

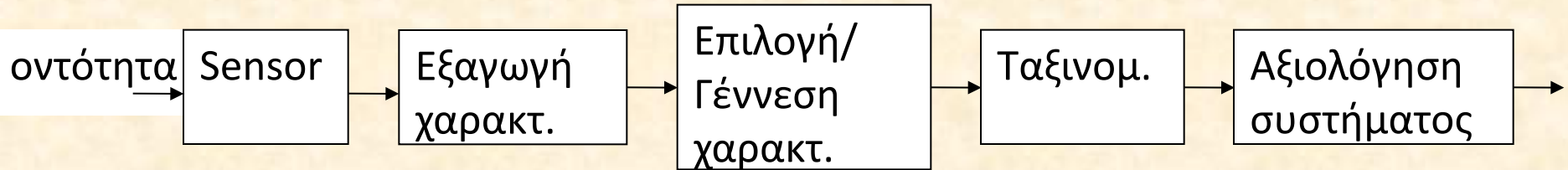
ΑΝΑΓΝΩΡΙΣΗ ΠΡΟΤΥΠΩΝ (PATTERN RECOGNITION)

Σέργιος Θεοδωρίδης

Κωνσταντίνος Κουτρούμπας

Επιλογή / γέννεση χαρακτηριστικών (feature selection / classification)

Τυπική αρχιτεκτονική συστήματος ταξινόμησης



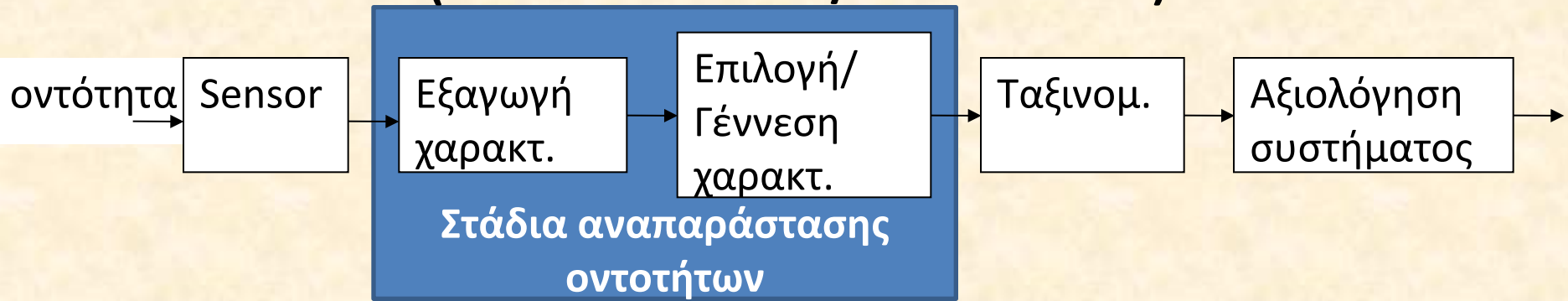
Επιβλεπόμενη ταξινόμηση (Supervised Classification):

Δεδομένα: Ένα πλήθος κλάσεων $\omega_1, \dots, \omega_M$ στις οποίες τα στοιχεία ενός **συνόλου** “οντοτήτων” θα ταξινομηθούν.

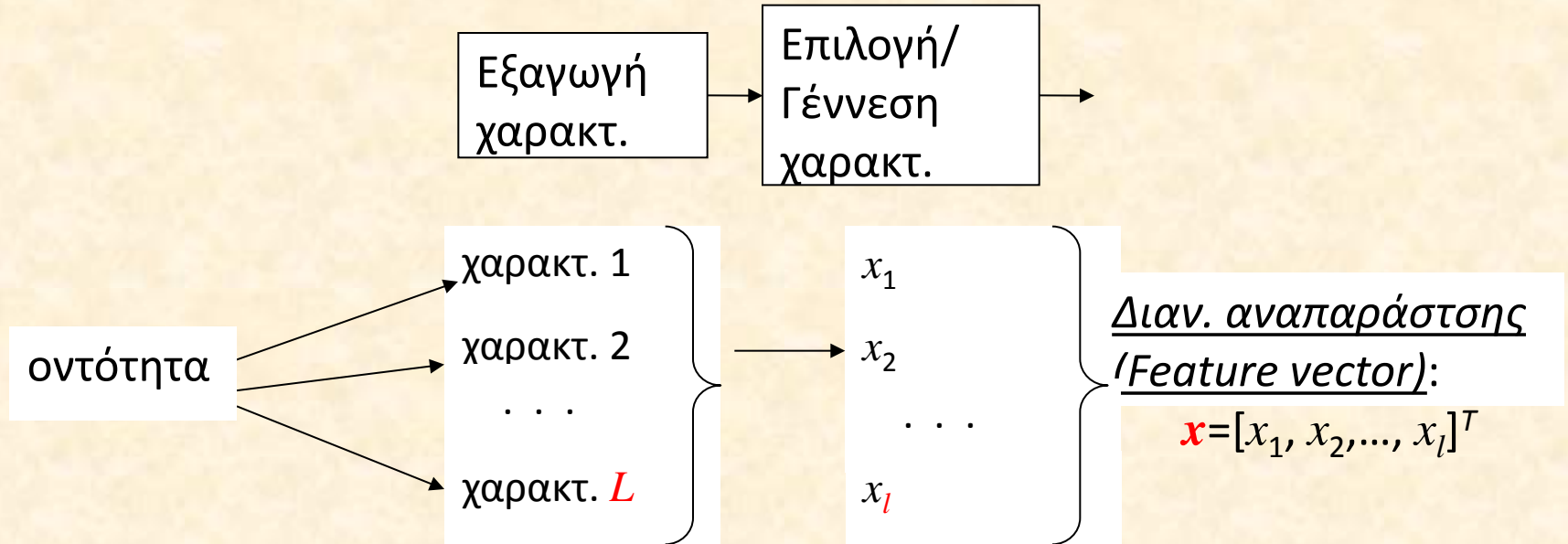
Στόχος:

Δοθείσης μιας **οντότητας**, **καταχώρησέ** την στην “**πιο κατάλληλη**” κλάση (με άλλα λόγια, **εκτίμησε** την “**πιο κατάλληλη**” κλάση γι’ αυτή).

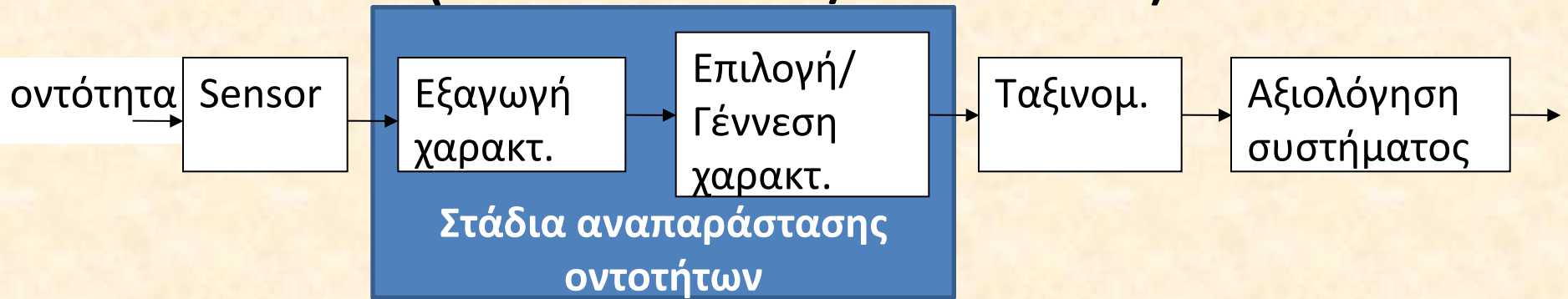
Επιλογή / γέννεση χαρακτηριστικών (feature selection / classification)



Αναπαράσταση οντοτήτων (πώς δουλεύει)



Επιλογή / γέννεση χαρακτηριστικών (feature selection / classification)



Επιλογή χαρακτηριστικών: Ανάμεσα από τα L αρχικά **χαρακτηριστικά** επέλεξε εκείνα που έχουν **αυξημένη διακριτική ικανότητα**.

Γιατί επιλογή χαρακτηριστικών;

Μεγάλες τιμές του L έχουν **τρία μειονεκτήματα:**

- Υψηλές υπολογιστικές απαιτήσεις
- Κακές εκτιμήσεις λάθους (**peaking phenomenon***)
- Απόδοση με **χαμηλή** ικανότητα γενίκευσης

Στόχοι της **επιλογής** χαρακτηριστικών:

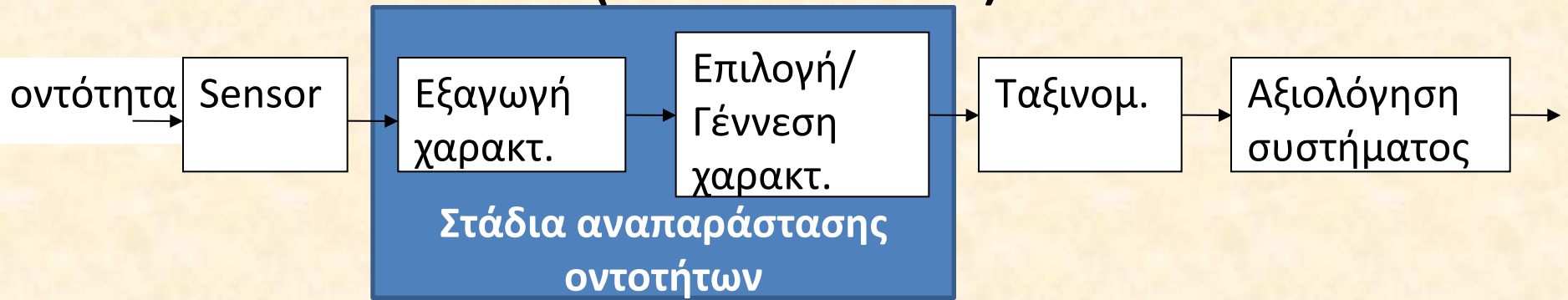
- Ο προσδιορισμός του **“βέλτιστου” αριθμού l** (από L) χαρακτηριστικά.
- Η **επιλογή** των **“βέλτιστων” l** χαρακτηριστικών.

ΣΗΜ.: Σε μερικές περιπτώσεις, οι παραπάνω διαδικασίες γίνονται ταυτόχρονα.

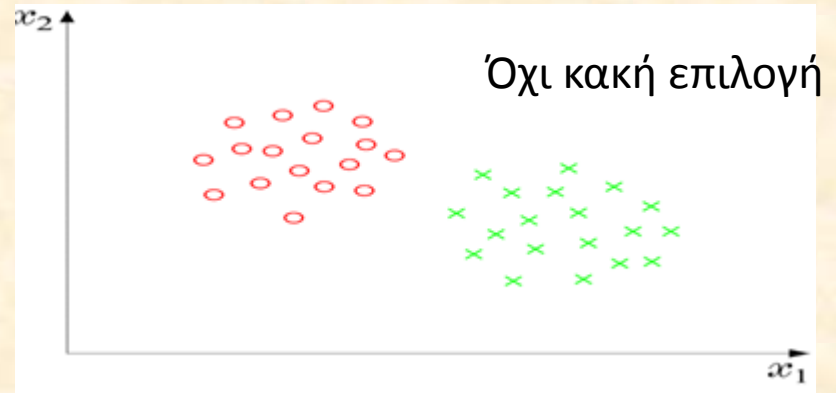
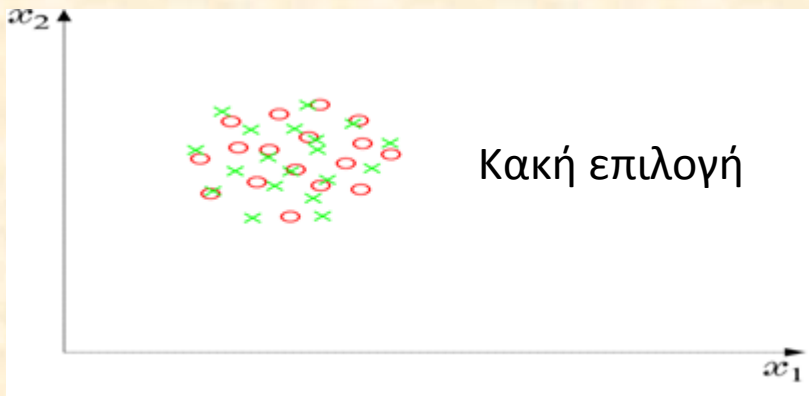
(*) Για δεδομένο N , η πρόσθεση χαρακτ. πάνω από συγκεκριμένο όριο οδηγεί σε **αύξηση** της πιθανότητας σφάλματος

$l < N/3$ είναι μια **λογική επιλογή** σε αρκετές περιπτώσεις.

Επιλογή χαρακτηριστικών (feature selection)



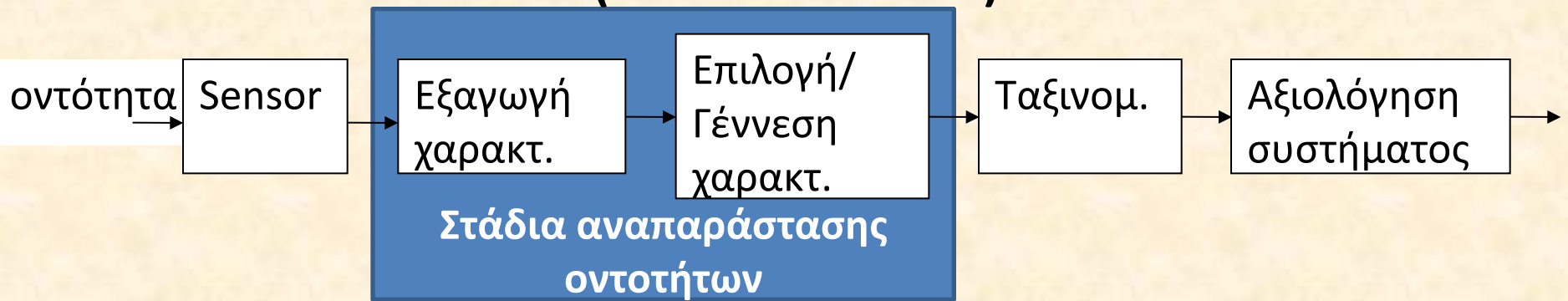
Επιλογή χαρακτηριστικών: Μεταξύ L (αρχικών) **features** επέλεξε εκείνα που έχουν **αυξημένη διακριτική ικανότητα για τις κλάσεις.**



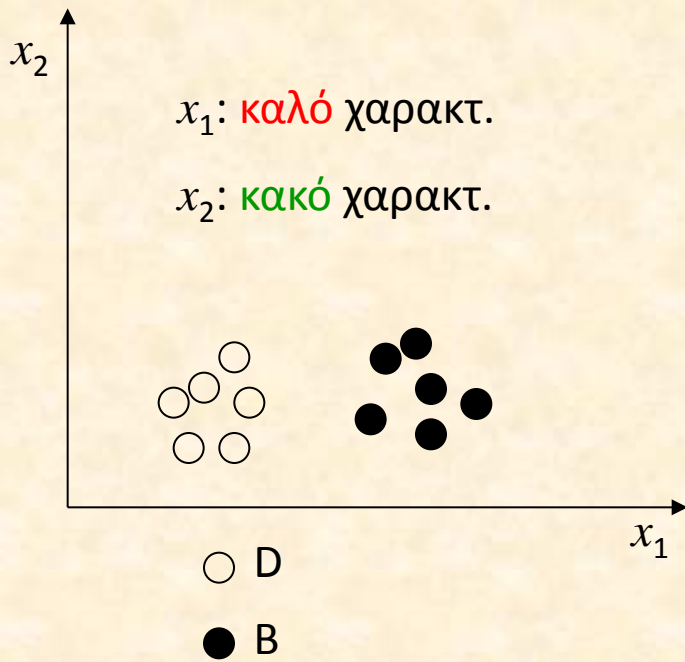
Ο βέλτιστος συνδυασμός χαρακτηριστικών αναπαριστά τις κλάσεις ώστε να έχουμε

- Μικρή διακύμανση μέσα στις κλάσεις (συμπαγείς κλάσεις)
- Μεγάλη απόσταση μεταξύ κλάσεων.

Επιλογή χαρακτηριστικών (feature selection)



Επιλογή χαρακτηριστικών: Μεταξύ L (αρχικών) **features επέλεξε** εκείνα που έχουν **αυξημένη διακριτική ικανότητα για τις κλάσεις.**



- Μέθοδοι επιλογής εξαρτώμενες από την εφαρμογή

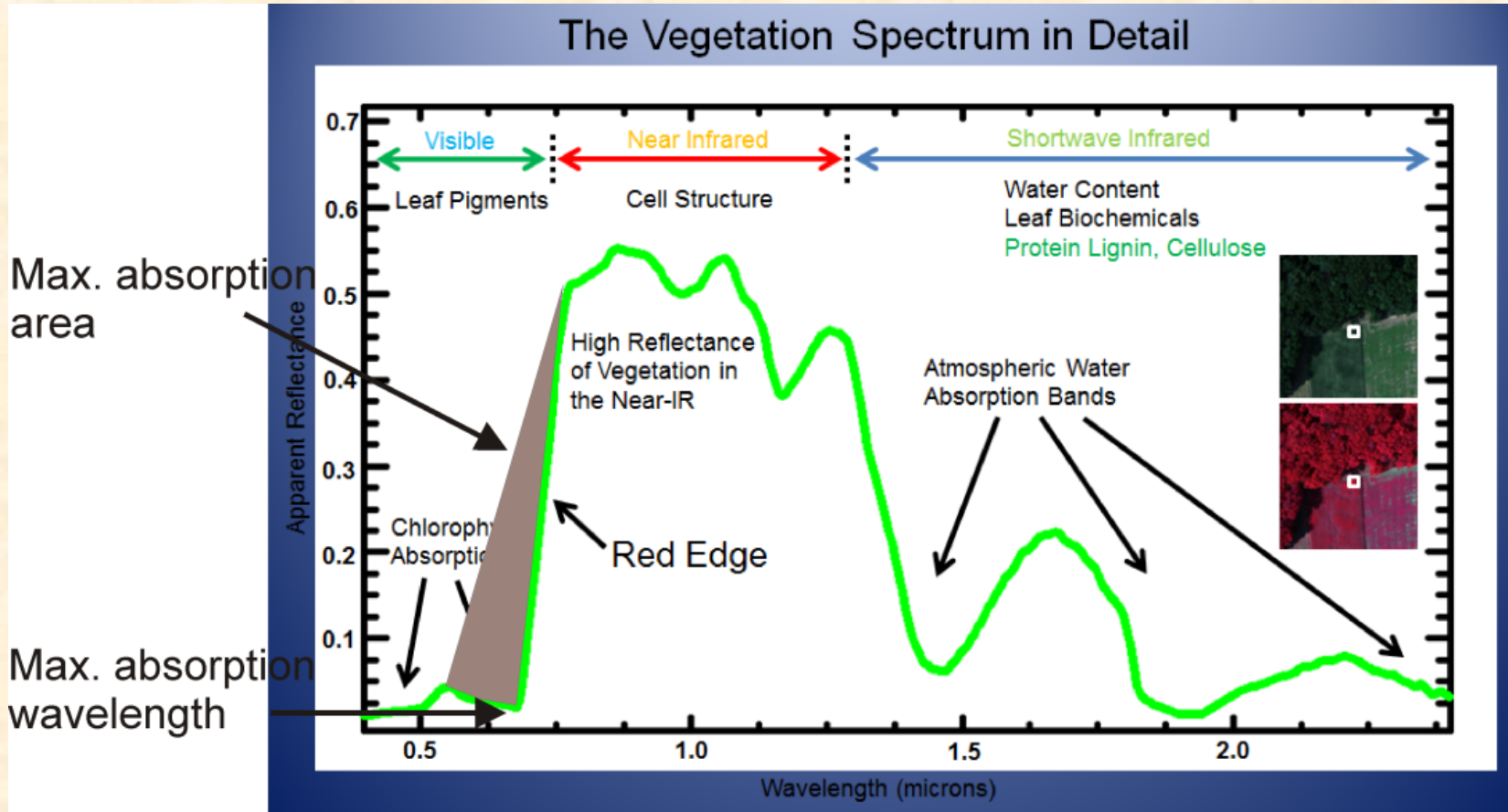
Βασίζονται στη “**φυσική**” του υπό μελέτη θέματος (π.χ., η επιλογή συγκεκριμένου μήκους κύματος για το διαχωρισμό δύο υλικών)

- Μέθοδοι επιλογής μη εξαρτώμενες από την εφαρμογή

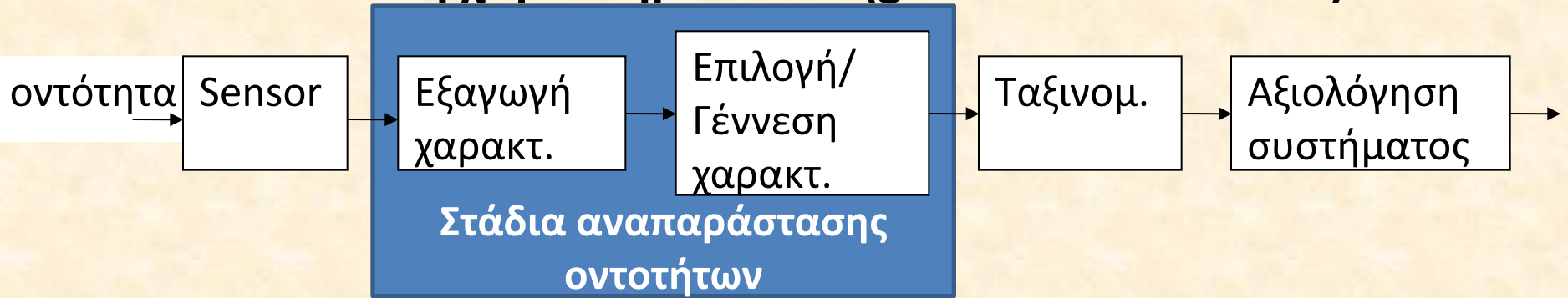
Βασίζονται στη χρήση **μαθηματικών εργαλείων** (στατιστικοί έλεγχοι, μέθοδοι βελτιστοποίησης) **Τα χαρακτ.** Μπορεί να **ελεγχθούν (a) ανεξάρτητα** το ένα από το άλλο ή **(b) σε ομάδες.**

Επιλογή χαρακτηριστικών (feature selection)

Παράδειγμα επιλογής φυσικών χαρακτ.



