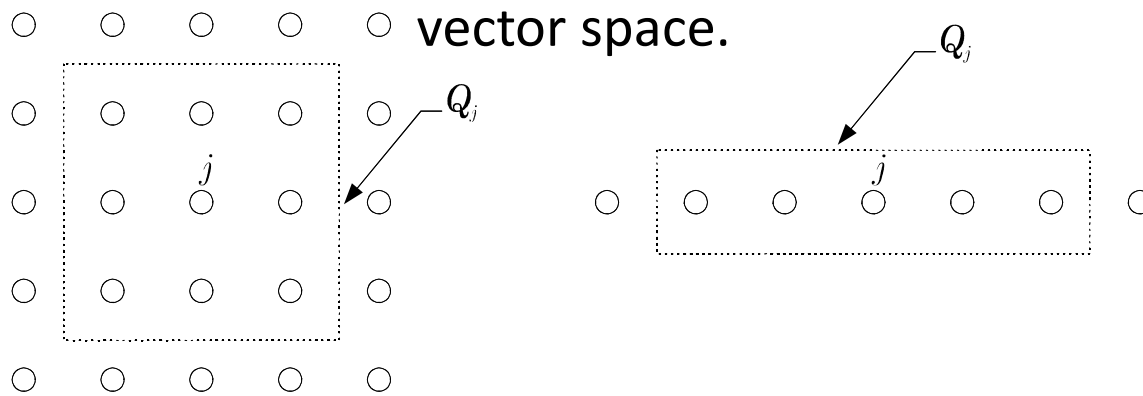
Konstantinos Koutroumbas

Unit 10

- SOM, LVQ
- Valley seeking clust. Algorithms
- Branch & bound clustering algorithms
- Simulated & Deterministic annealing
- Genetic algorithms
- Density-based algorithms for large data sets (DBSCAN)

Self-organizing maps

- It is used for **data visualization** (maps high dim. Data → 1-d or 2-d maps) and (“loose”) **clustering**.
- Here **interrelation between representatives** is assumed.
- For each representative w_j a **topological neighborhood** of **representatives** $Q_j(t)$ is defined, centered at w_j .
- As t (no. of iterations) **increases**, $Q_j(t)$ **shrinks** and concentrates around w_j .
- The **neighborhood** is defined with respect to the indices j and it is **independent** of the **geometrical distances** between representatives in the



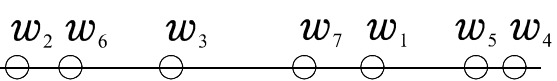
(a)

(b)

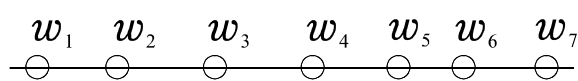
Assuming that in fig. (c) the neighborhood of w_j constitutes of w_{j-1} and w_{j+1} ,

w_2 and w_4 are the **topological neighbors** of w_3 although

w_6 and w_7 are closer in terms of the geometrical distance to w_3)



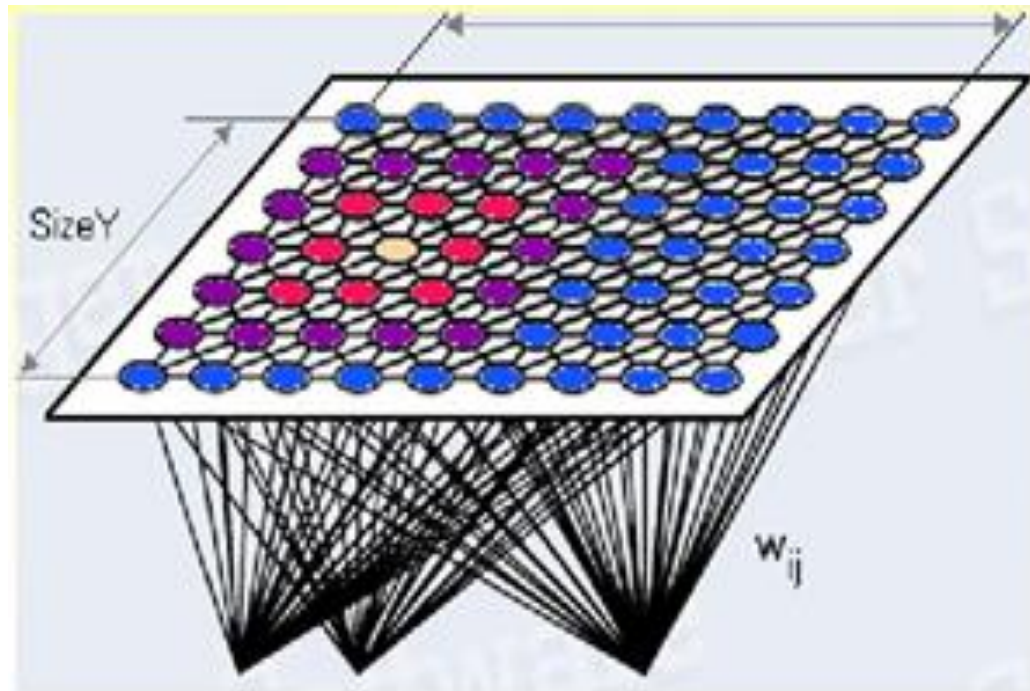
(c)



(d)

Self-organizing maps

- Here **interrelation** between representatives is assumed.
- For each representative w_j a **neighborhood of representatives $Q_j(t)$** is defined, centered at w_j .
- As t (number of iterations) **increases**, $Q_j(t)$ **shrinks** and concentrates around w_j .
- **The neighborhood is defined with respect to the indices j and it is independent of the distances between representatives in the vector space.**



Self-organizing maps

- If \mathbf{w}_j wins on the current input \mathbf{x} all the representatives in $Q_j(t)$ are updated (**Self Organizing Map (SOM) scheme**).
- SOM (in its simplest version) may be viewed as a special case of GCLS if
 - Parts (A), (B) and (C) are defined as in the basic competitive learning scheme.
 - In part (D), if \mathbf{w}_j wins on \mathbf{x} , the updating equation becomes:

$$\mathbf{w}_k(t) = \begin{cases} \mathbf{w}_k(t-1) + \eta_t^{k,j}(\mathbf{x} - \mathbf{w}_k(t-1)), & \text{if } \mathbf{w}_k \in Q_j(t) \\ \mathbf{w}_k(t-1), & \text{otherwise} \end{cases}$$

where $\eta_t^{k,j}$ is a variable learning rate, which decreases with t and with the topological distance between the k -th and the j -th representatives.

- After convergence, neighboring representatives also lie “close” in terms of their geometrical distance in the vector space (**topographical ordering**) (see fig. (d)).

Self-organizing maps

The algorithm

➤ $t = 0$

➤ **Repeat**

• $t = t + 1$

• **Present** a new randomly selected $\mathbf{x} \in X$ to the algorithm.

• (B) **Determine** the **winning** representative \mathbf{w}_j on \mathbf{x} as the one for which

$$d(\mathbf{x}, \mathbf{w}_j(t-1)) = \min_{k=1, \dots, m} d(\mathbf{x}, \mathbf{w}_k(t-1))$$

• (D) *Parameter updating*

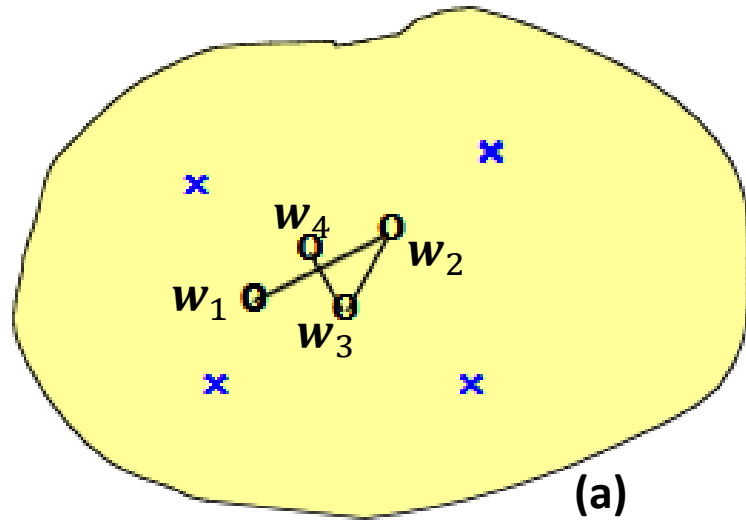
$$\mathbf{w}_k(t) = \begin{cases} \mathbf{w}_k(t-1) + \eta_t^{k,j} (\mathbf{x} - \mathbf{w}_k(t-1)), & \text{if } \mathbf{w}_k \in Q_j(t) \\ \mathbf{w}_k(t-1), & \text{otherwise} \end{cases}$$

• End

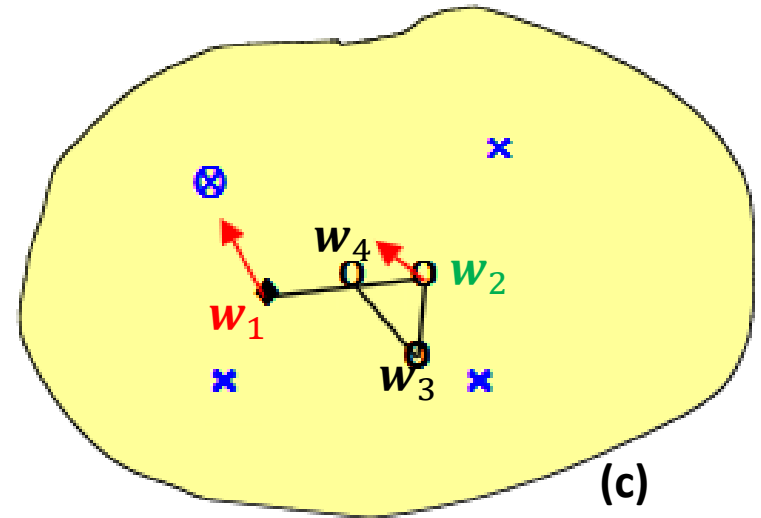
➤ (E) **Until** (convergence occurred) *OR* ($t > t_{max}$)

Self-organizing maps

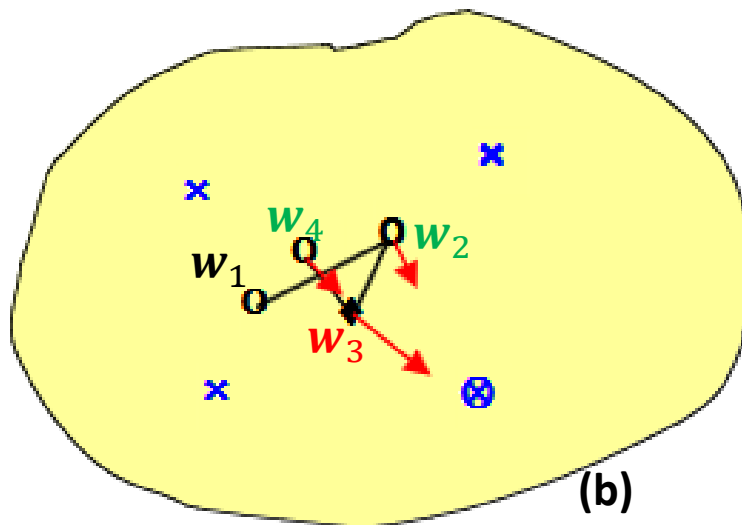
Example



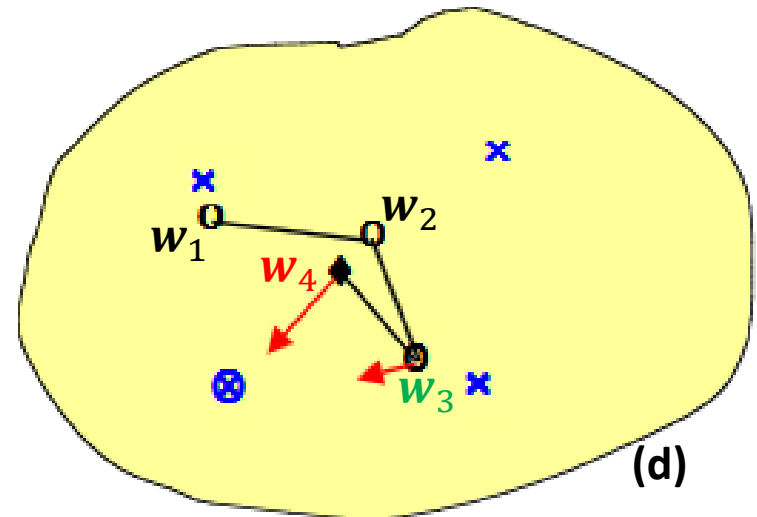
(a)



(c)



(b)



(d)

Self-organizing maps

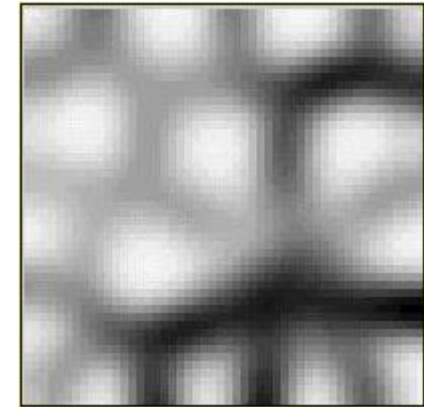
How to represent the result of a SOM

NOTE: After SOM convergence the topological ordering of the m **representatives** will comply with their “geometrical ordering”.