Γέννηση χαρακτηριστικών (generation selection)



Γέννηση χαρακτηριστικών: «Γέννησε» **νέα χαρακτηριστικά** από ήδη υπάρχοντα.

Στόχος: Η περαιτέρω **μείωση** του **αριθμού των χαρακτηριστικών**.

Δύο δημοφιλείς μέθοδοι:

Ανάλυση κυρίων συνιστωσών (Principal component analysis - PCA):

Ο **αρχικός χώρος** μετασχηματίζεται σε **νέο ορθογώνιο χώρο** όπου τα **χαρακτηριστικά** είναι **μη συσχετισμένα**. Συγκεκριμένα: **κατά μήκος** του (καλούμενου) **1^{ου} κύριου άξονα** διατηρείται η **μέγιστη δυνατή διακύμανση** του συνόλου δεδομένων, **κατά μήκος** του (καλούμενου) **2^{ου} κύριου άξονα** διατηρείται η **μέγιστη δυνατή** (από την **εναπομείνασα**) **διακύμανση** κλπ.

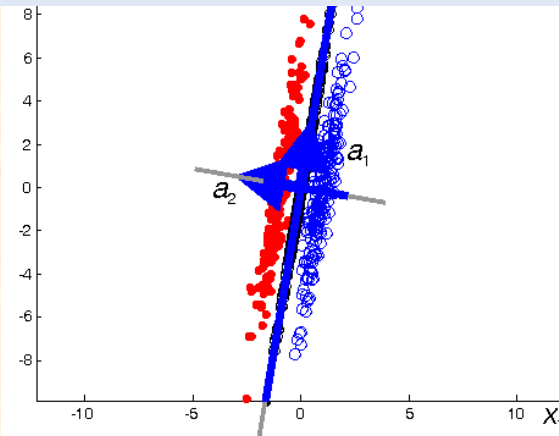
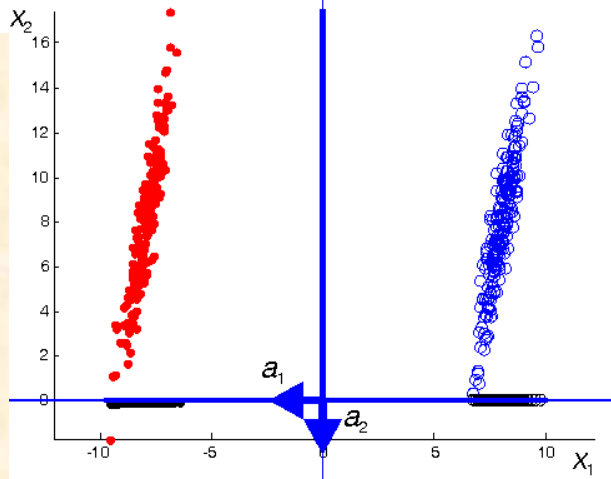
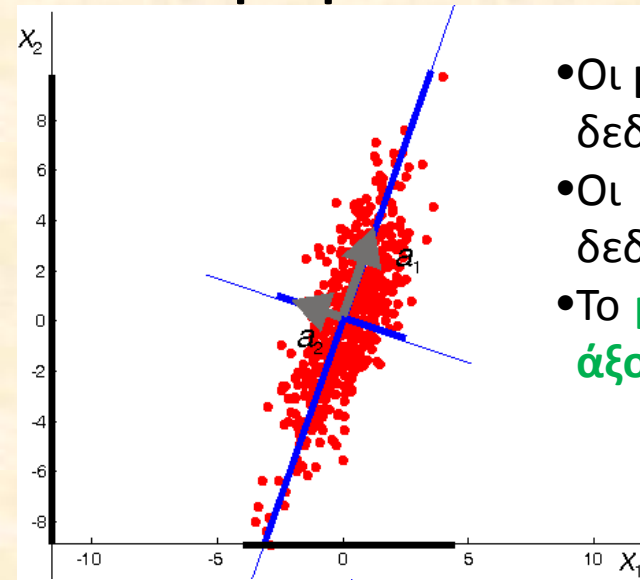
Προβάλλοντας στο χώρο που ορίζεται από τους **«πιο σημαντικούς» κύριους άξονες** επιτυγχάνουμε **μείωση** της **διάστασης**.

Γέννεση χαρακτηριστικών (generation selection)

Ανάλυση κύριων συνιστωσών - Principal Component Analysis - PCA

- Οι **μαύρες γραμμές** δείχνουν το **διάστημα τιμών** των διανυσμάτων δεδομένων κατά μήκος των **αρχικών αξόνων**.
- Οι **μπλε γραμμές** δείχνουν το **διάστημα τιμών** των διανυσμάτων δεδομένων κατά μήκος των **κύριων αξόνων**.
- Το **μέγιστο εύρος** τιμών παρατηρείται κατά μήκος του **1ου κύριου άξονα**.

ΣΗΜ: Η διατήρηση της **μέγιστης δυνατής διασποράς** του συνόλου δεδομένων **ΔΕΝ** εγγυάται απαραίτητα και τη **διατήρηση της διαχωρισιμότητας των κλάσεων**.



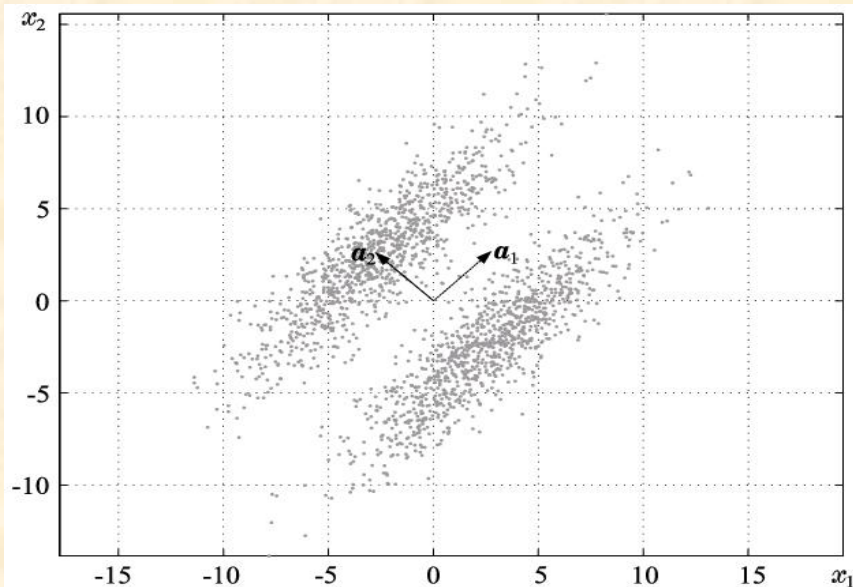
Η **προβολή κατά μήκος του 1ου κύρ. άξονα (a_1)** διατηρεί τη διαχωρισιμότητα των κλάσεων.

Η **προβολή κατά μήκος του 1ου κύρ. άξονα (a_1)** **ΔΕΝ** διατηρεί τη διαχωρισιμότητα των κλάσεων.

Γέννεση χαρακτηριστικών (generation selection)

Ανάλυση ανεξάρτητων συνιστωσών - Independent Component Analysis-ICA)

- Παράγει **στατιστικώς ανεξάρτητα** χαρακτηριστικά
- Δείχνει **προτίμηση** σε **χαρακτ.** των οποίων η κατανομή έχει τη **μικρότερη δυνατή ομοιότητα με την Gaussian pdf** (αυτά είναι πιο πιθανόν να **διατηρήσουν** την πληροφορία **διαχωρισμού των κλάσεων**)



$$K_4(y_2) = -1.7$$

$$K_4(y_1) = 0.1$$

- Κατά μήκος της a_2 η στατιστική είναι **bimodal**.
- Κατά μήκος της a_1 η στατιστική είναι **unimodal** (μοιάζει με **Gaussian**).

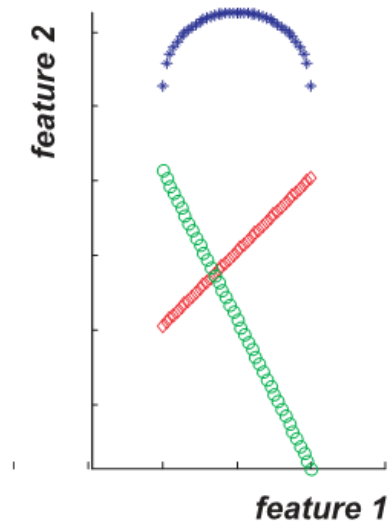
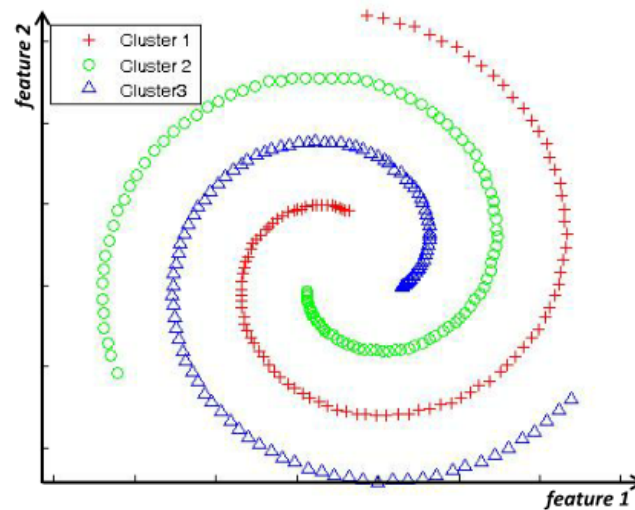
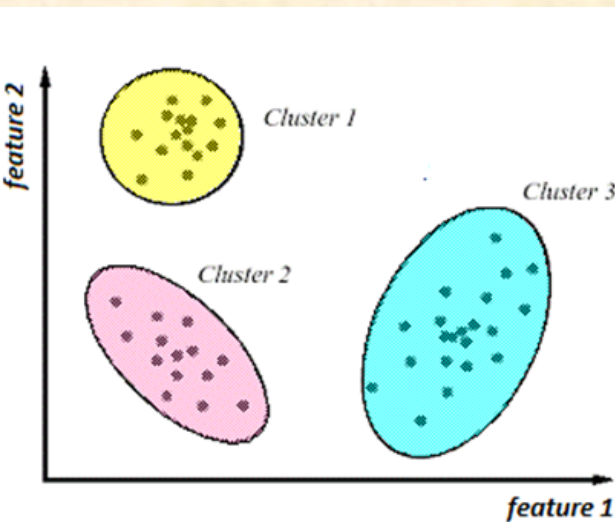
Ομαδοποίηση (Clustering) – Βασικές έννοιες

Ορισμός ομαδοποίησης

Ομαδοποίηση: Η διαδικασία κατά την οποία “**όμοιες**” οντότητες **καταχωρούνται** στο ίδιο σύνολο (“**ομάδα**” - “**cluster**”).

Ορισμός ομάδας: παρατηρήσεις

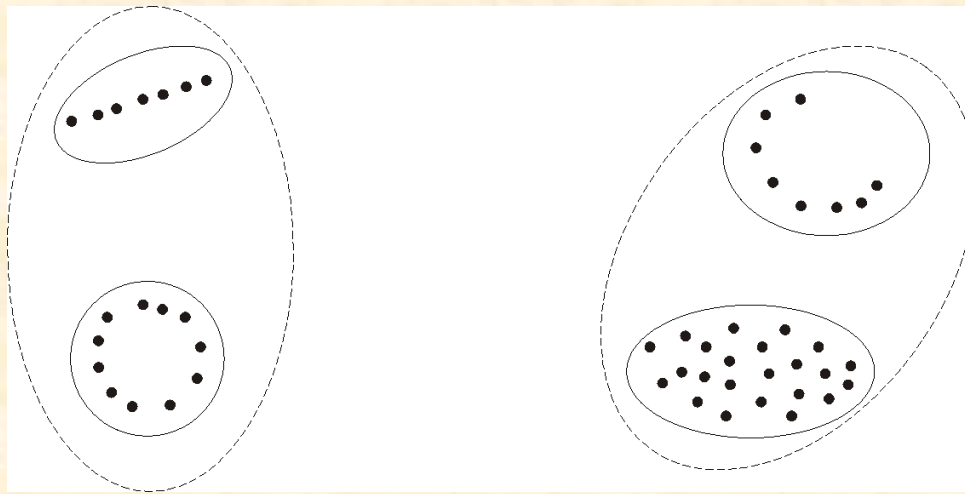
- ✓ Δεν υπάρχει αυστηρός ορισμός για την έννοια της **ομάδας**.
- ✓ Ωστόσο, εννοούμε συνήθως μια **συσσώρευση σημείων** γύρω από :
 - ένα **συγκεκριμένο σημείο** του χώρου χαρακτηριστικών (στην περίπτωση αυτή η ομάδα μοντελοποιείται συνήθως από **κανονική** κατανομή).
 - ένα **manifold** (π.χ. υπερεπίπεδο, υπερσφαίρα) στο χώρο χαρακτηριστικών.



Ομαδοποίηση (Clustering) – Βασικές έννοιες

Ομαδοποίηση είναι η διαδικασία της ταυτοποίησης συσσωρεύσεων σημείων σ' έναν I -διάστατο χώρο.

1^η Παρατήρηση: Για την ομαδοποίηση, η **υποκειμενικότητα** είναι μια αναπόφευκτη πραγματικότητα..



Ποιος είναι ο «σωστός» αριθμός ομάδων;;
2 ή 4;;

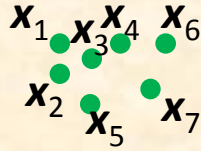
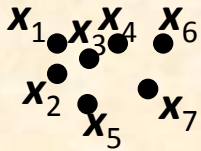
2^η παρατήρηση: Υπάρχει **τουλάχιστον μία “σφαιρική” παράμετρος** σε κάθε αλγόριθμο ομαδοποίησης (π.χ. ο αριθμός των ομάδων) που χρειάζεται να **καθορισθεί από το χρήστη** (**ελλιπώς ορισμένο (ill-posed)** πρόβλημα).

Ομαδοποίηση (Clustering) – Βασικές έννοιες

Ομαδοποίηση είναι η διαδικασία της ταυτοποίησης συσσωρεύσεων σημείων σ' έναν l -διάστατο χώρο.

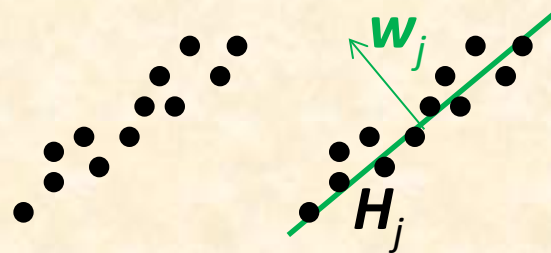
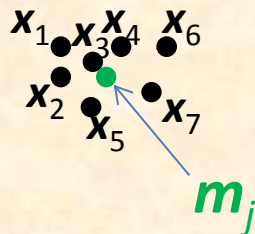
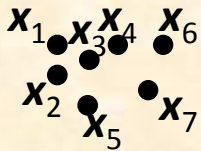
Αναπαράσταση ομάδας

✓ Μέσω όλων των σημείων της (μη-παραμετρική αναπαράσταση)



$\{x_1, x_2, x_3, x_4, x_5, x_6, x_7\}$

✓ Μέσω ενός συνόλου παραμέτρων (παραμετρική αναπαράσταση)



Παράμετροι:

$$m_j = [m_{j1}, m_{j2}, \dots, m_{jl}]^T$$

$$H_j: \theta_j^T x_i + \theta_{j0} = 0$$

Παράμετροι:

$$\theta_j = [\theta_{j1}, \theta_{j2}, \dots, \theta_{jl}]^T, \theta_{j0}$$

Ομαδοποίηση (Clustering) – Βασικές έννοιες

Ομαδοποίηση: Η έννοια της ομοιότητας

Μια σημαντική παρατήρηση: Η έννοια της **εγγύτητας (proximity)** μεταξύ **διανυσμάτων** (ή μεταξύ **ομάδων** σε κάποιες περιπτώσεις) πρέπει να ποσοτικοποιηθεί.

Μέτρα ομοιότητας μεταξύ x_i και x_j : Όσο **υψηλότερη** είναι η τιμή τους, τόσο **πιο όμοια** είναι τα x_i και x_j .

Μέτρα ανομοιότητας μεταξύ x_i και x_j : Όσο **υψηλότερη** είναι η τιμή τους, τόσο **λιγότερο όμοια** είναι τα x_i και x_j .

Σε αρκετές περιπτώσεις, χρειάζεται επίσης ο ορισμός ενός **κατωφλίου** για τον **ορισμό** των εννοιών “**όμοια**”, “**ανόμοια**”.

Παραδείγματα:

- (a) Μεταξύ **διανυσμάτων**: **τετραγωνική Ευκλείδ. απόσταση** $d(x_i, x_j) = \|x_i - x_j\|^2$
- (b) Μεταξύ συνόλων (**ομάδων**): $d(C_i, C_j) = \min_{x \in C_i, y \in C_j} d(x, y)$ (ή **max** ή **mean**)

ΣΗΜ: Η επιλογή του **μέτρου εγγύτητας** μπορεί να **επηρεάσει** το αποτέλεσμα ομαδοποίησης.

Ομαδοποίηση (Clustering) – Βασικές έννοιες

Σύμβαση:

Ομαδοποίηση

\mathcal{C}

\leftrightarrow

Σύνολο ομάδων

=

$\{C_1, C_2, \dots, C_m\}$

Ιεραρχία ομαδοποιήσεων

\mathcal{H}

\leftrightarrow

Σύνολο ομαδοποιήσεων

=

$\{\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_N\}$

Παράδειγμα:

Σημεία δεδομένων:

x_1, x_2, x_3, x_4, x_5

Ομάδες:

$C_1 = \{x_1, x_3\}, C_2 = \{x_2, x_4, x_5\}$

Ομαδοποίηση:

$\mathcal{C} = \{C_1, C_2\}$

Ιεραρχία ομαδοποιήσεων:

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

Κατηγοριοποίηση αλγορίθμων ομαδοποίησης

Ιεραρχικοί αλγόριθμοι:

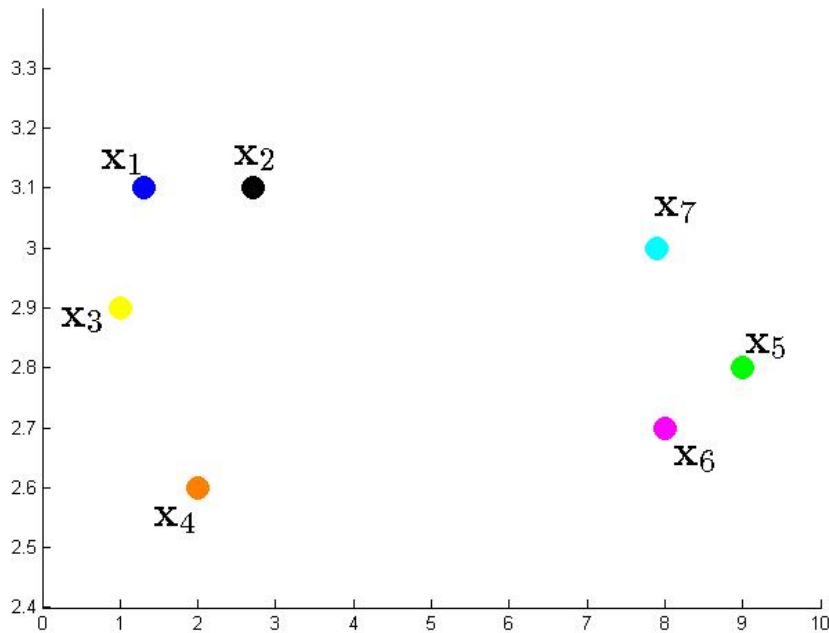
- Παράγουν μια **ιεραρχία ομαδοποιήσεων των δεδομένων**.
- Μπορεί να είναι **είτε παραμετρικοί είτε μη παραμετρικοί**.
- Κύριες φιλοσοφίες: **συσσωρευτική (agglomerative)**, **διαιρετική (divisive)**.