Produce an image A of size

- m for the 1-d case or
- $k \times k$ for the 2-d case ($m = k^2$)

As follows



For each representative (pixel of A):

Compute its **average distance** d_{avg} from its neighboring representatives

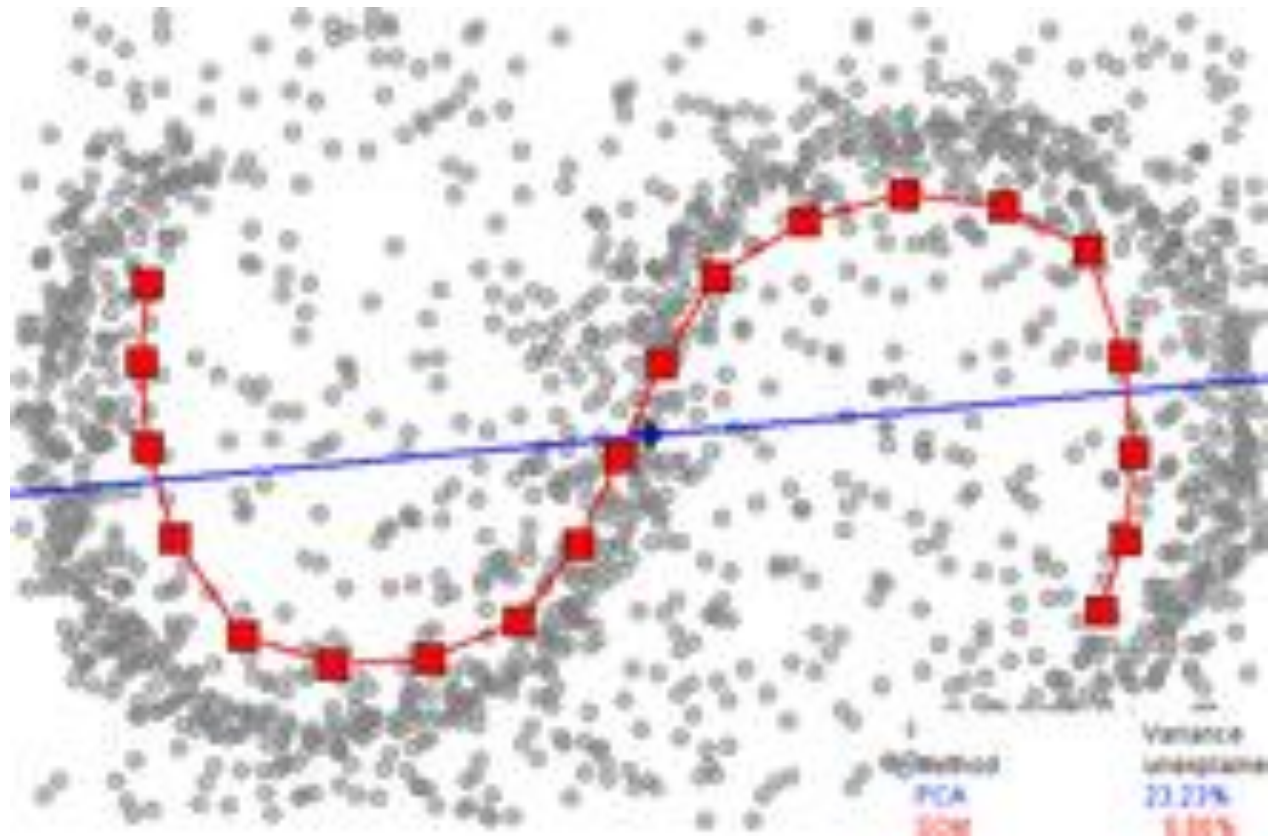
Draw the associate pixel of A with a color so that:

The larger the d_{avg} , the darker the color will be.

Then **lighter areas surrounded by darker areas** in A are indicative of clustering structure in the data.

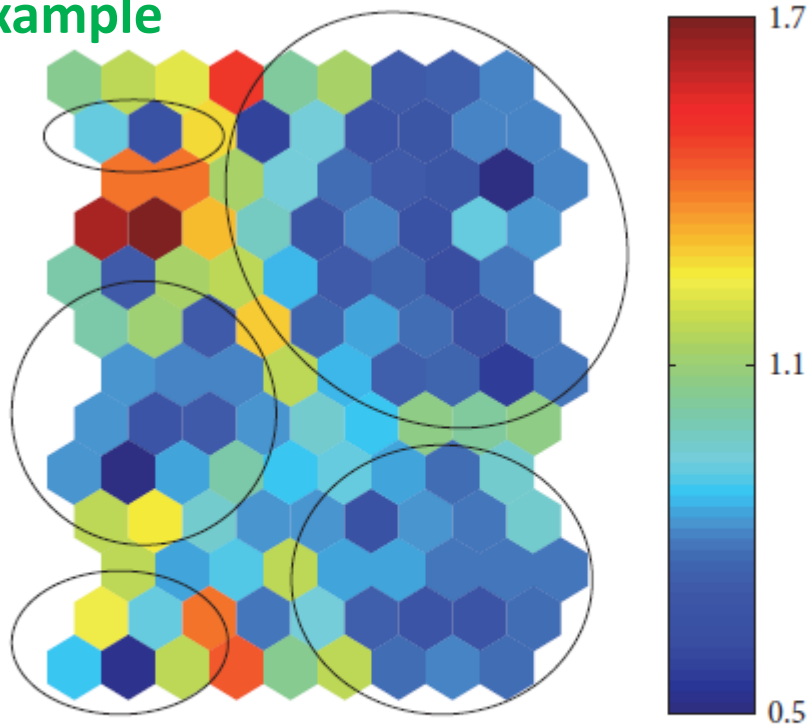
Self-organizing maps

Example

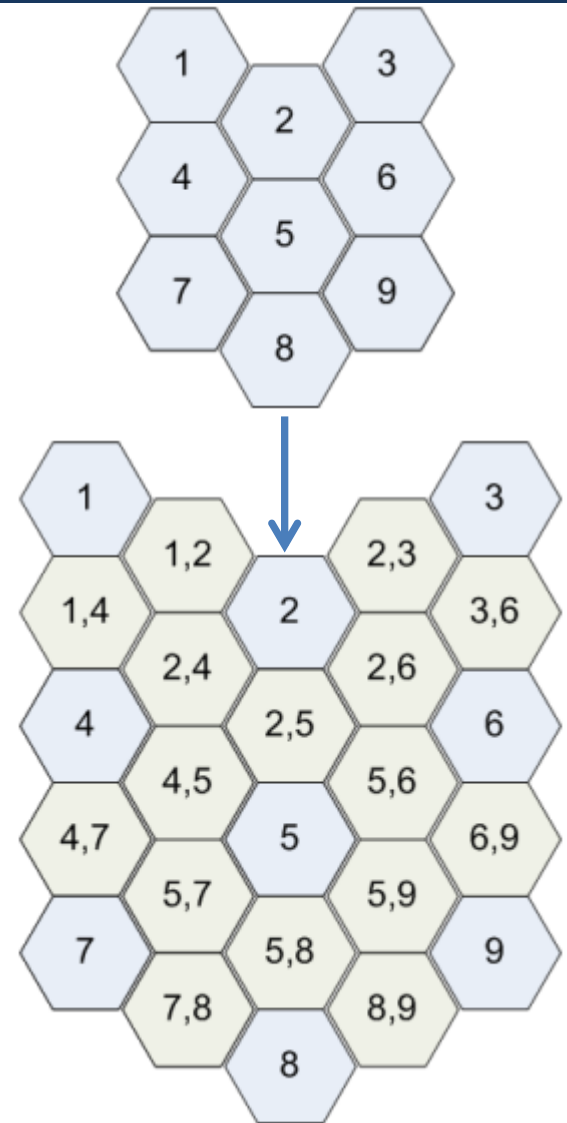


Self-organizing maps

Example



U-matrix of the 54 ERCs texts preprocessed through the VSM. Color scale is related to distances between map units. Red colors represent large distances and blue colors represent small distances.



From: Massimo Pacella, Antonio Grieco, Marzia Blaco, "On the Use of Self-Organizing Map for Text Clustering in Engineering Change Process Analysis: A Case Study", *Computational Intelligence and Neuroscience*, vol. 2016, Article ID 5139574, 11 pages, 2016. <https://doi.org/10.1155/2016/5139574>

Computation of the numbers in U -matrix via example:
 $\{2,3\} \rightarrow d(w_2, w_3)$
 $\{3\} \rightarrow \text{mean}\{d(w_2, w_3), d(w_3, w_6)\}$

Supervised Learning Vector Quantization (VQ)

In this case

- each **cluster** is **treated** as a **class** (m **compact** classes are assumed)
- the available vectors have **known class labels**.

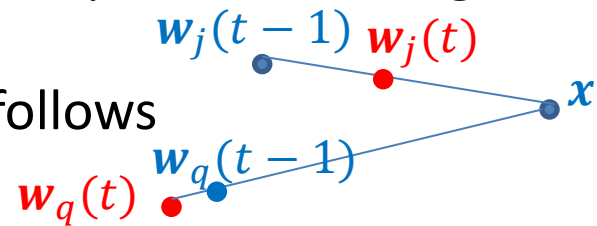
The goal:

Use a set of m representatives and place them in such a way so that each class is “optimally” represented.

The simplest version of VQ (LVQ1) may be obtained from GCLS as follows:

- Parts (A), (B) and (C) are the same with the basic competitive learning scheme.
- In part (D) the updating for w_j 's is carried out as follows

$$w_j(t) = \begin{cases} w_j(t-1) + \eta(t) (x - w_j(t-1)), & \text{if } w_j \text{ correctly wins on } x \\ w_j(t-1) - \eta(t) (x - w_j(t-1)), & \text{if } w_j \text{ wrongly wins on } x \\ w_j(t-1), & \text{otherwise} \end{cases}$$



Supervised Learning Vector Quantization (VQ)

The algorithm

- $t = 0$
- **Repeat**
 - $t = t + 1$
 - **Present** a new randomly selected $\mathbf{x} \in X$ to the algorithm.
 - (B) **Determine** the **winning** representative \mathbf{w}_j on \mathbf{x} as the one for which
$$d(\mathbf{x}, \mathbf{w}_j(t-1)) = \min_{k=1, \dots, m} d(\mathbf{x}, \mathbf{w}_k(t-1))$$
 - (D) *Parameter updating*
$$\mathbf{w}_j(t) = \begin{cases} \mathbf{w}_j(t-1) + \eta(t) (\mathbf{x} - \mathbf{w}_j(t-1)), & \text{if } \mathbf{w}_j \text{ correctly wins on } \mathbf{x} \\ \mathbf{w}_j(t-1) - \eta(t) (\mathbf{x} - \mathbf{w}_j(t-1)), & \text{if } \mathbf{w}_j \text{ wrongly wins on } \mathbf{x} \\ \mathbf{w}_j(t-1), & \text{otherwise} \end{cases}$$
- (E) **Until** (convergence occurred) *OR* ($t > t_{max}$) (max allowable no of iter.)

In words:

- \mathbf{w}_j is moved:
 - Towards \mathbf{x} if \mathbf{w}_j wins and \mathbf{x} belongs to the j -th class.
 - Away from \mathbf{x} if \mathbf{w}_j wins and \mathbf{x} does not belong to the j -th class.
- All other representatives remain unaltered.

Valley seeking clustering algorithms

Let $p(\mathbf{x})$ be the **density function** describing the **distribution** of the vectors in X .

➤ **Clusters** may be **viewed** as **peaks** of $p(\mathbf{x})$ **separated** by **valleys**.

Thus one may

- **Identify** these **valleys** and
- Try to **move** the **borders** of the clusters **in** these **valleys**.

A simple method in this spirit.

Preliminaries

➤ Let the **distance** $d(\mathbf{x}, \mathbf{y})$ be defined as

$$d(\mathbf{x}, \mathbf{y}) = (\mathbf{y} - \mathbf{x})^T A (\mathbf{y} - \mathbf{x})$$

where A is a **positive definite matrix**

➤ Let the **local region** of \mathbf{x} , $V(\mathbf{x})$, be defined as

$$V(\mathbf{x}) = \{\mathbf{y} \in X - \{\mathbf{x}\} : d(\mathbf{x}, \mathbf{y}) \leq a\}$$

where a is a user-defined parameter

➤ k_j^i be the **number of vectors** of the j **cluster** that belong to $V(\mathbf{x}_i) - \{\mathbf{x}_i\}$.

➤ $c_i \in \{1, \dots, m\}$ denote the **cluster** to which \mathbf{x}_i will be **assigned**.

Valley seeking clustering algorithms

Valley-Seeking algorithm

- **Fix** a .
- **Fix** the number of clusters m .
- **Define** an **initial clustering** X .
- **Repeat**
 - For $i = 1$ to N
 - Find** j : $k_j^i = \max_{q=1,\dots,m} k_q^i$
 - Set** $c_i = j$
 - End For

 - For $i = 1$ to N
 - Assign** x_i to cluster C_{c_i} .
 - End For
- **Until** **no reclustering** of vectors occurs.

Valley seeking clustering algorithms

The algorithm

- **Centers** a **window** defined by $d(x, y) \leq a$ at x and **counts** the **points from different clusters** in it.
- **Assigns** x to the cluster with the larger number of points in the window (the **cluster** that **corresponds** to the **highest local pdf**).

In other words:

- The **boundary** is **moved away** from the “**winning**” **cluster**.

Remarks:

- The **algorithm** is **sensitive** to a . It is suggested to perform several runs, for different values of a .
- The algorithm is of a **mode-seeking** nature (if more than enough clusters are initially appointed, some of them will become empty).

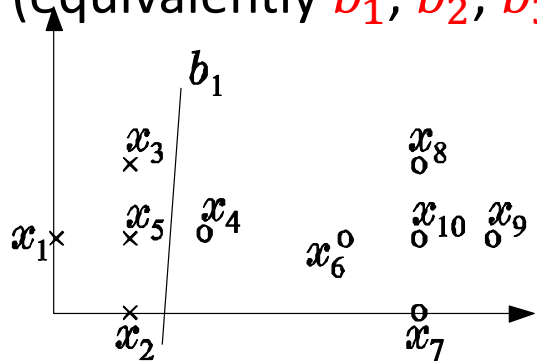
Valley seeking clustering algorithms

Example: Let $X = \{x_1, \dots, x_{10}\}$ and $a = 1.1415 (> \sqrt{2})$. X contains two physical clusters: $C_1 = \{x_1, \dots, x_5\}$, $C_2 = \{x_6, \dots, x_{10}\}$.

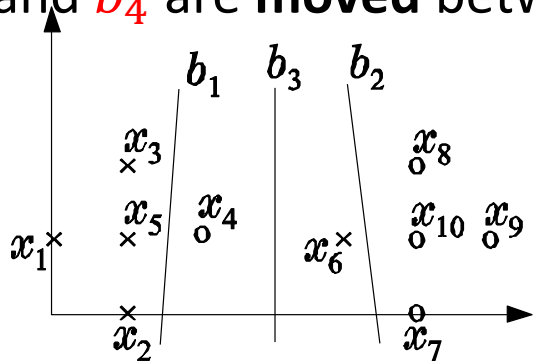
(a) **Initially two clusters** are considered **separated** by b_1 . After the convergence of the algorithm, C_1 and C_2 are identified (equivalently, b_1 is **moved** between x_4 and x_6).

(b) **Initially two clusters** are considered **separated** by b_1 , b_2 and b_3 . After the convergence of the algorithm, C_1 and C_2 are identified (equivalently b_1 and b_2 are **moved** to the **area** where b_3 lies).