Μη ιεραρχικοί αλγόριθμοι:

- Παράγουν **μια ομαδοποίηση των δεδομένων**.
- Βασίζονται σε ποικιλία ιδεών που προέρχονται από
 - Το πλαίσιο της **βελτιστοποίησης συναρτήσεων κόστους (Cost function optimization - CFO) framework** (κυρίως **παραμετρικοί αλγόριθμοι**)
 - ✓ **Πιθανοτικό** πλαίσιο (**παραμετρικοί αλγόριθμοι**)
 - ✓ **Μη πιθανοτικό** πλαίσιο (**παραμετρικοί αλγόριθμοι**)
 - Το Πλαίσιο **θεωρίας γράφων** (π.χ. spectral clustering, ελάχιστο δένδρο κάλυψης) (κυρίως **μη παραμετρικοί αλγόριθμοι**)
 - **Πυκνότητα ομάδων (Cluster-density)** (**μη παραμετρικοί αλγόριθμοι**)

Ιεραρχικοί αλγόριθμοι ομαδοποίησης

x_1 x_2 x_3 x_4 x_5 x_6 x_7



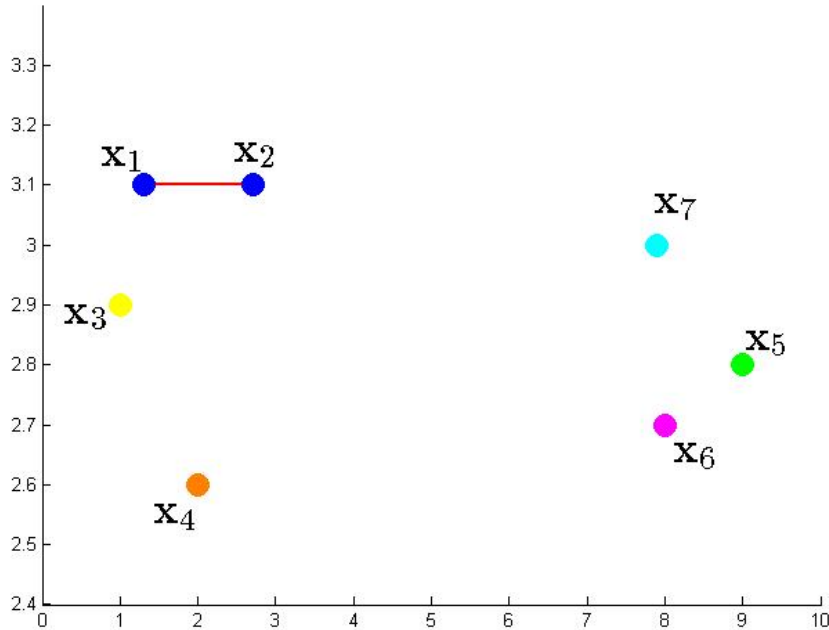
Συσσωρευτική φιλοσοφία (Agglomerative philosophy):

- Στην **αρχική ομαδοποίηση** όλα τα **διανύσματα δεδομένων ανήκουν** σε **διαφορετικές ομάδες**.
- Σε κάθε βήμα ορίζεται μια **νέα ομαδοποίηση συνενώνοντας** σε ένα τις **δύο όμοιες μεταξύ τους ομάδες**.
- Στην **τελική ομαδοποίηση** όλα τα **διανύσματα ανήκουν** στην **ίδια ομάδα**.

Ιεραρχικοί αλγόριθμοι ομαδοποίησης

x_1 x_2 x_3 x_4 x_5 x_6 x_7

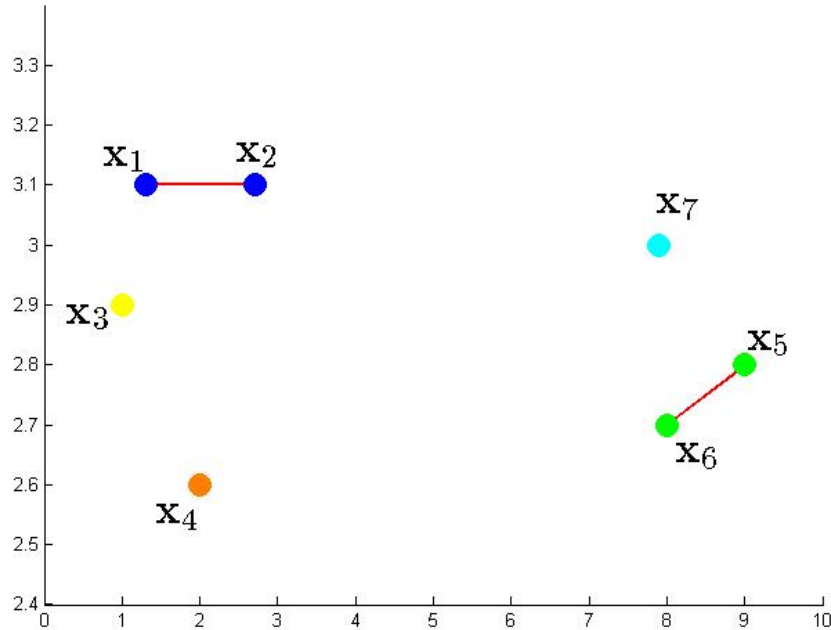
1.4



Συσσωρευτική φιλοσοφία (Agglomerative philosophy):

- Στην **αρχική ομαδοποίηση** όλα τα **διανύσματα δεδομένων ανήκουν** σε **διαφορετικές ομάδες**.
- Σε κάθε βήμα ορίζεται μια **νέα ομαδοποίηση συνενώνοντας** σε ένα τις **δύο όμοιες μεταξύ τους ομάδες**.
- Στην **τελική ομαδοποίηση** όλα τα **διανύσματα ανήκουν** στην **ίδια ομάδα**.

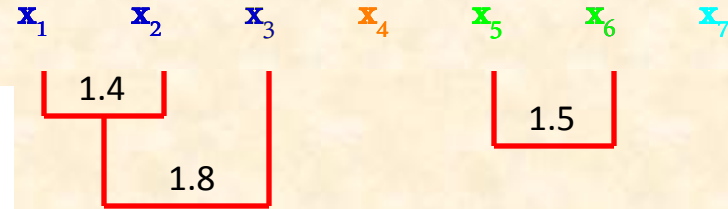
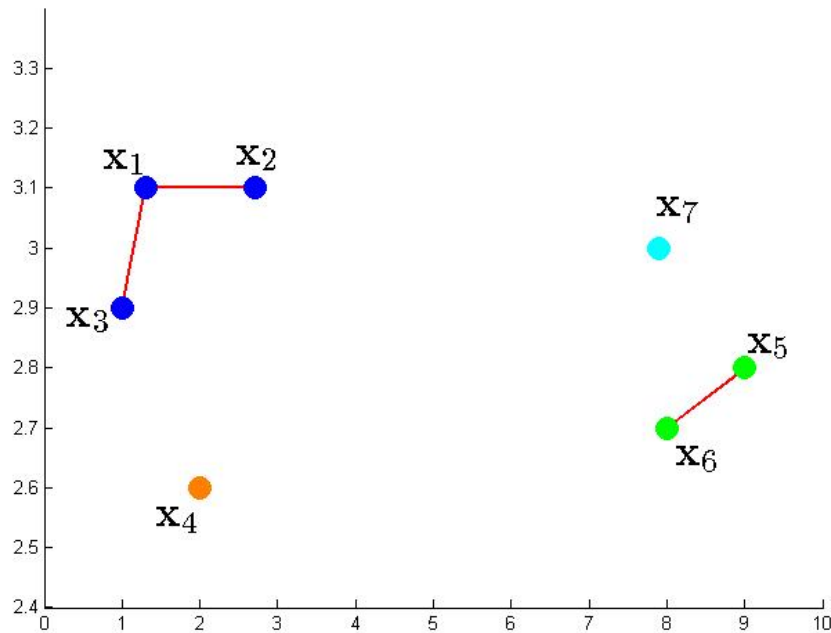
Ιεραρχικοί αλγόριθμοι ομαδοποίησης



Συσσωρευτική φιλοσοφία (Agglomerative philosophy):

- Στην **αρχική ομαδοποίηση** όλα τα **διανύσματα δεδομένων ανήκουν** σε **διαφορετικές ομάδες**.
- Σε κάθε βήμα ορίζεται μια **νέα ομαδοποίηση συνενώνοντας** σε ένα τις **δύο όμοιες μεταξύ τους ομάδες**.
- Στην **τελική ομαδοποίηση** όλα τα **διανύσματα ανήκουν** στην **ίδια ομάδα**.

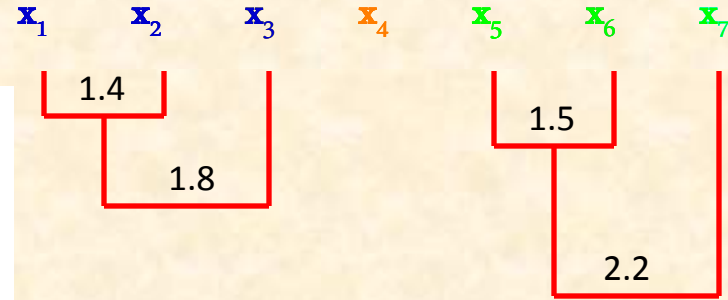
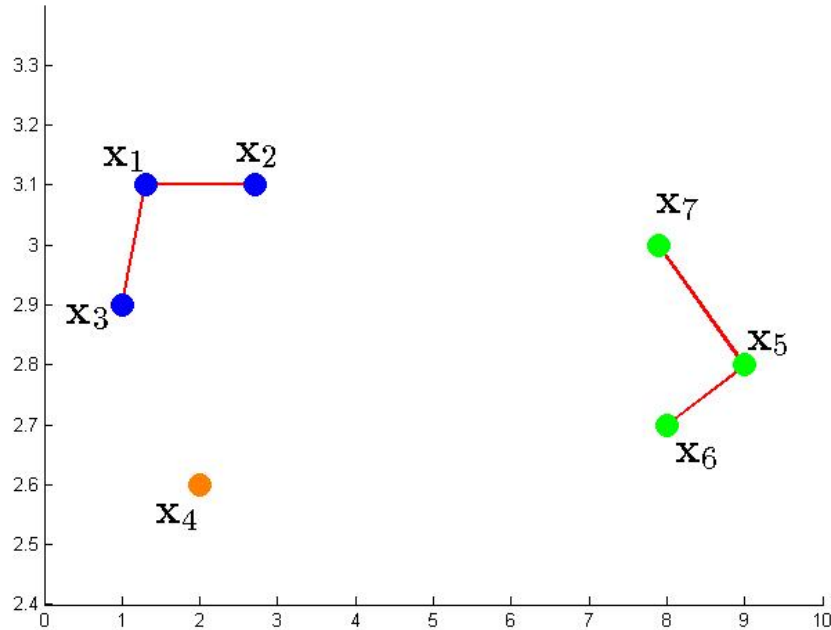
Ιεραρχικοί αλγόριθμοι ομαδοποίησης



Συσσωρευτική φιλοσοφία (Agglomerative philosophy):

- Στην αρχική ομαδοποίηση όλα τα διανύσματα δεδομένων ανήκουν σε διαφορετικές ομάδες.
- Σε κάθε βήμα ορίζεται μια νέα ομαδοποίηση συνενώνοντας σε ένα τις δύο όμοιες μεταξύ τους ομάδες.
- Στην τελική ομαδοποίηση όλα τα διανύσματα ανήκουν στην ίδια ομάδα.

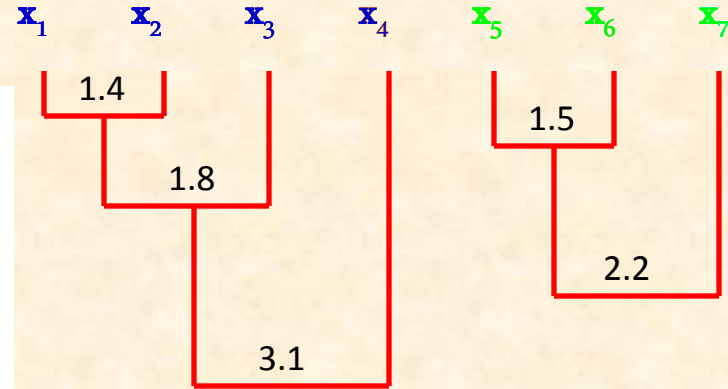
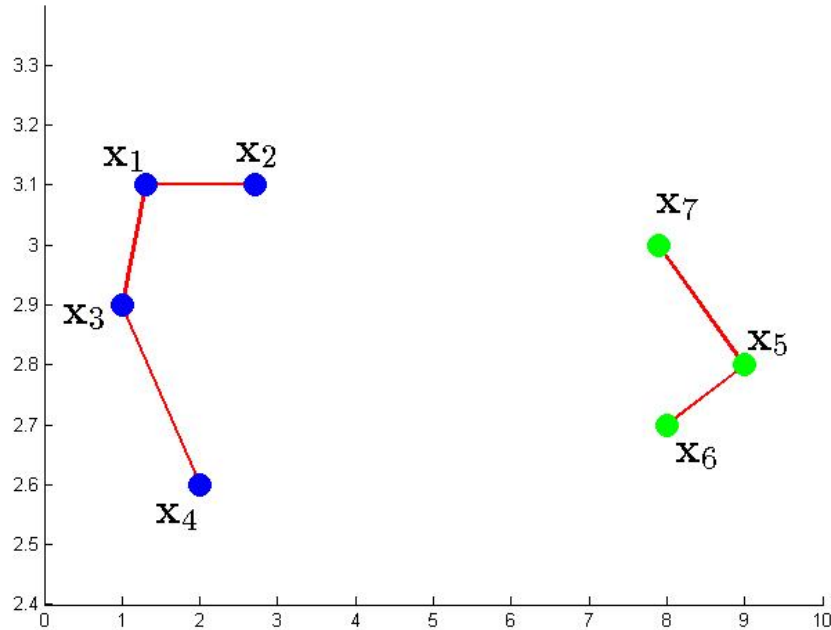
Ιεραρχικοί αλγόριθμοι ομαδοποίησης



Συσσωρευτική φιλοσοφία (Agglomerative philosophy):

- Στην αρχική ομαδοποίηση όλα τα διανύσματα δεδομένων ανήκουν σε διαφορετικές ομάδες.
- Σε κάθε βήμα ορίζεται μια νέα ομαδοποίηση συνενώνοντας σε ένα τις δύο όμοιες μεταξύ τους ομάδες.
- Στην τελική ομαδοποίηση όλα τα διανύσματα ανήκουν στην ίδια ομάδα.

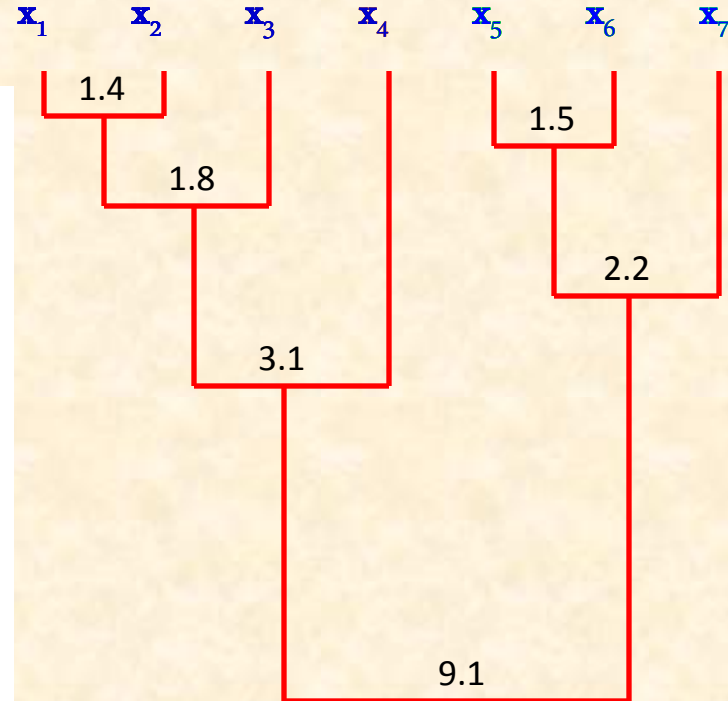
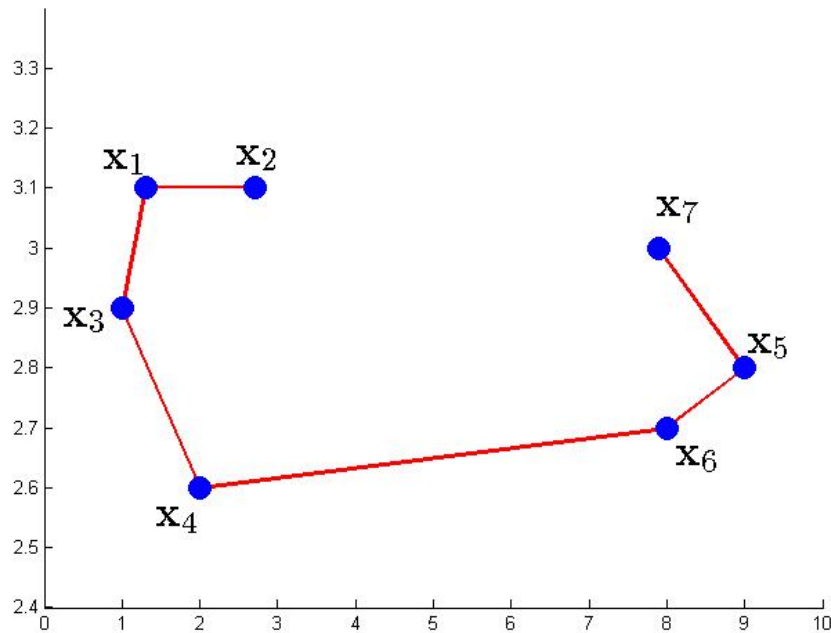
Ιεραρχικοί αλγόριθμοι ομαδοποίησης



Συσσωρευτική φιλοσοφία (Agglomerative philosophy):

- Στην αρχική ομαδοποίηση όλα τα διανύσματα δεδομένων ανήκουν σε διαφορετικές ομάδες.
- Σε κάθε βήμα ορίζεται μια νέα ομαδοποίηση συνενώνοντας σε ένα τις δύο όμοιες μεταξύ τους ομάδες.
- Στην τελική ομαδοποίηση όλα τα διανύσματα ανήκουν στην ίδια ομάδα.

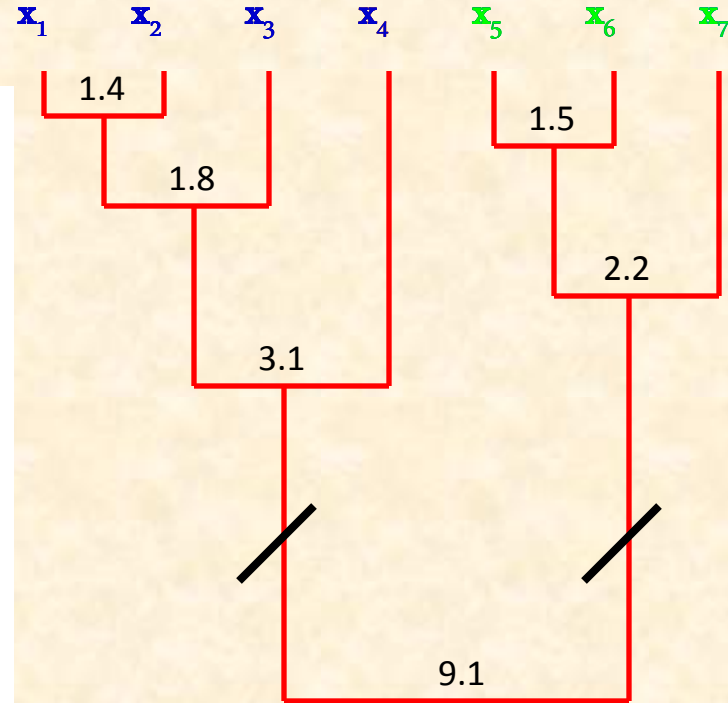
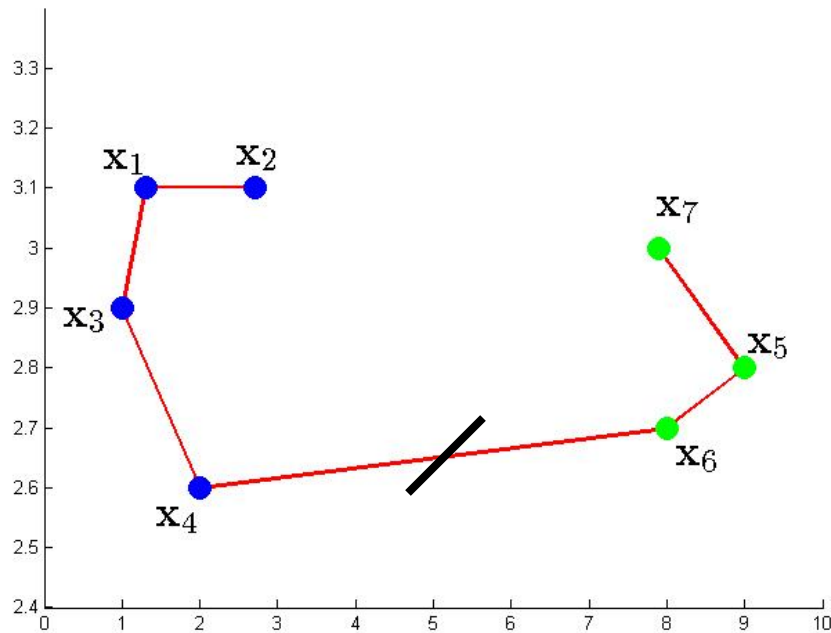
Ιεραρχικοί αλγόριθμοι ομαδοποίησης



Συσσωρευτική φιλοσοφία (Agglomerative philosophy):

- Στην αρχική ομαδοποίηση όλα τα διανύσματα δεδομένων ανήκουν σε διαφορετικές ομάδες.
- Σε κάθε βήμα ορίζεται μια νέα ομαδοποίηση συνενώνοντας σε ένα τις δύο όμοιες μεταξύ τους ομάδες.
- Στην τελική ομαδοποίηση όλα τα διανύσματα ανήκουν στην ίδια ομάδα.

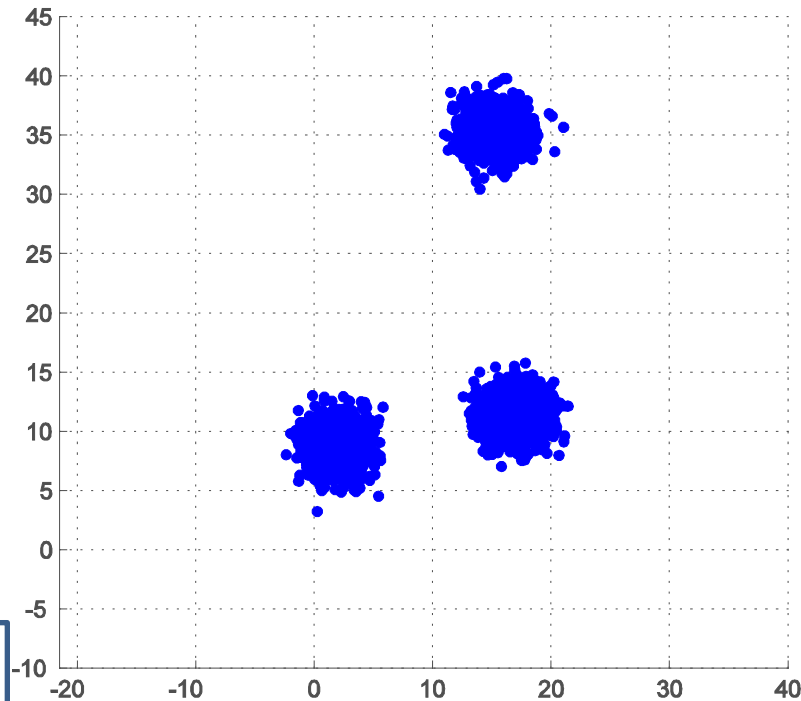
Ιεραρχικοί αλγόριθμοι ομαδοποίησης



Συσσωρευτική φιλοσοφία (Agglomerative philosophy):

- Στην αρχική ομαδοποίηση όλα τα διανύσματα δεδομένων ανήκουν σε διαφορετικές ομάδες.
- Σε κάθε βήμα ορίζεται μια νέα ομαδοποίηση συνενώνοντας σε ένα τις δύο όμοιες μεταξύ τους ομάδες.
- Στην τελική ομαδοποίηση όλα τα διανύσματα ανήκουν στην ίδια ομάδα.

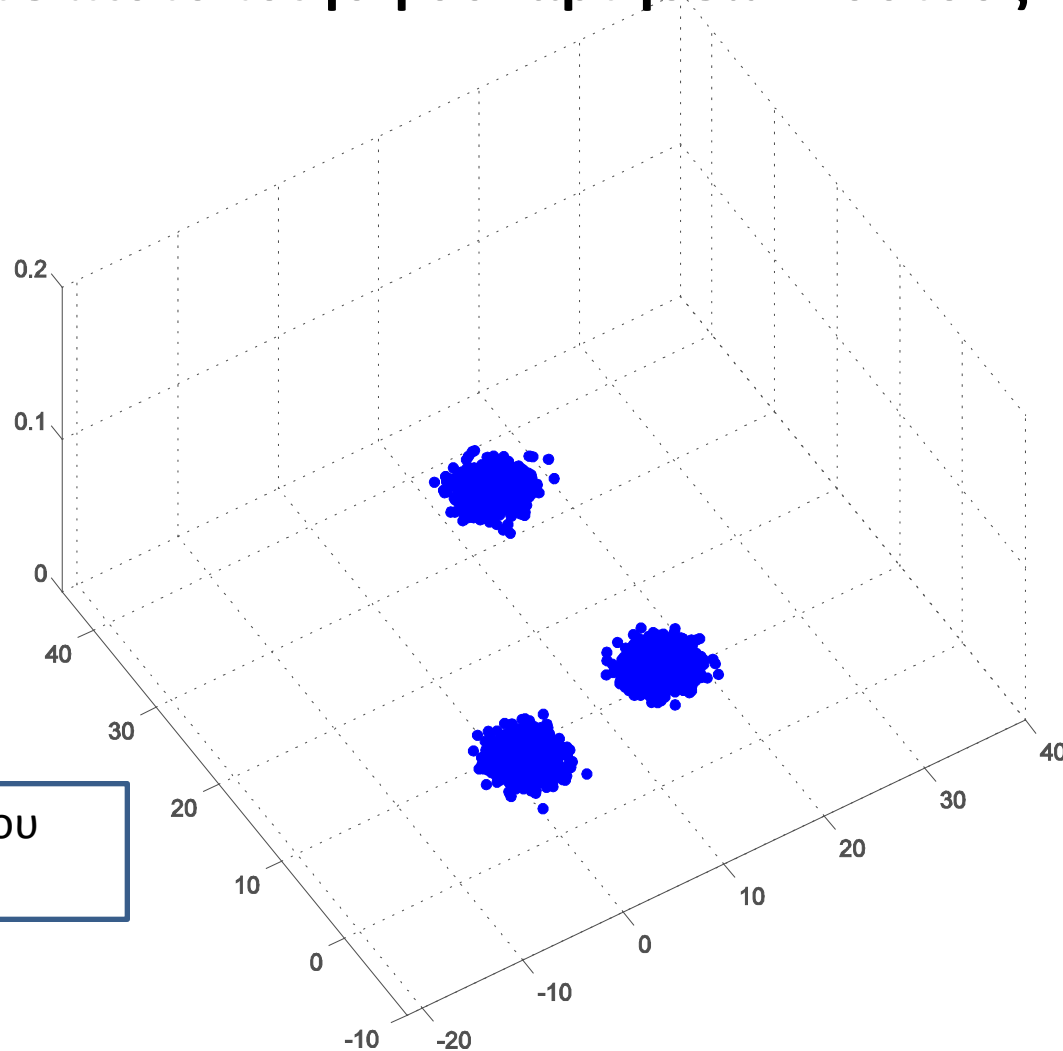
Πιθανοτικοί αλγόριθμοι ομαδοποίησης βασιζόμενοι στην βελτιστοποίηση συναρτήσεων κόστους



Προαπαίτηση: Γνώση του αριθμού των ομάδων.

- Υιοθέτησε μια παραμετρική μίξη κατανομών, κάθε μία από τις οποίες αντιστοιχεί σε μια ομάδα (π.χ., μίξη Gaussians), και οι παράμετροί της αρχικοποιούνται τυχαία.
- Μετακίνησε προοδευτικά κάθε κατανομή πάνω από μια ομάδα, μέσω της βελτιστοποίησης ενός κριτηρίου.

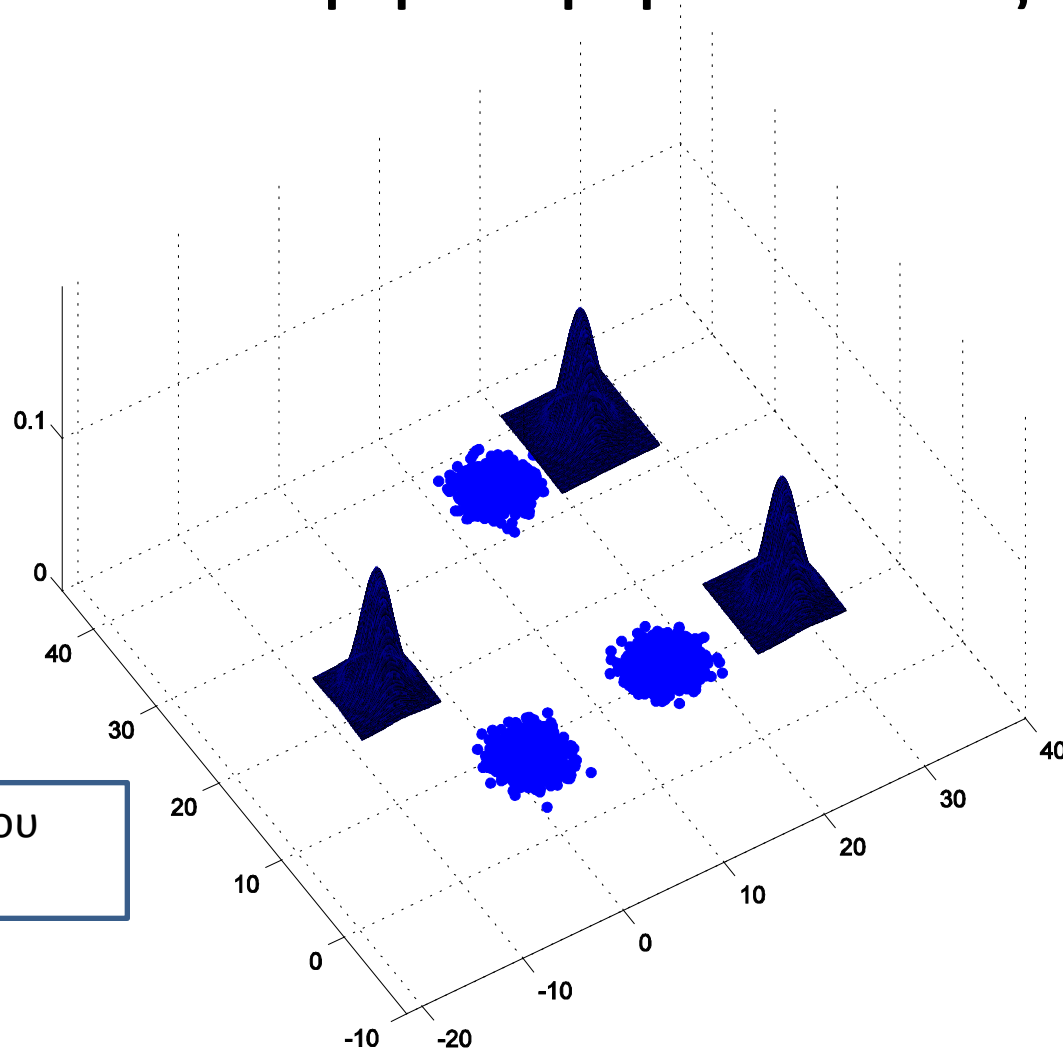
Πιθανοτικοί αλγόριθμοι ομαδοποίησης βασιζόμενοι στην βελτιστοποίηση συναρτήσεων κόστους



Προαπαίτηση: Γνώση του αριθμού των ομάδων.

- Υιοθέτησε μια παραμετρική μίξη κατανομών, κάθε μία από τις οποίες αντιστοιχεί σε μια ομάδα (π.χ., μίξη Gaussians), και οι παράμετροί της αρχικοποιούνται τυχαία.
- Μετακίνησε προοδευτικά κάθε κατανομή πάνω από μια ομάδα, μέσω της βελτιστοποίησης ενός κριτηρίου.

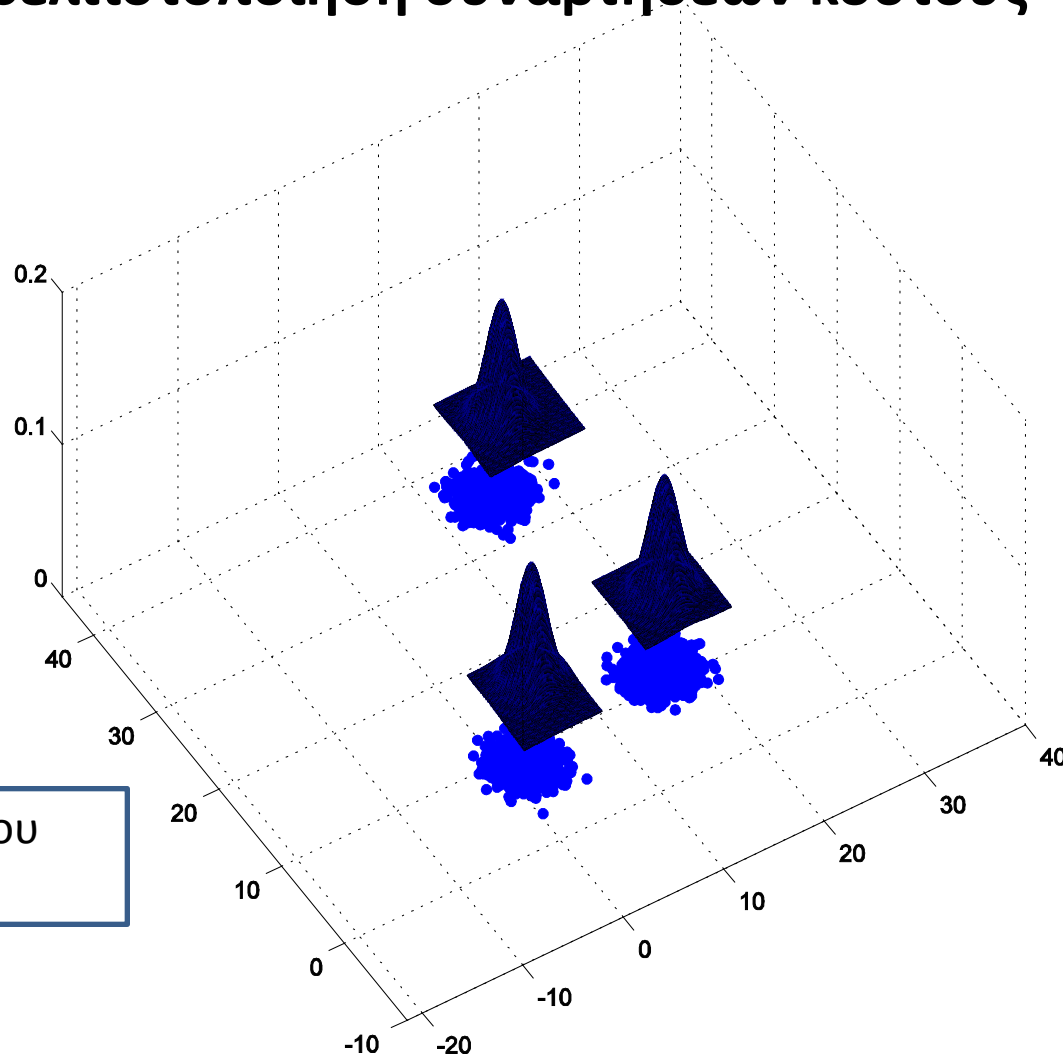
Πιθανοτικοί αλγόριθμοι ομαδοποίησης βασιζόμενοι στην βελτιστοποίηση συναρτήσεων κόστους



Προαπαίτηση: Γνώση του αριθμού των ομάδων.

- Υιοθέτησε μια παραμετρική μίξη κατανομών, κάθε μία από τις οποίες αντιστοιχεί σε μια ομάδα (π.χ., μίξη Gaussians), και οι παράμετροί της αρχικοποιούνται τυχαία.
- Μετακίνησε προοδευτικά κάθε κατανομή πάνω από μια ομάδα, μέσω της βελτιστοποίησης ενός κριτηρίου.

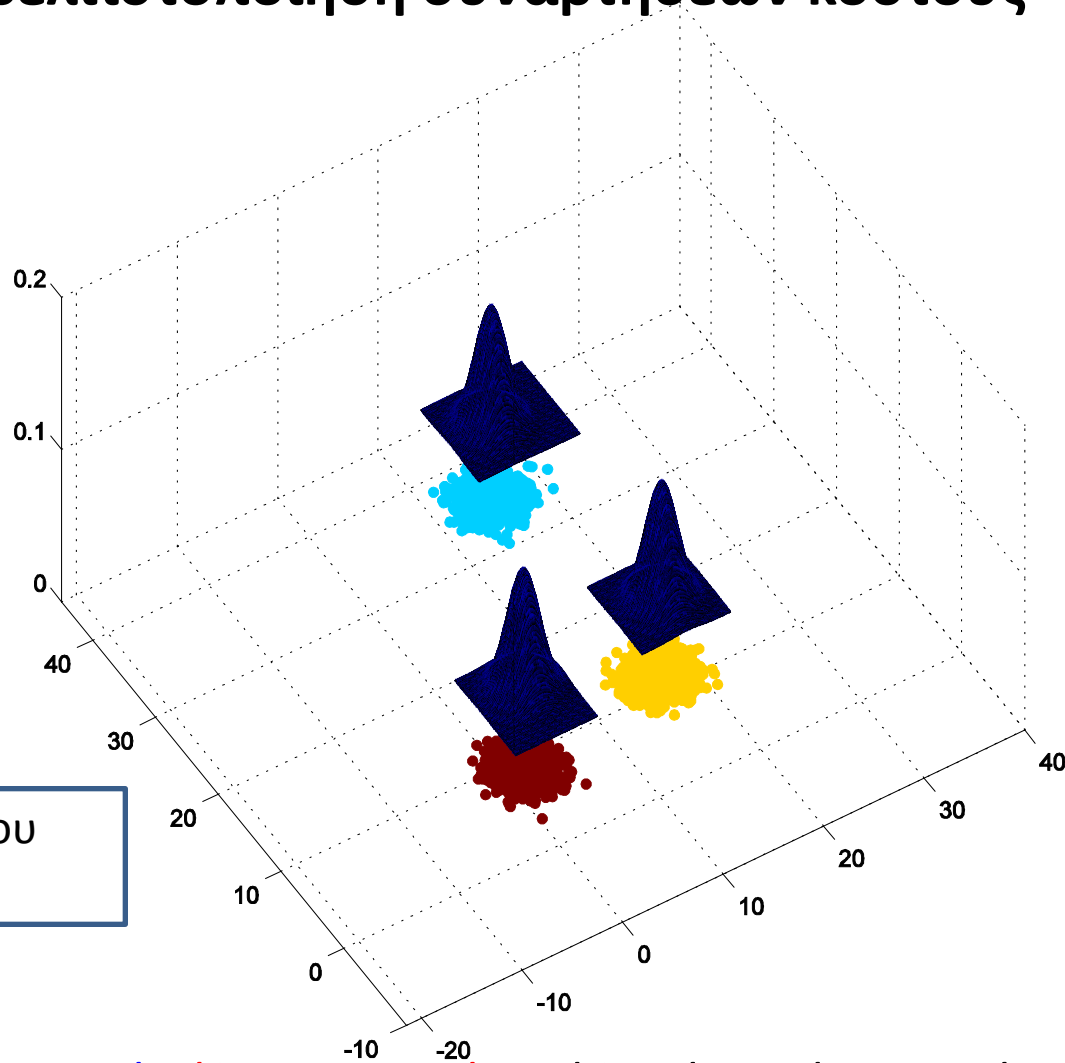
Πιθανοτικοί αλγόριθμοι ομαδοποίησης βασιζόμενοι στην βελτιστοποίηση συναρτήσεων κόστους



Προαπαίτηση: Γνώση του αριθμού των ομάδων.

- Υιοθέτησε μια παραμετρική μίξη κατανομών, κάθε μία από τις οποίες αντιστοιχεί σε μια ομάδα (π.χ., μίξη Gaussians), και οι παράμετροί της αρχικοποιούνται τυχαία.
- Μετακίνησε προοδευτικά κάθε κατανομή πάνω από μια ομάδα, μέσω της βελτιστοποίησης ενός κριτηρίου.