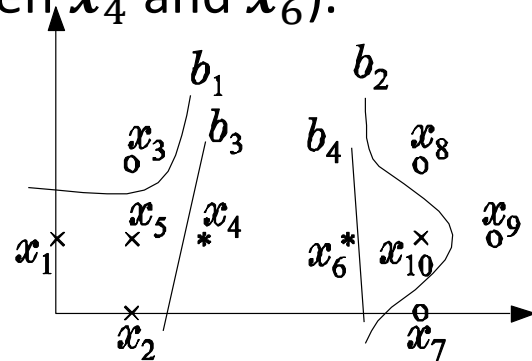
(c) **Initially three clusters** are considered **separated** by b_1 , b_2 , b_3 , b_4 . After the convergence of the algorithm, only two clusters are identified, C_1 and C_2 (equivalently b_1 , b_2 , b_3 and b_4 are **moved** between x_4 and x_6).



(a)



(b)



(c)

Branch and Bound Clustering algorithms

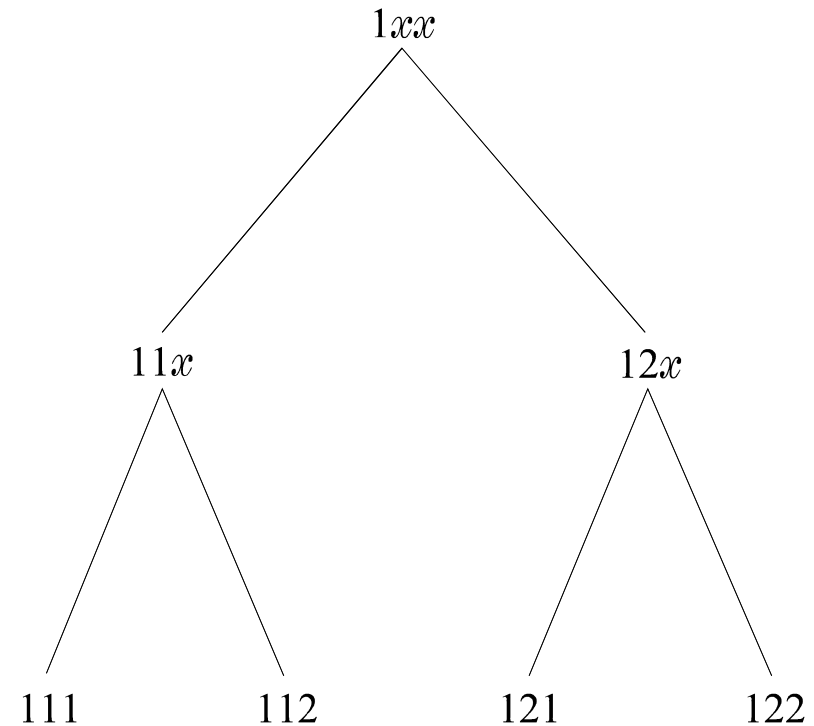
- They compute the **globally optimal solution** to combinatorial problems.
- They avoid exhaustive search via the employment of a **monotonic criterion** J .

Monotonic criterion J : if k vectors of X have been assigned to clusters, the assignment of an extra vector to a cluster **does not decrease** the value of J .

Consider the following 3-vectors, 2-class case:

121: 1st, 3rd vectors belong to class 1
2nd vector belongs to class 2.
(**leaf** of the **tree**)

12 x : 1st vector belongs to class 1
2nd vector belongs to class 2
3rd vector is unassigned
(**Partial clustering- node** of the **tree**).



Branch and Bound Clustering algorithms

How exhaustive search is avoided

- Let B be the **best value** for criterion J computed **so far**.
- **If** at a **node** of the tree, the **corresponding value** of J is **greater than B** , **no further search** is **performed** for **all subsequent** descendants springing from this node.
- Let $\mathbf{C}_r = [c_1, \dots, c_r]$, $1 \leq r \leq N$, denotes a **partial clustering** where $c_i \in \{1, 2, \dots, m\}$, $c_i = j$ if the vector \mathbf{x}_i belongs to cluster C_j and $\mathbf{x}_{r+1}, \dots, \mathbf{x}_N$ are yet unassigned.
- For **compact clusters** and **fixed** number of clusters, m , a suitable cost function is

$$J(\mathbf{C}_r) = \sum_{i=1}^r \|\mathbf{x}_i - \mathbf{m}_{c_i}(\mathbf{C}_r)\|^2$$

where \mathbf{m}_{c_i} is the **mean vector** of the cluster C_{c_i}

$$\mathbf{m}_j(\mathbf{C}_r) = \frac{1}{n_j(\mathbf{C}_r)} \sum_{\{q=1, \dots, r, c_q=j\}} \mathbf{x}_q, \quad j = 1, \dots, m$$

with $n_j(\mathbf{C}_r)$ being the number of vectors $\mathbf{x} \in \{\mathbf{x}_1, \dots, \mathbf{x}_r\}$ that belong to cluster C_j .

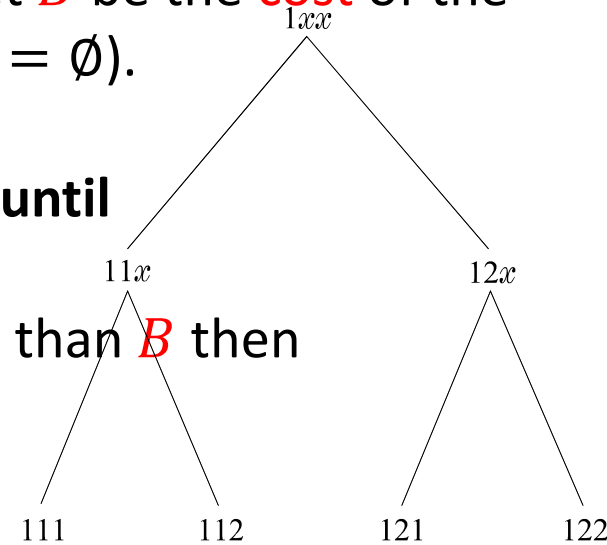
Branch and Bound Clustering algorithms

Initialization

- **Start** from the **initial node** and **go down** to a **leaf**. Let B be the **cost** of the **corresponding clustering C** (initially set $B = +\infty$, $C = \emptyset$).

Main stage

- **Start** from the **initial node** of the tree and **go down until**
 - **Either** (i) A **leaf** is encountered.
 - o If the cost B' of the corr. clustering C' is **smaller** than B then
 - * $B = B'$
 - * $C = C'$ is the best clustering found so far
 - o End if
 - **Or** (ii) a **node q** with **value** of J **greater** than B is encountered. Then
 - o **No** subsequent **clustering** branching **from q** is **considered**.
 - o **Backtrack** to the parent of q , q^{par} , in order to **span** a **different path**.
 - o If **all paths** branching from q^{par} have been **considered** then
 - * **Move** to the **grandparent** of q .
 - o End if
 - **End if**



Terminate when **all possible paths** have been **considered explicitly** or **implicitly**.

Branch and Bound Clustering algorithms

Remarks

- **Variations** of the above algorithm, where **much tighter bounds** of B are **used** (that is, many more clusterings are rejected without explicit consideration) have also been proposed.
- A **disadvantage** of the algorithm is the **excessive** (and **unpredictable**) amount of required **computational time**.