Πιθανοτικοί αλγόριθμοι ομαδοποίησης βασιζόμενοι στην βελτιστοποίηση συναρτήσεων κόστους



Προαπαίτηση: Γνώση του αριθμού των ομάδων.

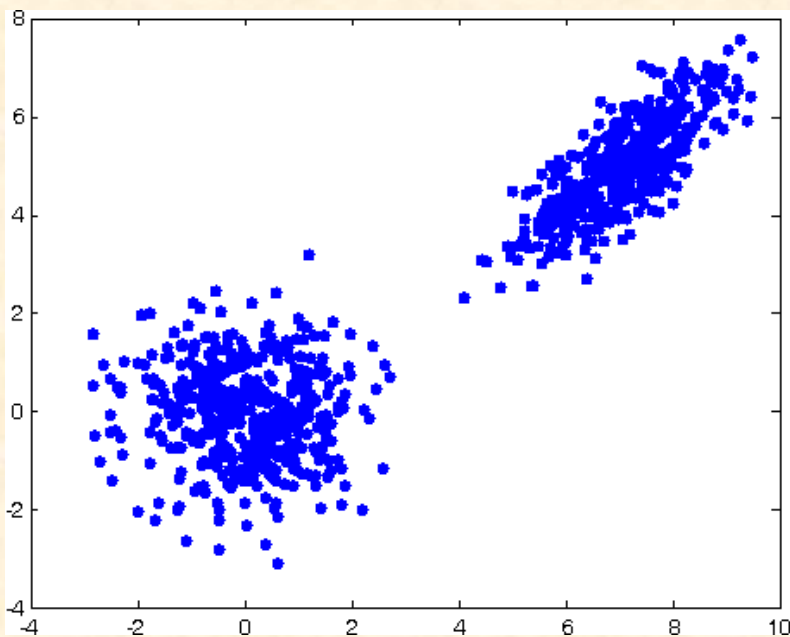
- Υιοθέτησε μια παραμετρική μίξη κατανομών, κάθε μία από τις οποίες αντιστοιχεί σε μια ομάδα (π.χ., μίξη Gaussians), και οι παράμετροί της αρχικοποιούνται τυχαία.
- Μετακίνησε προοδευτικά κάθε κατανομή πάνω από μια ομάδα, μέσω της βελτιστοποίησης ενός κριτηρίου.

Πιθανοτικοί αλγόριθμοι ομαδοποίησης βασιζόμενοι στην βελτιστοποίηση συναρτήσεων κόστους

Ο αλγόριθμος Αναμονής-Μεγιστοποίησης (Expectation – Maximization - EM) (μοντέλα μίξης - mixture models)

Μοντέλο μίξης: Σταθμισμένο άθροισμα pdfs γνωστής παραμετρικής μορφής

$$p(x) = \sum_{j=1}^m P_j p(x | j), \quad \sum_{j=1}^m P_j = 1, \quad \int_{-\infty}^{+\infty} p(x | j) = 1$$



Παρατηρήσεις:

- Απαιτείται εκ των προτέρων γνώση του αριθμού των ομάδων, m .
- Κάθε μία pdf αντιστοιχεί και σε μία φυσική ομάδα.

Συμβολισμός:

$$p(x | j) = p(x | j; \theta_j)$$

$$\Theta = (\theta_1, \dots, \theta_m), \quad P = [P_1, \dots, P_m]^T$$

$X = \{x_1, \dots, x_N\}$: **ελλιπή** (παρατηρούμενα) δεδομένα.

$X^c = \{(x_1, j_1), \dots, (x_N, j_N)\}$: **πλήρη** δεδομένα

Πιθανοτικοί αλγόριθμοι ομαδοποίησης βασιζόμενοι στην βελτιστοποίηση συναρτήσεων κόστους

Ο **EM** αλγόριθμος (μίξη μοντέλων)

Στόχος:

Εκτίμηση των Θ και P μέσω **μεγιστοποίησης** της **αναμενόμενης τιμής** της συνάρτησης **log-likelihood** **δεσμευμένης** στα **παρατηρούμενα δεδομένα**.

$$\ln p(X; \Theta, P) = \sum_{n=1}^N \sum_{j=1}^m P(j | \mathbf{x}_n; \theta_j) \ln p(\mathbf{x}_n, j; \theta_j) =$$

$$\sum_{n=1}^N \sum_{j=1}^m P(j | \mathbf{x}_n; \theta_j) \ln \left(p(\mathbf{x}_n | j; \theta_j) P_j \right)$$

Μίξη κανονικών κατανομών

$$p(\mathbf{x}) = \sum_{j=1}^m P_j p(\mathbf{x} | j; \boldsymbol{\mu}_j, \Sigma_j),$$

$$p(\mathbf{x} | j; \boldsymbol{\mu}_j, \Sigma_j) = \frac{1}{(2\pi)^{l/2} |\Sigma_j|^{1/2}} \exp\left(-0.5 \cdot (\mathbf{x} - \boldsymbol{\mu}_j)^T \Sigma_j^{-1} (\mathbf{x} - \boldsymbol{\mu}_j)\right)$$

Πιθανοτικοί αλγόριθμοι ομαδοποίησης βασιζόμενοι στην βελτιστοποίηση συναρτήσεων κόστους

Ο **EM** αλγόριθμος για την περίπτωση της **μίξης κανονικών κατανομών**.

- Αρχικοποίηση: $\mu_j = \mu_j^{(0)}$, $\Sigma_j = \Sigma_j^{(0)}$, $P = P^{(0)}$

- $t=0$

- Επανάλαβε

$$P(j | \mathbf{x}_n; \Theta^{(t)}, P^{(t)}) = \frac{p(\mathbf{x}_n | j; \theta_k^{(t)}) P_j^{(t)}}{\sum_{q=1}^m p(\mathbf{x}_n | q; \theta_q^{(t)}) P_q^{(t)}} \equiv \gamma_{jn}^{(t)}$$

Expectation
step

➤ $t=t+1$

Maximization
step

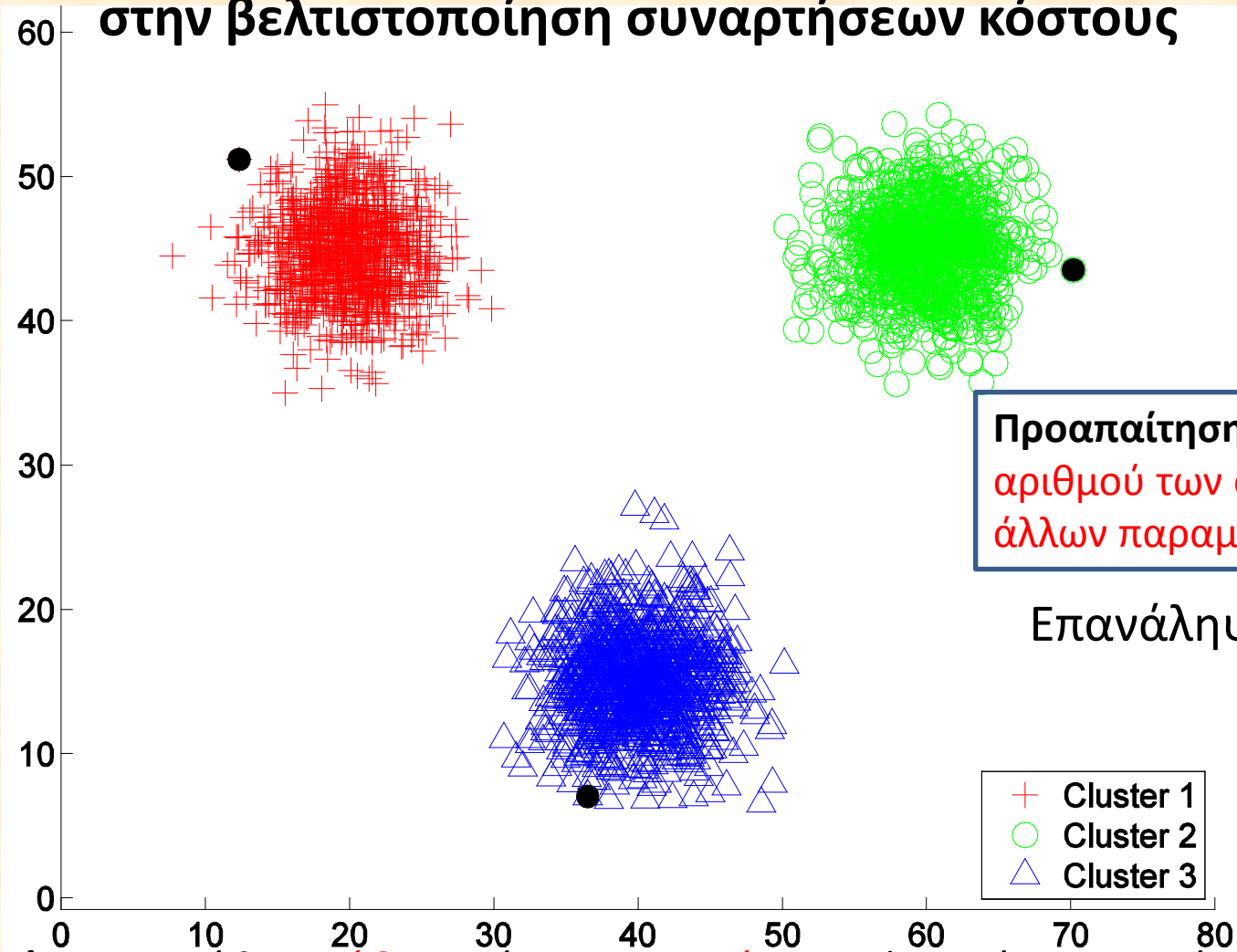
$$\mu_j^{(t+1)} = \frac{\sum_{n=1}^N \gamma_{jn}^{(t)} \mathbf{x}_n}{\sum_{n=1}^N \gamma_{jn}^{(t)}}$$

$$\Sigma_j^{(t+1)} = \frac{\sum_{n=1}^N \gamma_{jn}^{(t)} (\mathbf{x}_n - \mu_j^{(t+1)}) (\mathbf{x}_n - \mu_j^{(t+1)})^T}{\sum_{n=1}^N \gamma_{jn}^{(t)}}$$

$$P_j^{(t+1)} = \frac{1}{N} \sum_{n=1}^N \gamma_{jn}^{(t)}$$

- Έως ότου ικανοποιηθεί ένα κατάλληλο κριτήριο

Μη πιθανοτικοί αλγόριθμοι ομαδοποίησης βασιζόμενοι στην βελτιστοποίηση συναρτήσεων κόστους

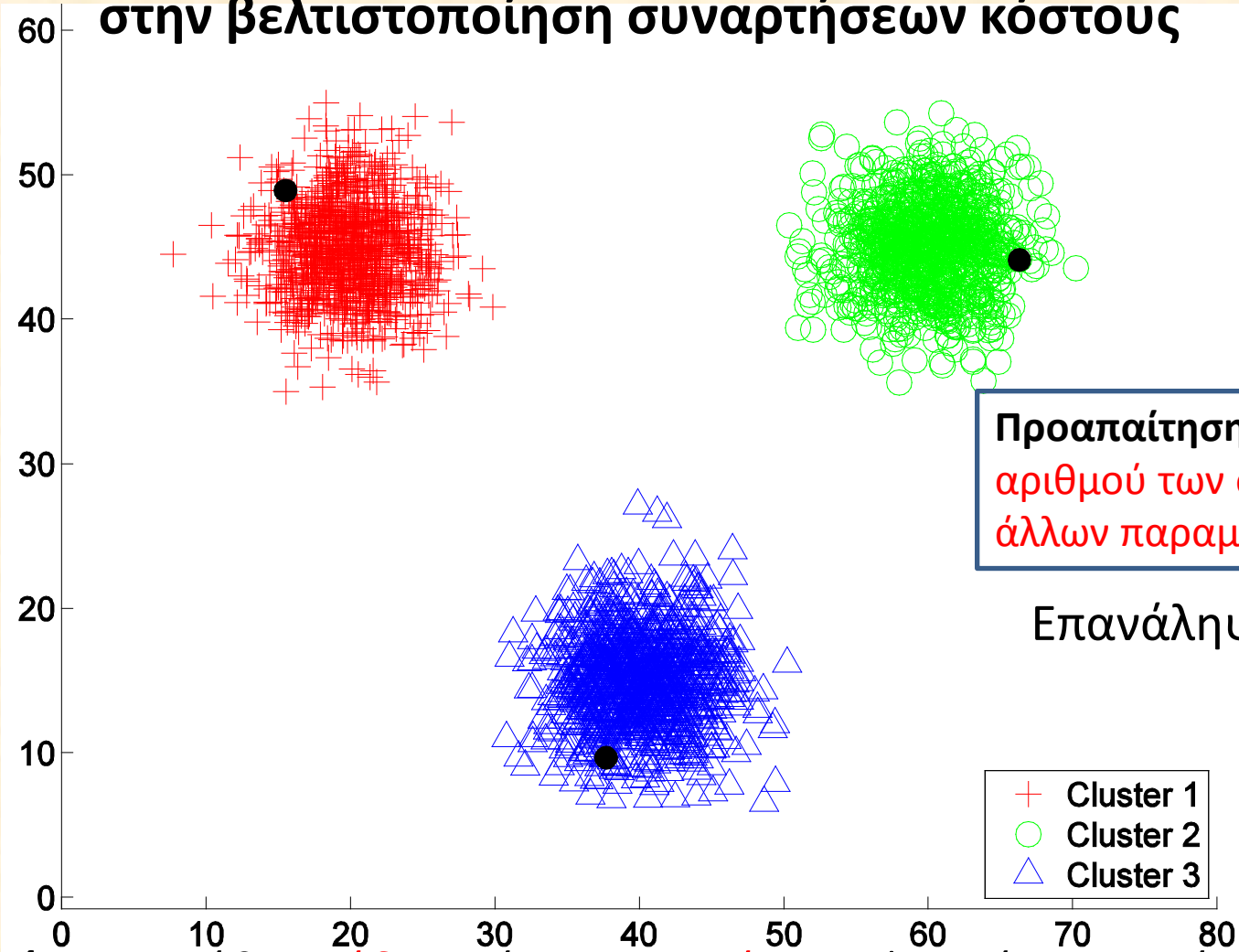


Προαπαίτηση: Γνώση του αριθμού των ομάδων ή άλλων παραμέτρων.

Επανάληψη 1

- Αντιπροσώπευσε κάθε ομάδα με έναν αντιπρόσωπο (σημείο, γραμμή, κλπ)
- Αρχικοποίησε τους αντιπροσώπους σε τυχαίες θέσεις
- Μετακίνησε προοδευτικά τους αντιπροσώπους προς τις φυσικές ομάδες που σχηματίζουν τα δεδομένα, βελτιστοποιώντας ένα κριτήριο.

Μη πιθανοτικοί αλγόριθμοι ομαδοποίησης βασιζόμενοι στην βελτιστοποίηση συναρτήσεων κόστους

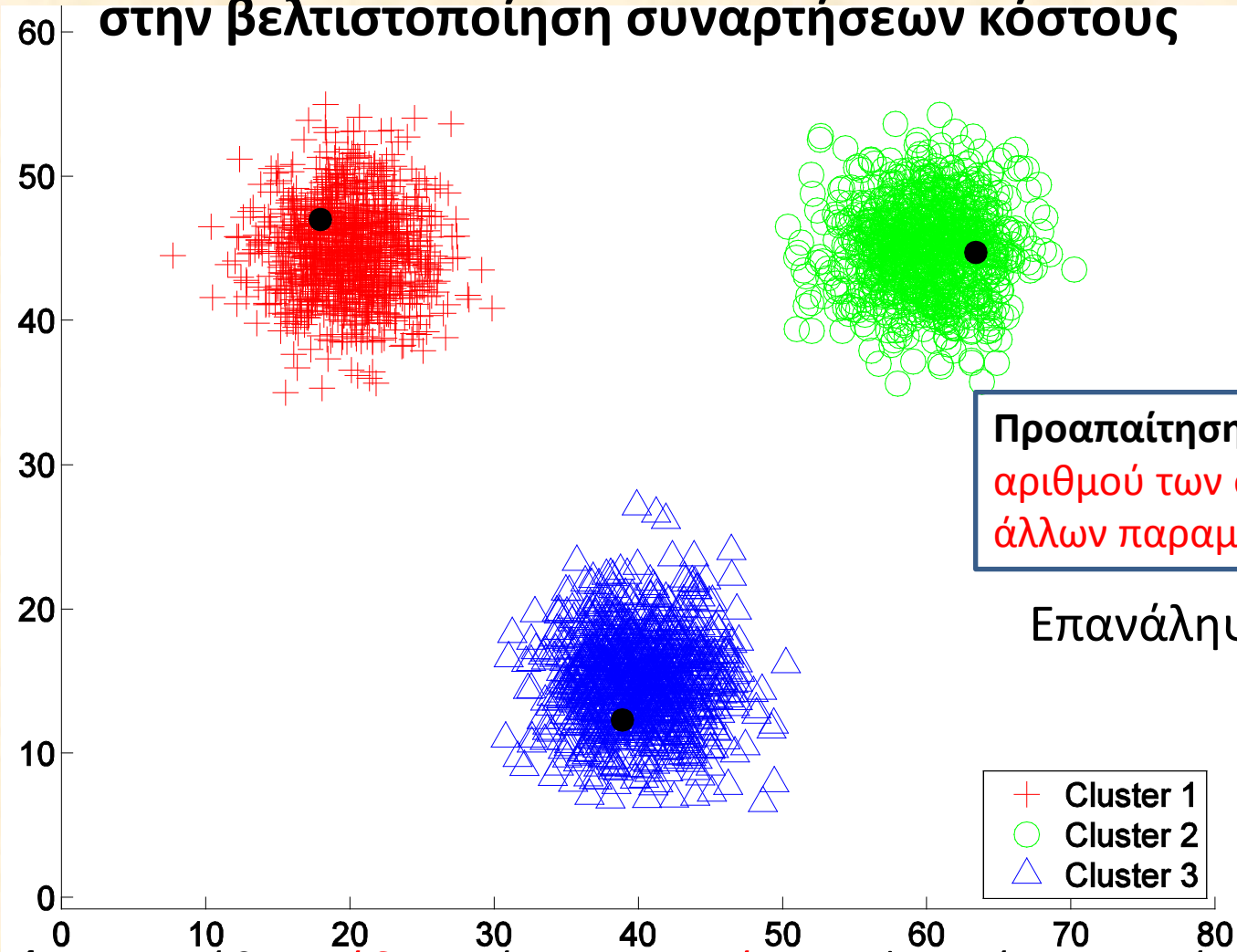


Προαπαίτηση: Γνώση του αριθμού των ομάδων ή άλλων παραμέτρων.

Επανάληψη 2

- Αντιπροσώπευσε κάθε ομάδα με έναν αντιπρόσωπο (σημείο, γραμμή, κλπ)
- Αρχικοποίησε τους αντιπροσώπους σε τυχαίες θέσεις
- Μετακίνησε προοδευτικά τους αντιπροσώπους προς τις φυσικές ομάδες που σχηματίζουν τα δεδομένα, βελτιστοποιώντας ένα κριτήριο.