Simulated Annealing

- It **guarantees** (under certain conditions) **in probability**, the **determination** of the **globally optimal solution** of the problem at hand via the **minimization of a cost function J** .
- It **may escape** from **local minima** since it **allows** moves that **temporarily** may **increase** the value of J .

Definitions

- An important parameter of the algorithm is the “**temperature**” T , which **starts** at a **high value** and **reduces gradually**.
- A **sweep** is the **time** the algorithm **spends at a given temperature** so that the system can enter the “**thermal equilibrium**” in this temperature.

Notation

- T_{max} is the **initial value** of the temperature T .
- C_{init} is the **initial clustering**.
- C is the **current clustering**.
- t is the **current sweep**.

Simulated Annealing

The algorithm:

- **Set** $T = T_{max}$ and $C = C_{init}$.
- $t = 0$
- **Repeat**
 - $t = t + 1$
 - Repeat
 - o **Compute** $J(C)$
 - o **Produce** a new clustering, C' , by **assigning** a **randomly chosen vector** from X to a **different cluster**.
 - o **Compute** $J(C')$
 - o **If** $\Delta J = J(C') - J(C) < 0$ then
 - * (A) $C = C'$
 - o **Else**
 - * (B) $C = C'$, with **probability** $P(\Delta J) = e^{-\Delta J/T}$.
 - o **End if**
 - Until an **equilibrium state** is **reached** at this temperature.
 - $T = f(T_{max}, t)$
- **Until** a predetermined value T_{min} for T is reached

Simulated Annealing

Remarks:

- For $T \rightarrow \infty$, it is $p(\Delta J) \approx 1$. Thus almost all movements of vectors between clusters are allowed.
- For lower values of T fewer moves of type (B) (from lower to higher cost clusterings) are allowed.
- As $T \rightarrow 0$ the probability of moves of type (B) tends to zero.
- Thus as T decreases, it becomes more probable to reach clusterings that correspond to lower values of J .
- Keeping T positive, we ensure a nonzero probability for escaping from a local minimum.
- We assume that the equilibrium state has been reached
"If for k successive random reassignments of vectors, C remains unchanged."
- A schedule for lowering T that guarantees convergence to the global minimum with probability 1, is

$$T = \frac{T_{max}}{\ln(1+t)}$$

- The method is computationally demanding.

Deterministic Annealing (DA)

➤ It is inspired by the **phase transition phenomenon** observed when the temperature of a material changes. It **involves** the parameter $\beta = 1/T$, where T is defined as in simulated annealing.

➤ **The Goal of DA: Locate** a set of representatives $\mathbf{w}_j, j = 1, \dots, m$ (m is fixed) in appropriate **positions** so that a distortion function J is **minimized**.

J is defined as

$$J = -\frac{1}{\beta} \sum_{i=1}^N \ln \left(\sum_{j=1}^m e^{-\beta d(\mathbf{x}_i, \mathbf{w}_j)} \right)$$

Assumption: $d(\mathbf{x}, \mathbf{w})$ is a convex function of \mathbf{w} for fixed \mathbf{x} .

➤ Then, the **optimal value** of a specific \mathbf{w}_r satisfies the following condition:

$$\frac{\partial J}{\partial \mathbf{w}_r} = \sum_{i=1}^N P_{ir} \frac{\partial d(\mathbf{x}_i, \mathbf{w}_r)}{\partial \mathbf{w}_r} = 0$$

where

$$P_{ir} = \frac{e^{-\beta d(\mathbf{x}_i, \mathbf{w}_r)}}{\sum_{j=1}^m e^{-\beta d(\mathbf{x}_i, \mathbf{w}_j)}}$$

➤ P_{ir} may be **interpreted** as the **probability** that \mathbf{x}_i belongs to $C_r, r = 1, \dots, m$.

Deterministic Annealing

Assumption: $d(\mathbf{x}, \mathbf{w})$ is a **convex function** of \mathbf{w} for **fixed** \mathbf{x} .

Stages of the algorithm

- For $\beta \rightarrow 0$, all P_{ij} 's are almost **equal** to $\frac{1}{m}$, for all \mathbf{x}_i 's, $i = 1, \dots, N$. Thus

$$\sum_{i=1}^N \frac{\partial d(\mathbf{x}_i, \mathbf{w}_r)}{\partial \mathbf{w}_r} = 0$$

Since $d(\mathbf{x}, \mathbf{w})$ is a convex function, $d(\mathbf{x}_1, \mathbf{w}_r) + \dots + d(\mathbf{x}_N, \mathbf{w}_r)$ is a convex function. All representatives coincide with its unique global minimum (all the data belong to a single cluster).

- As β **increases**, it **reaches** a **critical value** where P_{ir} 's “**depart sufficiently**” from the **uniform model**. Then the representatives split up in order to provide an optimal presentation of the data set at the new phase.
- The **increase** of β **continues** until P_{ij} **approach** the **hard clustering model** (for all \mathbf{x}_i , $P_{ir} \approx 1$ for a specific r , and $P_{ij} \approx 0$, for $j \neq r$).

Deterministic Annealing

Application: For the squared Euclidean distance $d(\mathbf{x}, \mathbf{w}) = (\mathbf{x} - \mathbf{w})^T (\mathbf{x} - \mathbf{w})$ it is

$$\frac{\partial J}{\partial \mathbf{w}_r} = \sum_{i=1}^N P_{ir} \frac{\partial d(\mathbf{x}_i, \mathbf{w}_r)}{\partial \mathbf{w}_r} = 2 \sum_{i=1}^N P_{ir} (\mathbf{x}_i - \mathbf{w}_r) = 0 \Leftrightarrow \mathbf{w}_r = \frac{\sum_{i=1}^N P_{ir} \mathbf{x}_i}{\sum_{i=1}^N P_{ir}}$$

This is **coupled** wrt \mathbf{w}_r

Remarks:

- It is **not guaranteed** that it **reaches** the **globally optimum clustering**.
- If m is chosen **greater than** the “**actual**” number of clusters, the **algorithm** has the ability to **represent** the **data properly**.

Clustering using genetic algorithms (GA)

A few hints concerning genetic algorithms

- They have been **inspired** by the **natural selection mechanism** (Darwin).
- They consider a **population of solutions** of the problem at hand and they **perform certain operators** on this, so that **the new population of the same size is improved compared to the previous one** (wrt a critierion function F).
- The **solutions are coded** and the **operators are applied** on the **coded** versions of the **solutions**.

The most well-known operators are:

Reproduction:

- It ensures that, in probability, the **better** (**worse**) a **solution** in the current population is, the **more** (**less**) **replicates** it has in the next population.
- A simple implementation:
 - For each solution s_i , out of the population of the p solutions, compute the associated criterion function value $F(s_i)$.
(it is assumed that the **higher** the value of F , the **better** the **solution**)
 - Assign to each s_i a probability $p_i = F(s_i) / \sum_{j=1}^p F(s_j)$.
 - Perform sampling with replacement to produce the next solution population.

Clustering using genetic algorithms (GA)

Crossover:

- It applies to the temporary population produced after the application of the reproduction operator.
- It **selects pairs of solutions** *randomly*, **splits** them at a *random position* and **exchanges** their **second parts**.