Μη πιθανοτικοί αλγόριθμοι ομαδοποίησης βασιζόμενοι στην βελτιστοποίηση συναρτήσεων κόστους

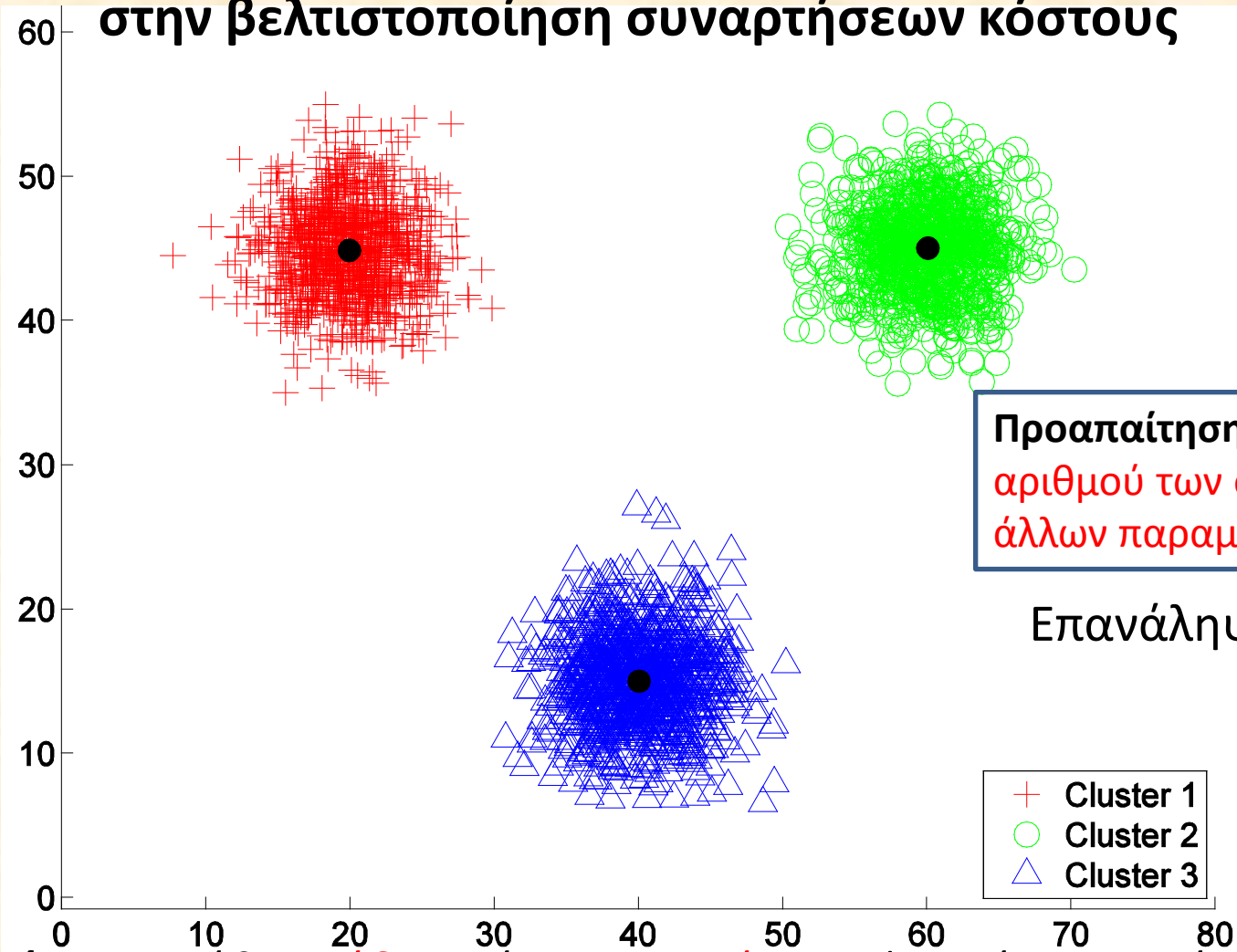


Προαπαίτηση: Γνώση του αριθμού των ομάδων ή άλλων παραμέτρων.

Επανάληψη 3

- Αντιπροσώπευσε κάθε ομάδα με έναν αντιπρόσωπο (σημείο, γραμμή, κλπ)
- Αρχικοποίησε τους αντιπροσώπους σε τυχαίες θέσεις
- Μετακίνησε προοδευτικά τους αντιπροσώπους προς τις φυσικές ομάδες που σχηματίζουν τα δεδομένα, βελτιστοποιώντας ένα κριτήριο.

Μη πιθανοτικοί αλγόριθμοι ομαδοποίησης βασιζόμενοι στην βελτιστοποίηση συναρτήσεων κόστους



Προαπαίτηση: Γνώση του αριθμού των ομάδων ή άλλων παραμέτρων.

Επανάληψη 4 (τελική)

- Αντιπροσώπευσε κάθε ομάδα με έναν αντιπρόσωπο (σημείο, γραμμή, κλπ)
- Αρχικοποίησε τους αντιπροσώπους σε τυχαίες θέσεις
- Μετακίνησε προοδευτικά τους αντιπροσώπους προς τις φυσικές ομάδες που σχηματίζουν τα δεδομένα, βελτιστοποιώντας ένα κριτήριο.

Μη πιθανοτικοί αλγόριθμοι ομαδοποίησης βασιζόμενοι στην βελτιστοποίηση συναρτήσεων κόστους – Ο **αλγ. k-means**

- Πραγματοποιεί **επεξεργασία κατά δέσμες** και επιστρέφει **μια ομαδοποίηση**
- Είναι αλγόριθμος **αυστηρής (hard) ομαδοποίησης**, που χρησιμοποιεί **σημεία αντιπροσώπων** (θ_j) για την αντιπροσώπευση των ομάδων (C_j).
- Προκύπτει από τη βελτιστοποίηση της ακόλουθης συνάρτησης κόστους

$$J(\Theta, U) = \sum_{i=1}^N \sum_{j=1}^m u_{ij} \| \mathbf{x}_i - \theta_j \|^2$$

N : αριθμός σημείων
 m : αριθμός ομάδων

όπου $\theta = \{\theta_1, \dots, \theta_m\}$ και $U = [u_{ij}]$, με $u_{ij} = \begin{cases} 1, & \text{αν } \mathbf{x}_i \in C_j \\ 0, & \text{διαφορετικά} \end{cases}$

- Είναι **επαναληπτικός** αλγόριθμος.
- **Αρχικά** τοποθετεί τους αντιπροσώπους θ_j σε **τυχαίες θέσεις** στο χώρο.
- Προοδευτικά **μετακινεί τους αντιπροσώπους** προς τα **κέντρα** των **φυσικών ομάδων** που σχηματίζουν τα δεδομένα.
- **Τερματίζει** όταν τα θ_j 's δεν αλλάζουν μεταξύ δύο διαδοχικών επαναλήψεων.
- Η **υπολογιστική πολυπλοκότητα** είναι **$O(kN)$** (k = αριθμός επαναλήψεων).
- Απαιτεί **εκ των προτέρων γνώση** του **αριθμού** των **ομάδων m** .

Μη πιθανοτικοί αλγόριθμοι ομαδοποίησης βασιζόμενοι στην βελτιστοποίηση συναρτήσεων κόστους – Ο αλγ. k-means

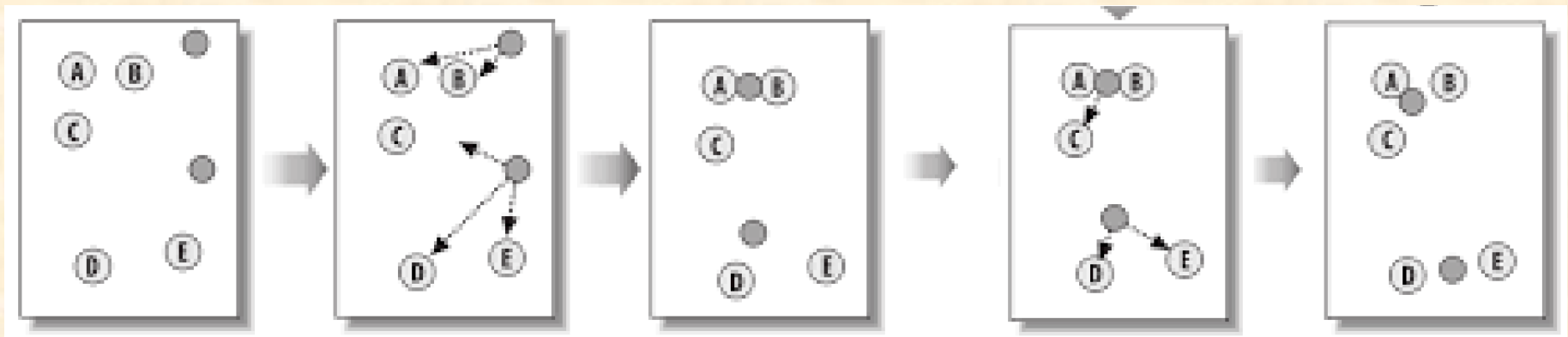
Ο αλγόριθμος k-means

Επέλεξε αυθαίρετες αρχικές εκτιμήσεις $\theta_j(0)$ για τα $\theta_j', j=1, \dots, m$.

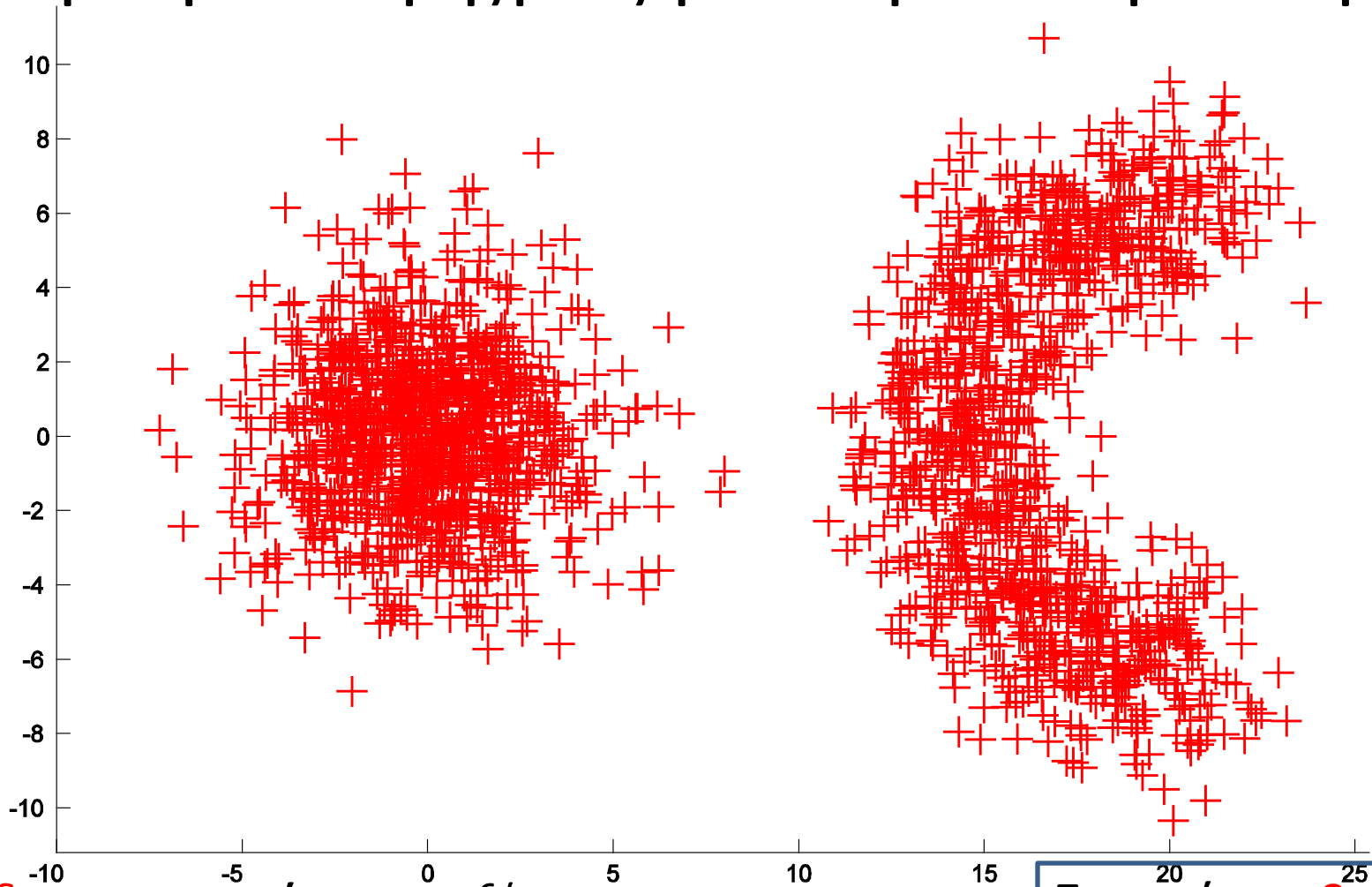
Επανάλαβε

- Για $i=1$ έως N *Προσδιορισμός ομάδων*
ο Προσδιόρισε τον πλησιέστερο αντιπρόσωπο, έστω θ_j , για το x_i
ο Θέσε $u_{ij}=1$ και $u_{iq}=0$, για $q=1, \dots, m, q \neq j$.
- Τέλος {Για}
- Για $j=1$ έως m *Ενημέρωση παραμέτρων*
ο Επανεκτίμησε το θ_j ως το μέσο διάνυσμα των $x_i \in X$ με $u_{ij}=1$.
- Τέλος {Για}

Έως ότου οι τιμές των θ_j παραμείνουν ίδιες μεταξύ δύο διαδοχικών επαναλήψεων



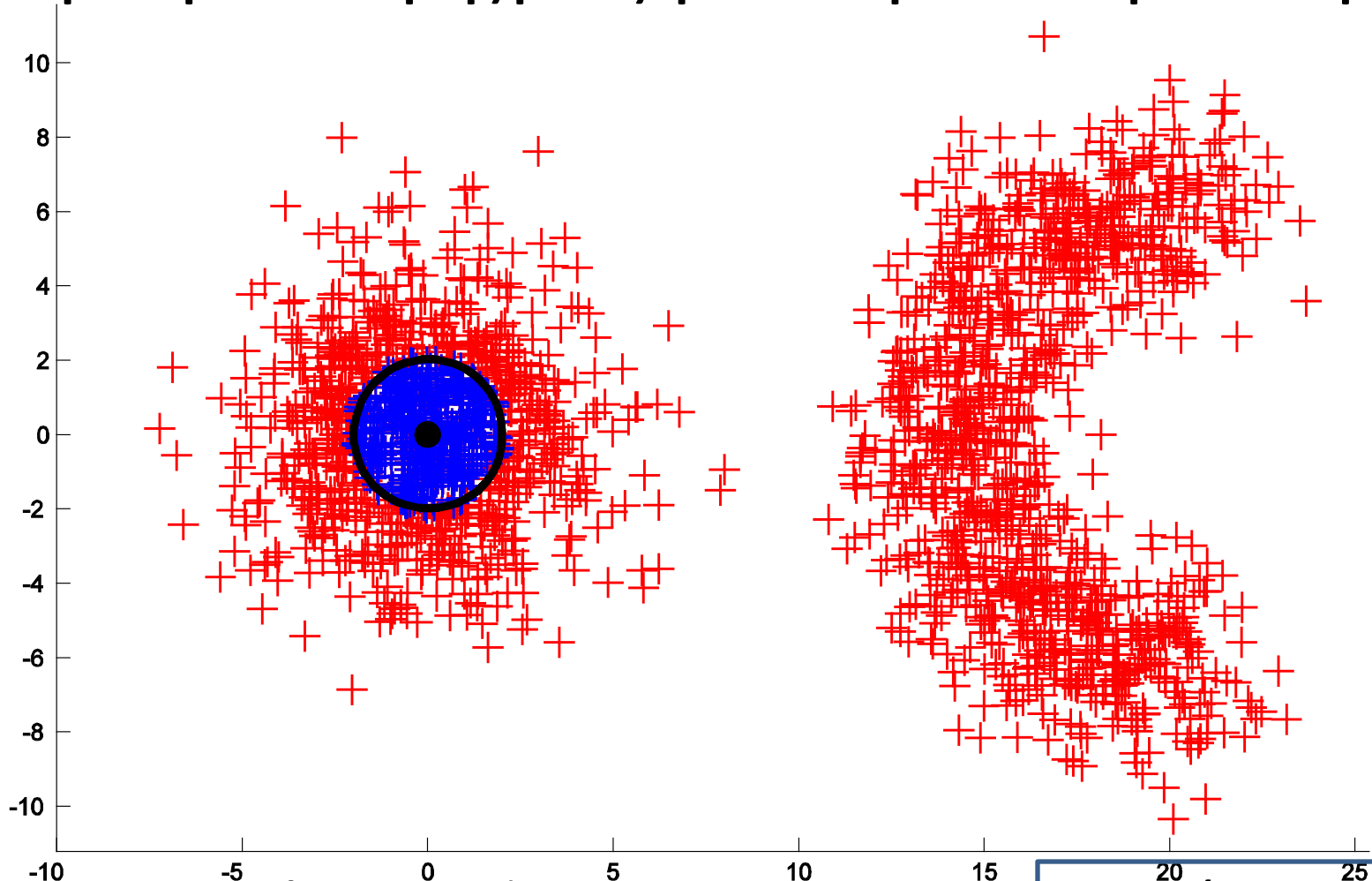
Αλγόριθμοι ομαδοποίησης βασισμένοι στην πυκνότητα των ομάδων



Οι ομάδες ταυτοποιούνται ως εξής:

- Όρισε μια νέα ομάδα C επιλέγοντας ένα σημ. δεδομένων x .
- Καταχώρησε όλα τα σημεία που βρίσκονται στη γειτονιά του x στην ίδια ομάδα.
- Επανάλαβε αναδρομικά το προηγούμενο βήμα έως ότου όλα τα γειτονικά σημεία ΟΛΩΝ των $x \in C$ καταχωρηθούν στο C .

Αλγόριθμοι ομαδοποίησης βασισμένοι στην πυκνότητα των ομάδων

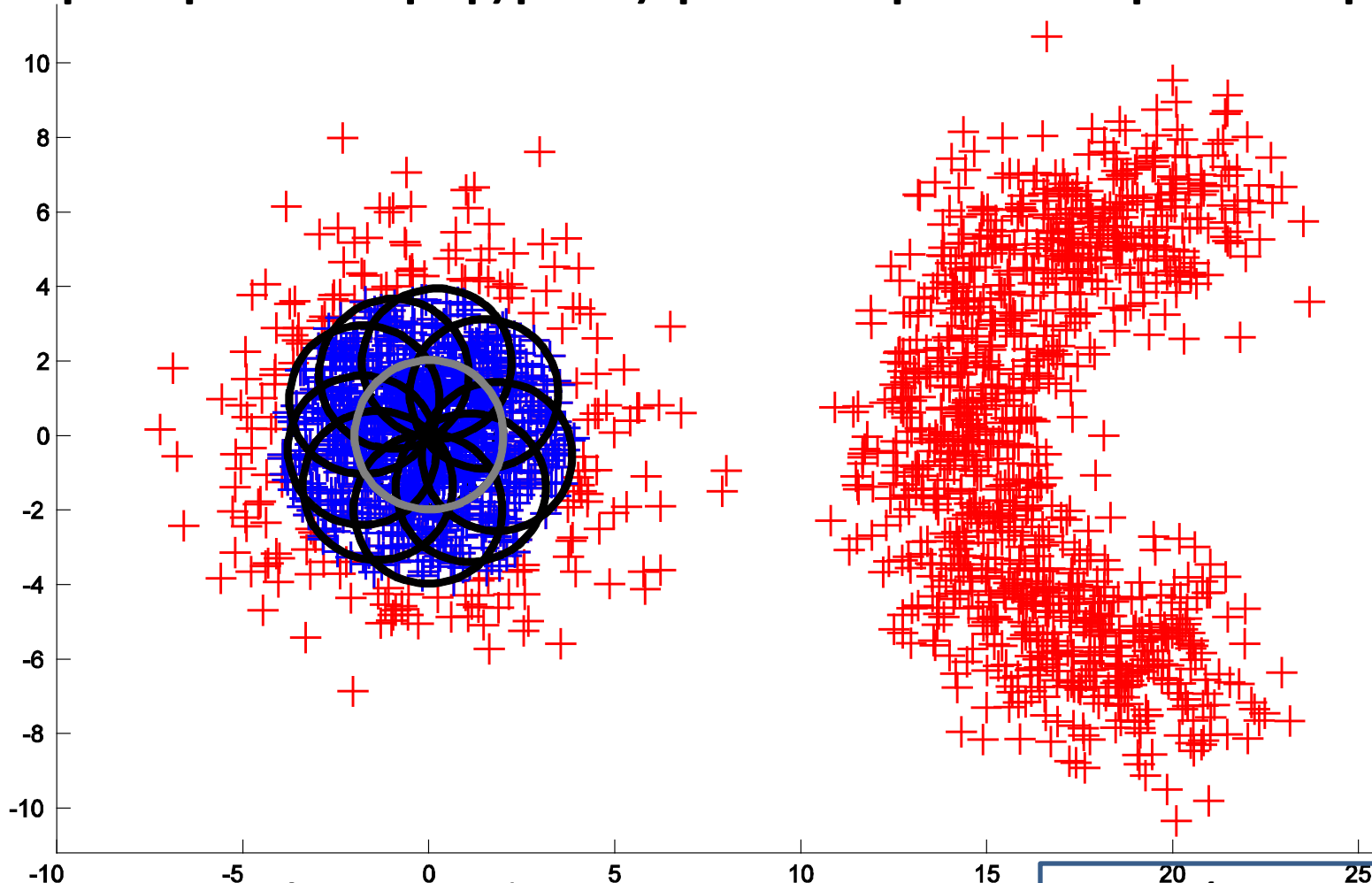


Οι ομάδες ταυτοποιούνται ως εξής:

- Όρισε μια νέα ομάδα C επιλέγοντας ένα σημ. δεδομένων x .
- Καταχώρησε όλα τα σημεία που βρίσκονται στη γειτονιά του x στην ίδια ομάδα.
- Επανάλαβε αναδρομικά το προηγούμενο βήμα έως ότου όλα τα γειτονικά σημεία ΟΛΩΝ των $x \in C$ καταχωρηθούν στο C .

Προαπαίτηση: Ορισμός
μεγέθους γειτονιάς

Αλγόριθμοι ομαδοποίησης βασισμένοι στην πυκνότητα των ομάδων

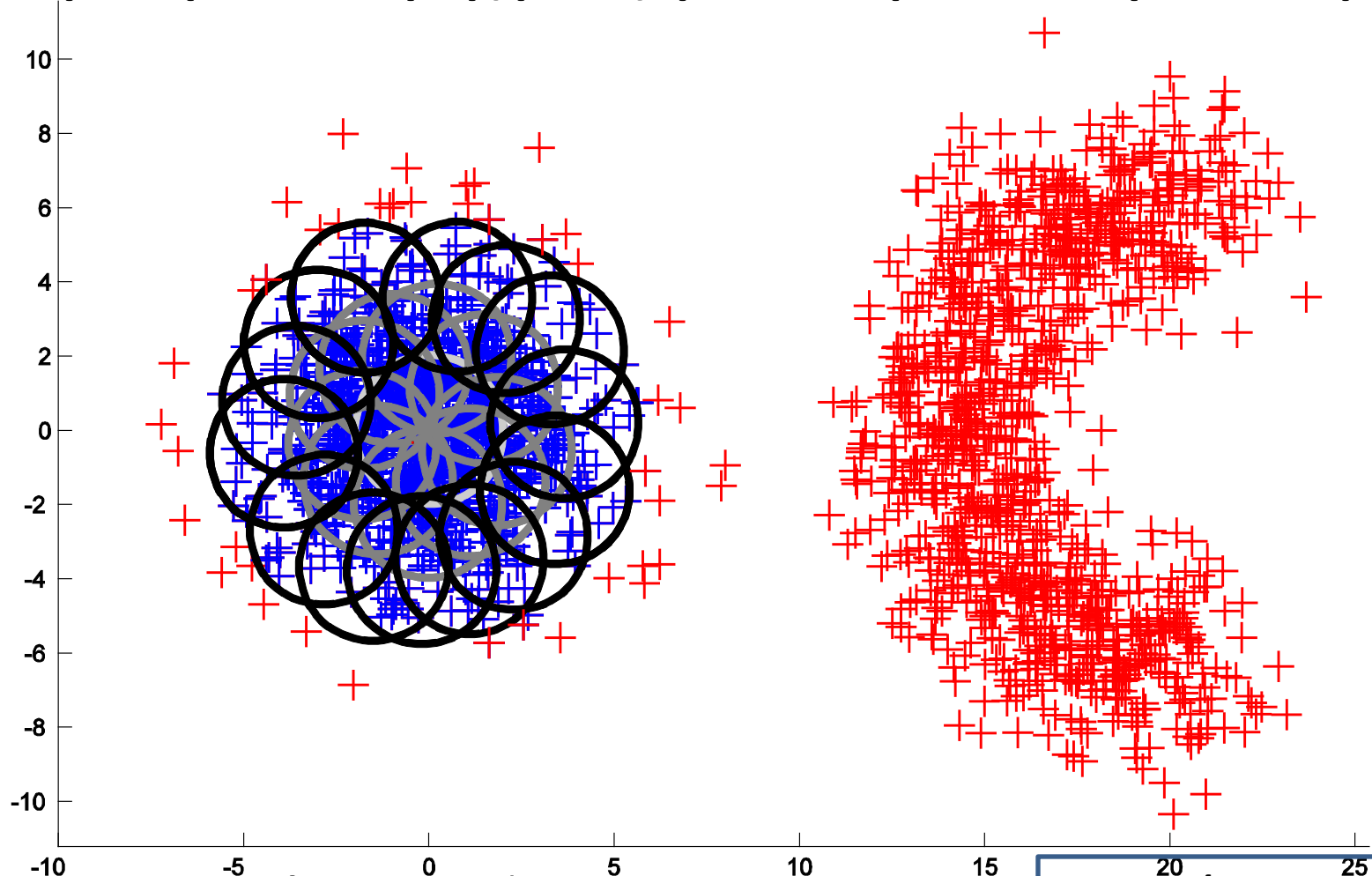


Οι ομάδες ταυτοποιούνται ως εξής:

- Όρισε μια νέα ομάδα C επιλέγοντας ένα σημ. δεδομένων x .
- Καταχώρησε όλα τα σημεία που βρίσκονται στη γειτονιά του x στην ίδια ομάδα.
- Επανάλαβε αναδρομικά το προηγούμενο βήμα έως ότου όλα τα γειτονικά σημεία **ΟΛΩΝ** των $x \in C$ καταχωρηθούν στο C .

Προαπαίτηση: Ορισμός
μεγέθους γειτονιάς

Αλγόριθμοι ομαδοποίησης βασισμένοι στην πυκνότητα των ομάδων

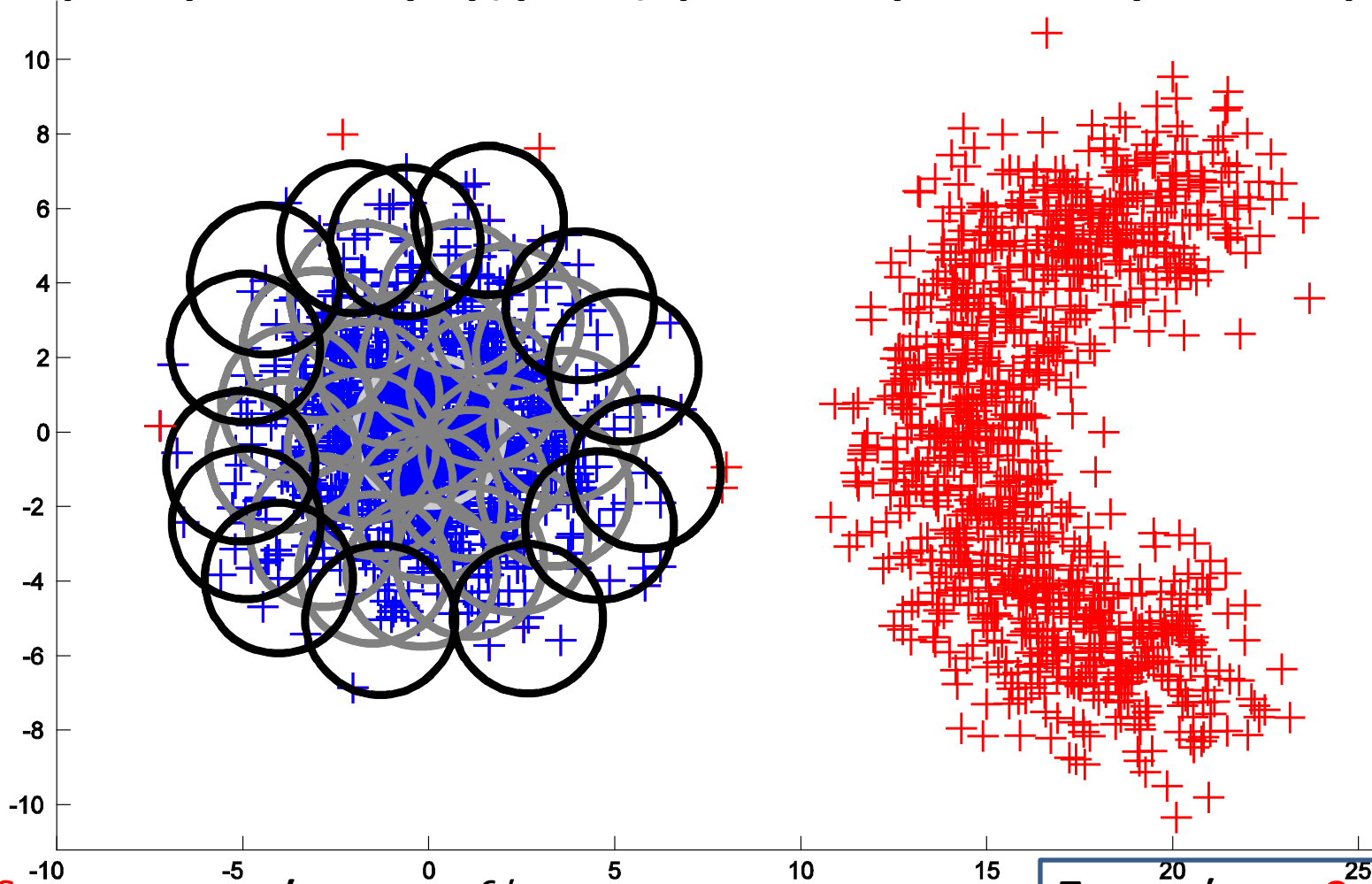


Οι ομάδες ταυτοποιούνται ως εξής:

- Όρισε μια νέα ομάδα C επιλέγοντας ένα σημ. δεδομένων x .
- Καταχώρησε όλα τα σημεία που βρίσκονται στη γειτονιά του x στην ίδια ομάδα.
- Επανάλαβε αναδρομικά το προηγούμενο βήμα έως ότου όλα τα γειτονικά σημεία **ΟΛΩΝ** των $x \in C$ καταχωρηθούν στο C .

Προαπαίτηση: Ορισμός
μεγέθους γειτονιάς

Αλγόριθμοι ομαδοποίησης βασισμένοι στην πυκνότητα των ομάδων

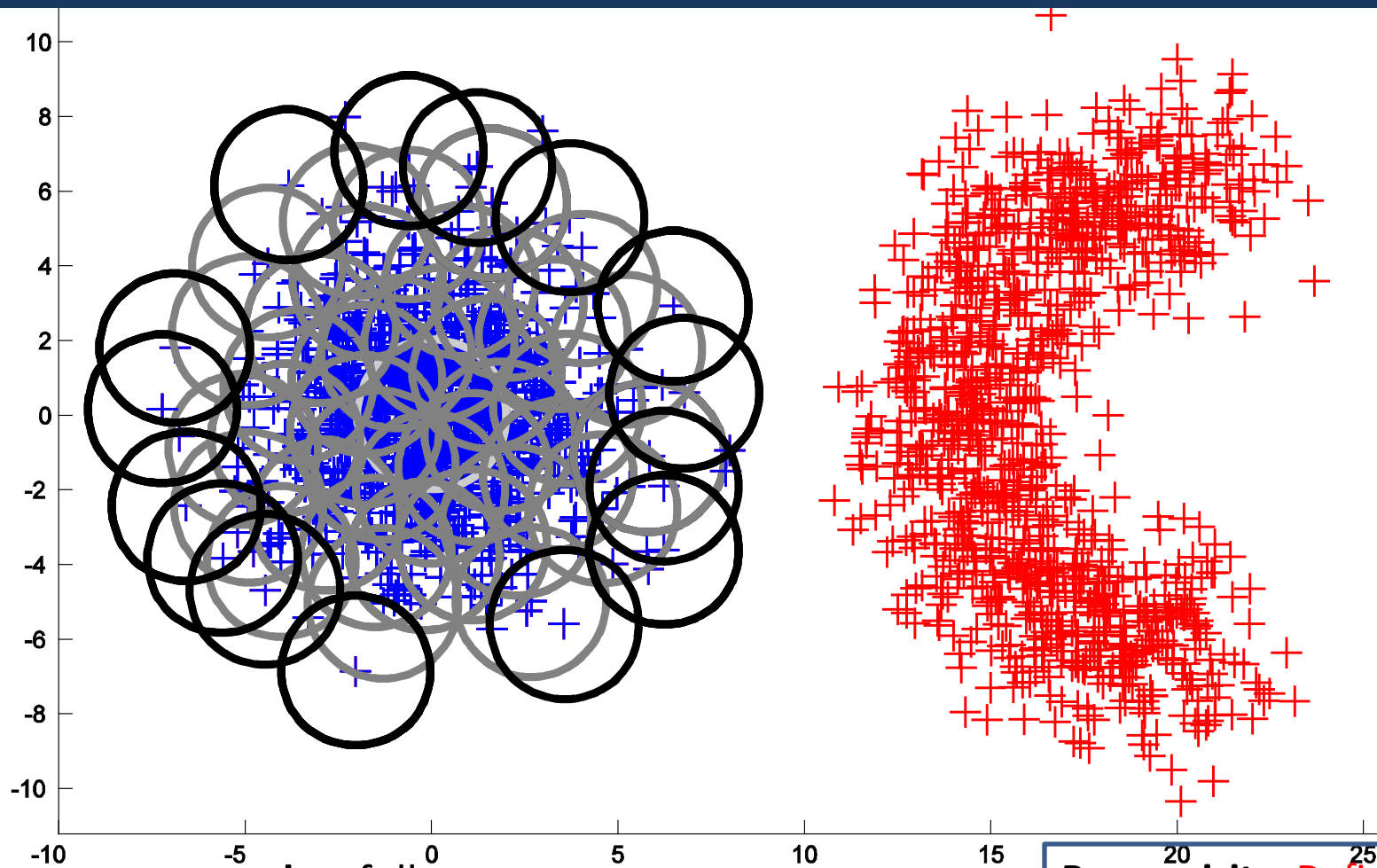


Οι ομάδες ταυτοποιούνται ως εξής:

- Όρισε μια νέα ομάδα C επιλέγοντας ένα σημ. δεδομένων x .
- Καταχώρησε όλα τα σημεία που βρίσκονται στη γειτονιά του x στην ίδια ομάδα.
- Επανάλαβε αναδρομικά το προηγούμενο βήμα έως ότου όλα τα γειτονικά σημεία $O\Lambda\Omega N$ των $x \in C$ καταχωρηθούν στο C .

Προαπαίτηση: Ορισμός
μεγέθους γειτονιάς

Clustering – Density-based algorithms

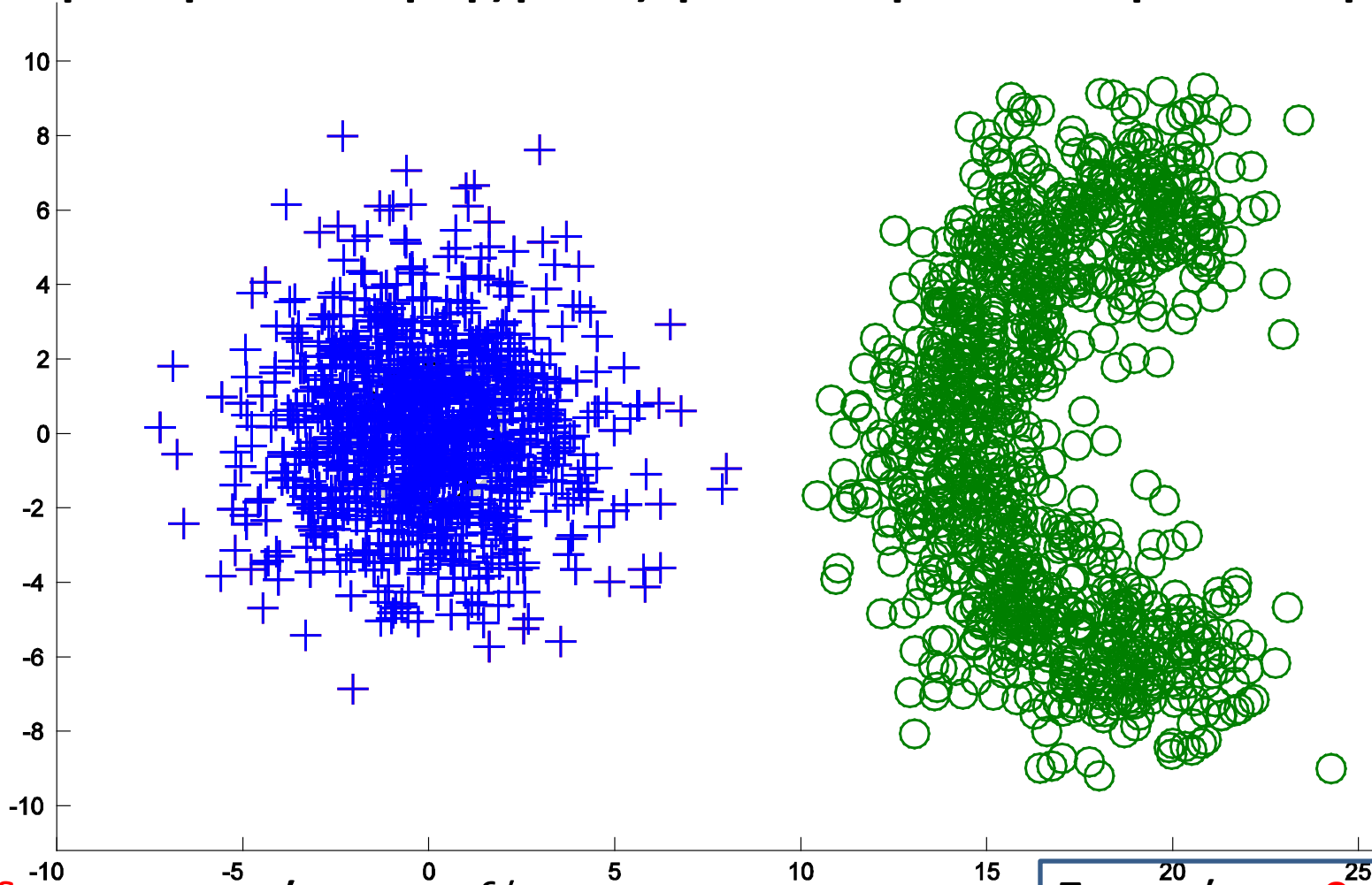


Clusters are **recovered** as follows:

- 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**

Αλγόριθμοι ομαδοποίησης βασισμένοι στην πυκνότητα των ομάδων



Οι ομάδες ταυτοποιούνται ως εξής:

- Όρισε μια νέα ομάδα C επιλέγοντας ένα σημ. δεδομένων x .
- Καταχώρησε όλα τα σημεία που βρίσκονται στη γειτονιά του x στην ίδια ομάδα.
- Επανάλαβε αναδρομικά το προηγούμενο βήμα έως ότου όλα τα γειτονικά σημεία ΟΛΩΝ των $x \in C$ καταχωρηθούν στο C .

Προαπαίτηση: Ορισμός
μεγέθους γειτονιάς

Αλγόριθμοι ομαδοποίησης που βασίζονται σε έννοιες γράφων

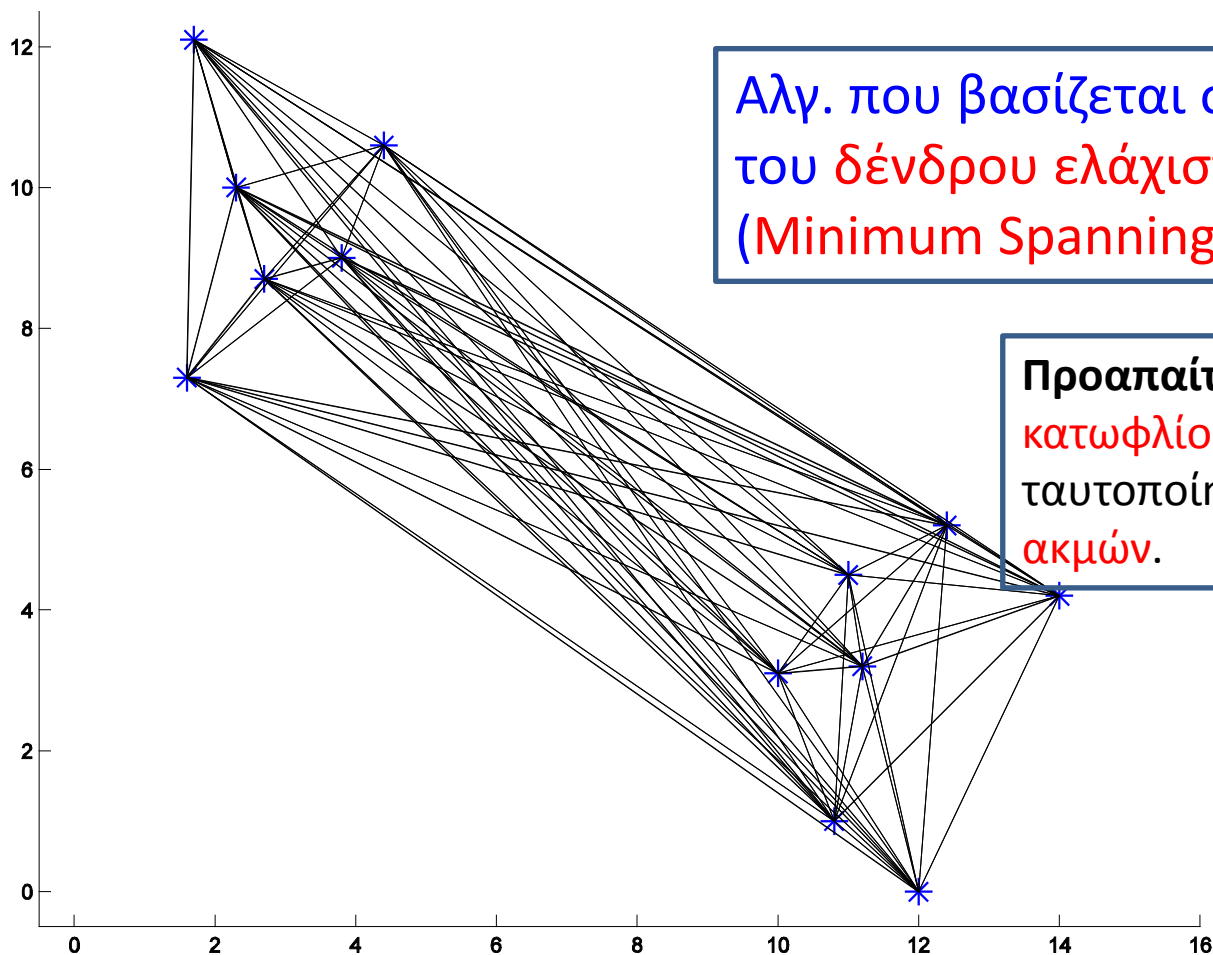


Αλγ. που βασίζεται στην έννοια του δένδρου ελάχιστης κάλυψης (Minimum Spanning Tree - MST).