Mutation:

- It applies to the temporary population produced after the application of the crossover operator.
- It **selects** *randomly* an **element** of a solution and **alters** it with **some probability**.
- It may be viewed as **a way out** of **getting stuck** in **local minima**.

Aspects/Parameters that affect the performance of the algorithm

The **coding** of the **solutions**.

The **number of solutions** in a population, p .

The **probability** with which **two solutions** are **selected** for **crossover**.

The **probability** with which an **element** of a solution is **mutated**.

Clustering using genetic algorithms (GA)

GA Algorithmic scheme

$t = 0$

Choose an **initial population** \wp_t of solutions.

Repeat

- **Apply reproduction** on \wp_t and let \wp'_t be the resulting **temporary population**.
- **Apply crossover** on \wp'_t and let \wp''_t be the resulting **temporary population**.
- **Apply mutation** on \wp''_t and let \wp_{t+1} be the **resulting population**.
- $t = t + 1$

Until a termination condition is met.

Return

- **either** the **best solution** of the **last population**,
- **or** the **best solution** found **during the evolution** of the algorithm.

Clustering using genetic algorithms (GA)

Genetic Algorithms in Clustering

The characteristics of a simple **GA hard clustering algorithm** suitable for **compact clusters**, whose number m is fixed, is discussed next.

A (not unique) way to **code a solution** is **via the cluster representatives**.

The **cost function** in use is

$$J = \sum_{i=1}^N u_{ij} d(\mathbf{x}_i, \mathbf{w}_j)$$

$[\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_m]$

The critierion function can be defined e.g., as
 $F(s_i) = e^{-J(s_i)}$

where

$$u_{ij} = \begin{cases} 1, & \text{if } d(\mathbf{x}_i, \mathbf{w}_j) = \min_{k=1, \dots, m} d(\mathbf{x}_i, \mathbf{w}_k), \\ 0, & \text{otherwise} \end{cases}, i = 1, \dots, N$$

The **allowable cut points** for the **crossover** operator **are between different representatives**.

The **mutation operator** **selects** randomly **a coordinate** and **decides** randomly to add a small random number to it.

Clustering using genetic algorithms (GA)

Remark:

- An **alternative** to the above scheme results if **prior to the application of the reproduction operator**, the hard clustering algorithm (GHAS), described in a previous lecture, **runs p times**, each time **using a different solution** of the current population **as the initial state**. The p resulting solutions constitute the population on which the reproduction operator will be applied.

Density-based algorithms for large data sets

These algorithms:

- Consider **clusters** as **regions** in the l -dimensional space that are “**dense**” in **points** of X .
- Have, in principle, the ability to **recover arbitrarily shaped clusters** (however, difficulties may arise in the case where the clusters differ in terms of their density).
- **Handle** efficiently **outliers**.
- Have **time complexity less** than $O(N^2)$.

Typical density-based algorithms are:

- The **DBSCAN** algorithm.
- The **DBCLASD** algorithm.
- The **DENCLUE** algorithm.

Density-based algorithms for large data sets

Density-Based Spatial Clustering of Applications with Noise (DBSCAN)

Algorithm

The “density” around a point \mathbf{x} is **estimated** as the **number of points in X** that **fall inside a specific region** of the l -dimensional space **surrounding \mathbf{x}** .

Notation

- $V_\varepsilon(\mathbf{x})$ is the **hypersphere** of **radius ε** (user-defined parameter) **centered** at \mathbf{x} .
- $N_\varepsilon(\mathbf{x})$ the **number of points of X lying** in $V_\varepsilon(\mathbf{x})$.
- q is the **minimum number** of **points of X** that must be contained **in $V_\varepsilon(\mathbf{x})$** , in order for \mathbf{x} to be considered an “**interior**” **point** of a cluster.

Definitions

1. A point \mathbf{y} is **directly density reachable** from a point $\mathbf{x} \in X$ if
 - (i) $\mathbf{y} \in V_\varepsilon(\mathbf{x})$
 - (ii) $N_\varepsilon(\mathbf{x}) \geq q$ (fig. (a)).
2. A point \mathbf{y} is **density reachable** from a point $\mathbf{x} \in X$ if there is a **sequence** of points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p \in X$, with $\mathbf{x}_1 = \mathbf{x}$, $\mathbf{x}_p = \mathbf{y}$, such that \mathbf{x}_{i+1} is **directly density reachable** from \mathbf{x}_i (fig. (b)).

Density-based algorithms for large data sets

DBSCAN Algorithm (cont.)

3. A point x is **density connected** to a point $y \in X$ if there exists $z \in X$ such that both x and y are **density reachable** from z (fig. (c)).

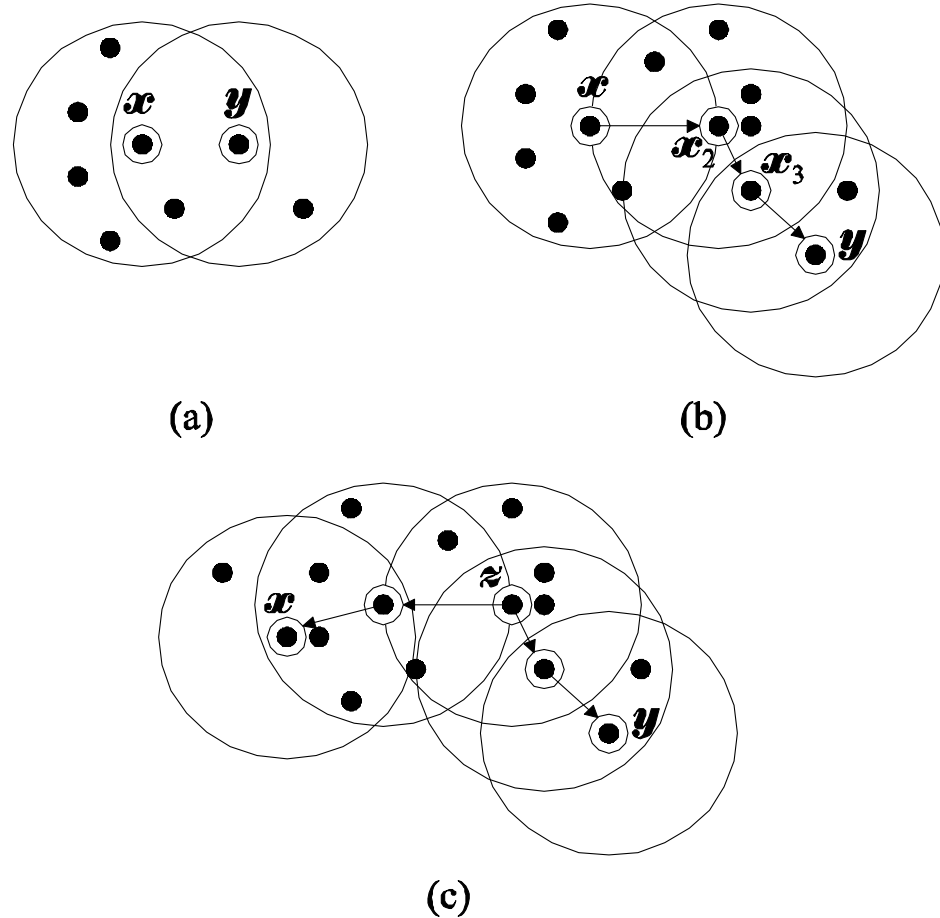
Example:

Assuming that $q = 5$,

(a) y is **directly density reachable** from x , but not vice versa,

(b) y is **density reachable** from x , but not vice versa, and

(c) x and y are **density connected** (in addition, y is **density reachable** from x , but not vice versa).



Density-based algorithms for large data sets

DBSCAN Algorithm (cont.)

4. A **cluster** C in DBSCAN is defined as a nonempty subset of X satisfying the following conditions:
 - If x belongs to C and $y \in X$ is **density reachable** from x , then $y \in C$.
 - For each pair $(x, y) \in C$, x and y are **density connected**.
5. Let C_1, \dots, C_m be the clusters in X . The set of **points** that are **not connected** in any of the C_1, \dots, C_m is known as **noise**.
6. A point x is called a **core** (**noncore**) **point** if it has **at least** (**less than**) q points in its neighborhood.
A **noncore point** may be either
 - a **border point** of a cluster (that is, density reachable from a core point) or
 - a **noisy point** (that is, not density reachable from other points in X).

Density-based algorithms for large data sets

DBSCAN Algorithm (cont.)

Proposition 1: If x is a **core point** and D is the **set** of points in X that are **density reachable from x** , then D is a **cluster**.

Proposition 2: If C is a **cluster** and x is a **core point** in C , then C **equals** to the **set** of the points $y \in X$ that are **density reachable** from x .

Therefore: A **cluster** is **uniquely determined** by any of its **core points**.

Notation

- X_{un} is the set of **points** in X that have **not been considered yet**.
- m denotes the **number of clusters**.

Density-based algorithms for large data sets

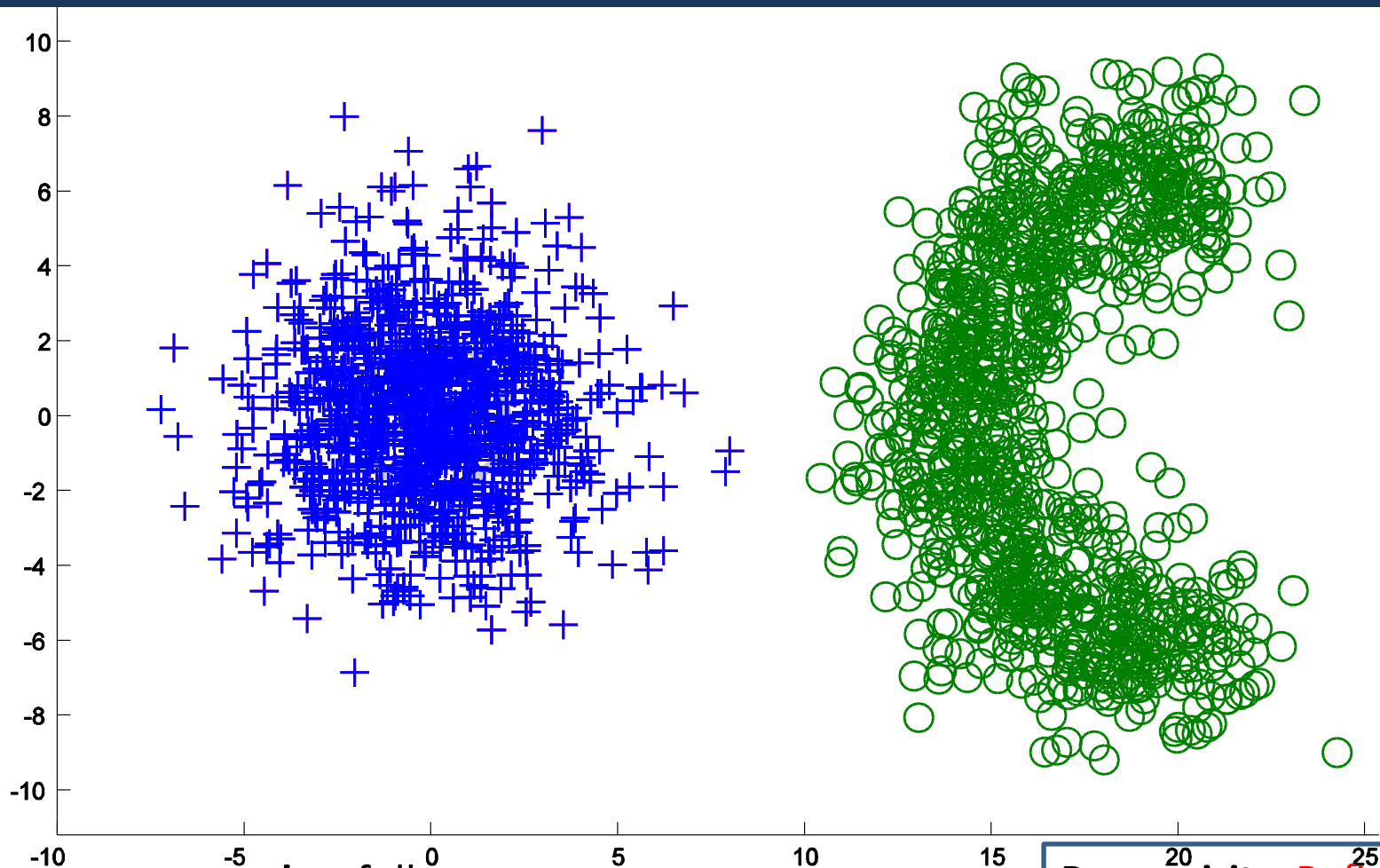
DBSCAN Algorithm (cont.)

DBSCAN Algorithm

- Set $X_{un} = X$
- Set $m = 0$
- While $X_{un} \neq \emptyset$ do
 - Arbitrarily **select** a $x \in X_{un}$
 - If x is a **noncore point** then
 - **Mark** x as **noise point**
 - $X_{un} = X_{un} - \{x\}$
 - End if

 - If x is a **core point** then
 - $m = m + 1$
 - **Determine all density-reachable** points $y \in X$ from x .
 - **Assign** x and the **previous points** to the cluster C_m . The **border points** among them that may have been **marked** as “**noise**” are also **assigned** to C_m .
 - $X_{un} = X_{un} - C_m$
 - End {if}
- **End** {while}

Clustering – Density-based algorithms



Clusters are recovered as follows:

- Start a new cluster C by choosing a data point x .
- Assign all the data points that lie in the neighborhood of x to the same cluster.
- Repeat recursively the previous step until all neighboring points of ALL $x \in C$ are assigned to C .

Prerequisite: Definition of the neighborhood size

Density-based algorithms for large data sets

DBSCAN Algorithm (cont.)

Important notes:

- If a **border point** y of a cluster C is selected, it will be **marked initially** as a **noise point**. However, **when (a)** a core point x in C is selected later on, and **(b)** y is identified as a density-reachable point from x **then** y will assigned to C .
- If an **actual noise point** y is selected, it will be marked as such and **since it is not density reachable** by any of the core points in X , its “**noise**” label will **remain unaltered**.

Remarks:

- The parameters ε and q influence significantly the performance of DBSCAN. These should **selected** such that the algorithm is **able to detect the least “dense” cluster** (experimentation with several values for ε and q should be carried out).
- Implementation of DBSCAN using R^* -tree data structure can achieve time complexity of $O(N \log_2 N)$ for low-dimensional data sets.
- **DBSCAN is not well suited** for cases where **clusters exhibit significant differences in density** as well as for applications of **high-dimensional data**.