Προαπαίτηση: Ορισμός κατωφλίου για την ταυτοποίηση “μεγάλων” ακμών.

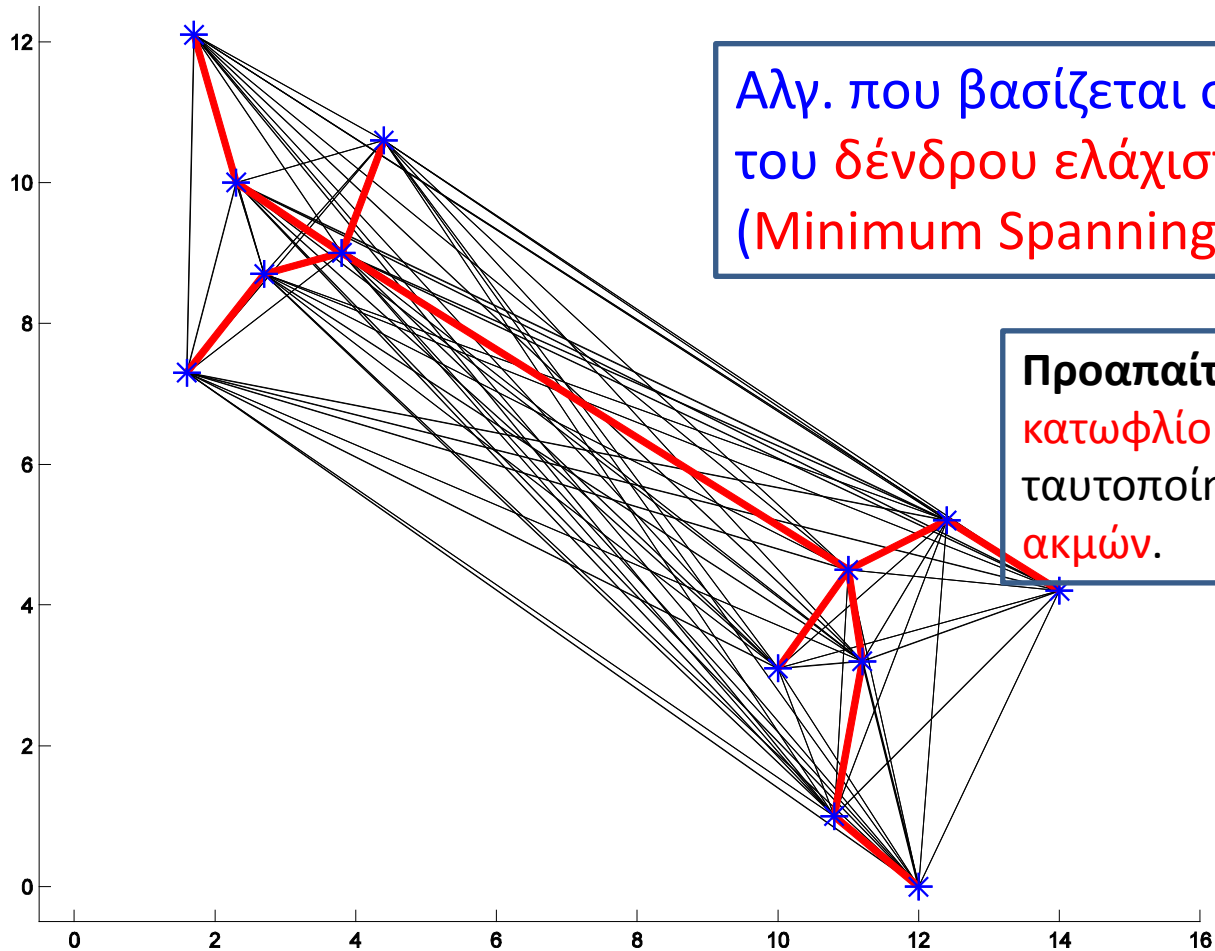
- Όρισε τον **πλήρη γράφο** με **κορυφές** τα **σημεία δεδομένων** και **ακμές** τα **τμήματα** που συνδέουν **κάθε ζεύγος κορυφών**.
- **Στάθμισε** κάθε **ακμή** με την **απόσταση** μεταξύ των **δύο κορυφών** που αυτή συνδέει.
- Όρισε το **MST** του γράφου **και αφάιρεσε** τις **“ασυνήθιστα μεγάλες”** ακμές.
- Οι **υπό-γράφοι** που δημιουργούνται **αντιστοιχούν** στις **ομάδες**.

Αλγόριθμοι ομαδοποίησης που βασίζονται σε έννοιες γράφων



- Όρισε τον **πλήρη γράφο** με **κορυφές** τα **σημεία δεδομένων** και **ακμές** τα **τμήματα** που συνδέουν **κάθε ζεύγος κορυφών**.
- Στάθμισε κάθε **ακμή** με την **απόσταση** μεταξύ των **δύο κορυφών** που αυτή συνδέει.
- Όρισε το **MST** του γράφου **και αφαίρεσε** τις **“ασυνήθιστα μεγάλες”** ακμές.
- Οι **υπό-γράφοι** που δημιουργούνται **αντιστοιχούν** στις **ομάδες**.

Αλγόριθμοι ομαδοποίησης που βασίζονται σε έννοιες γράφων

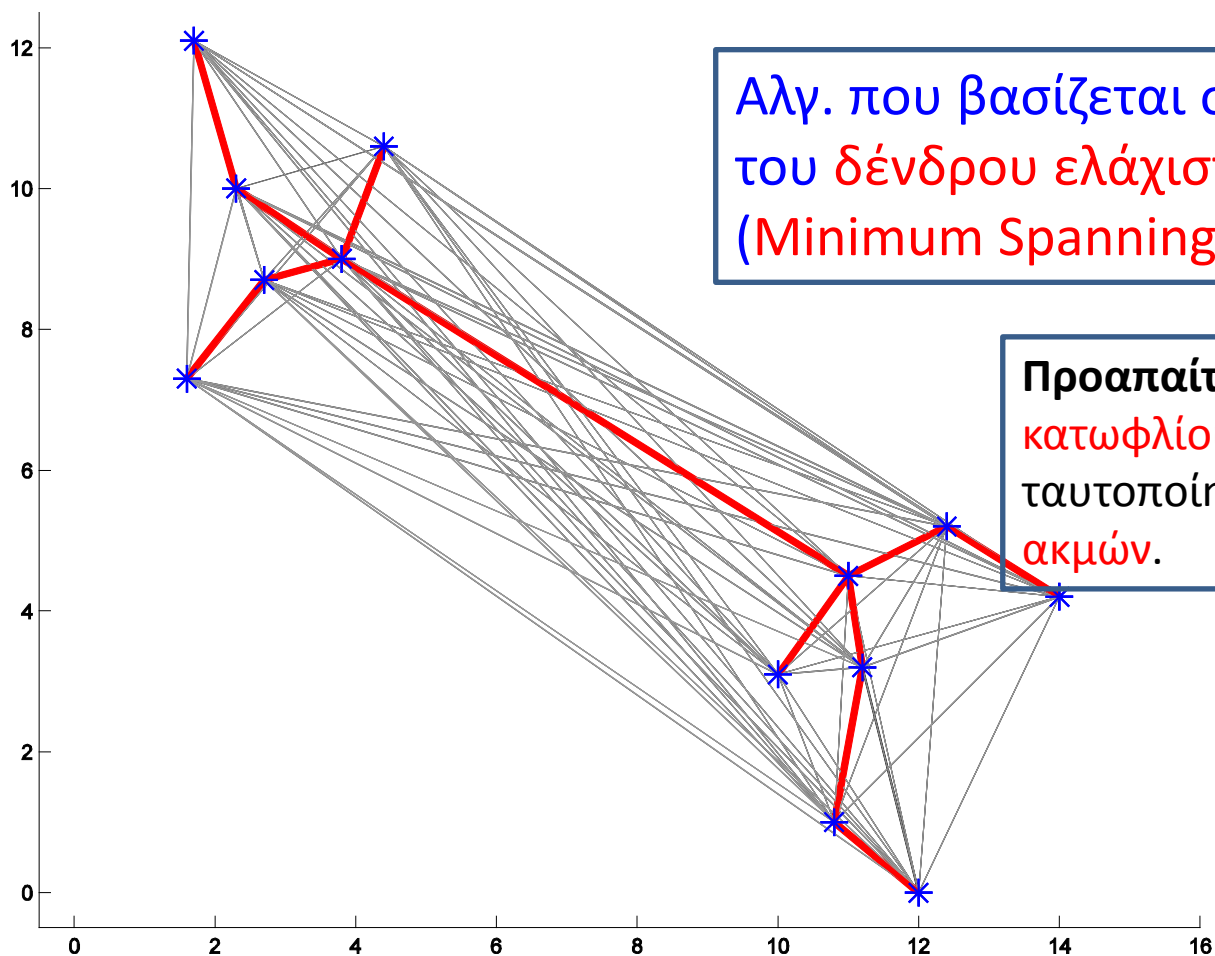


Αλγ. που βασίζεται στην έννοια του δένδρου ελάχιστης κάλυψης (Minimum Spanning Tree - MST).

Προαπαίτηση: Ορισμός κατωφλίου για την ταυτοποίηση “μεγάλων” ακμών.

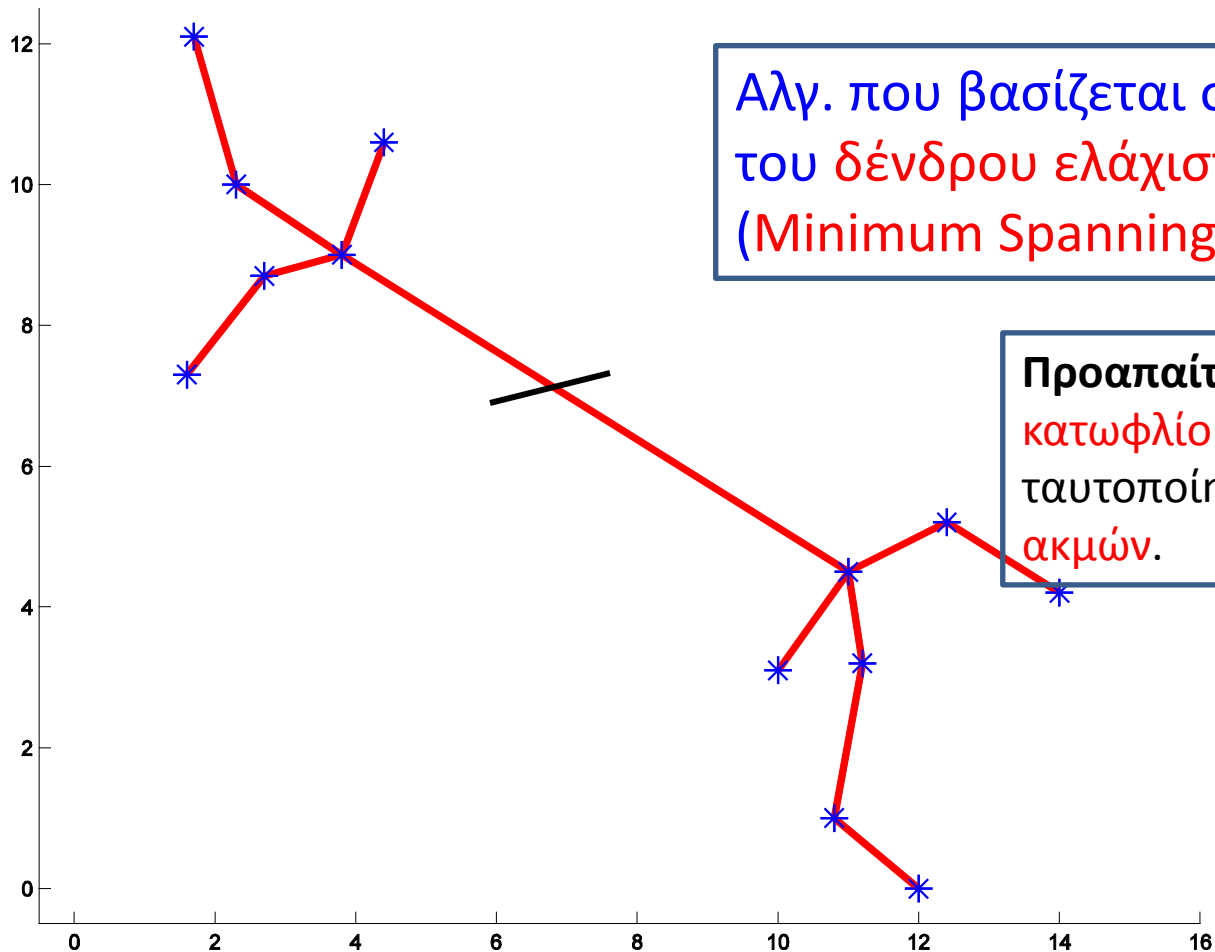
- Όρισε τον **πλήρη γράφο** με **κορυφές** τα **σημεία δεδομένων** και **ακμές** τα **τμήματα** που συνδέουν **κάθε ζεύγος κορυφών**.
- **Στάθμισε** κάθε **ακμή** με την **απόσταση** μεταξύ των **δύο κορυφών** που αυτή συνδέει.
- Όρισε το **MST** του γράφου **και αφάιρεσε** τις **“ασυνήθιστα μεγάλες”** ακμές.
- Οι **υπό-γράφοι** που δημιουργούνται **αντιστοιχούν** στις **ομάδες**.

Αλγόριθμοι ομαδοποίησης που βασίζονται σε έννοιες γράφων



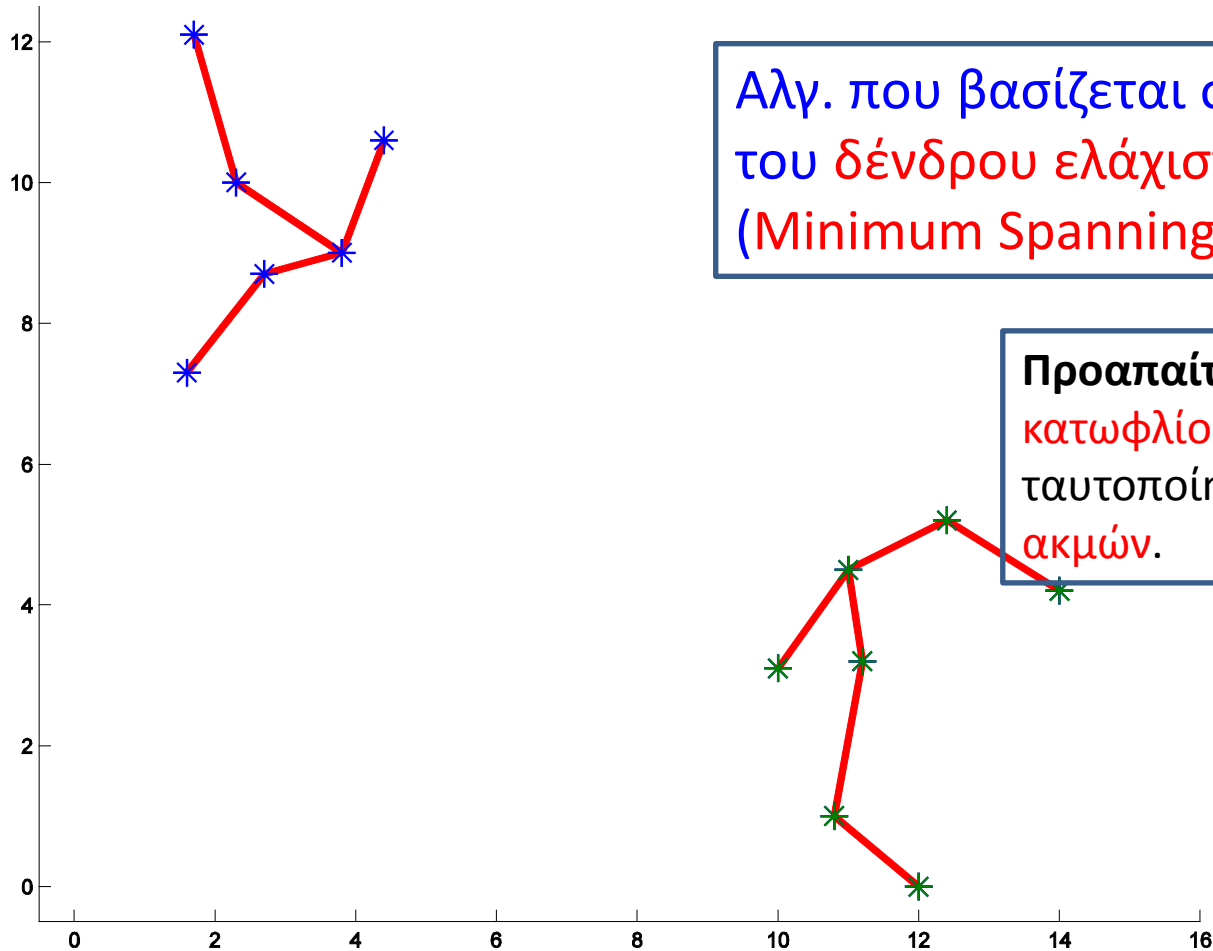
- Όρισε τον **πλήρη γράφο** με **κορυφές** τα **σημεία δεδομένων** και **ακμές** τα **τμήματα** που συνδέουν **κάθε ζεύγος κορυφών**.
- **Στάθμισε** κάθε **ακμή** με την **απόσταση** μεταξύ των **δύο κορυφών** που αυτή συνδέει.
- Όρισε το **MST** του γράφου **και αφαίρεσε** τις **“ασυνήθιστα μεγάλες”** ακμές.
- Οι **υπό-γράφοι** που δημιουργούνται **αντιστοιχούν** στις **ομάδες**.

Αλγόριθμοι ομαδοποίησης που βασίζονται σε έννοιες γράφων



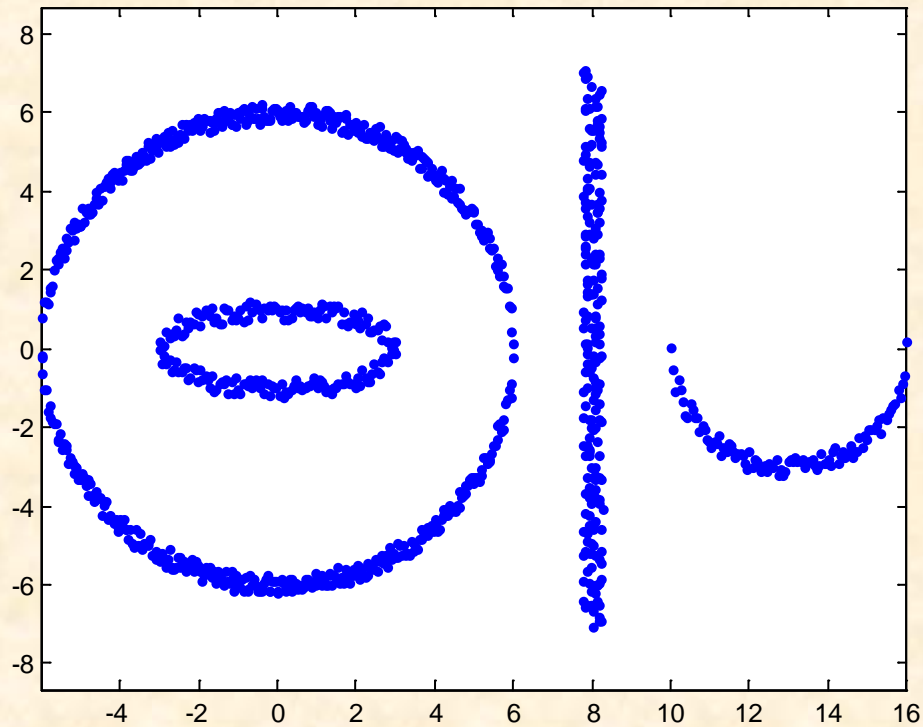
- Όρισε τον **πλήρη γράφο** με **κορυφές** τα **σημεία δεδομένων** και **ακμές** τα **τμήματα** που συνδέουν **κάθε ζεύγος κορυφών**.
- **Στάθμισε** κάθε **ακμή** με την **απόσταση** μεταξύ των **δύο κορυφών** που αυτή συνδέει.
- Όρισε το **MST** του γράφου **και αφάιρεσε** τις **“ασυνήθιστα μεγάλες”** ακμές.
- Οι **υπό-γράφοι** που δημιουργούνται **αντιστοιχούν** στις **ομάδες**.

Αλγόριθμοι ομαδοποίησης που βασίζονται σε έννοιες γράφων



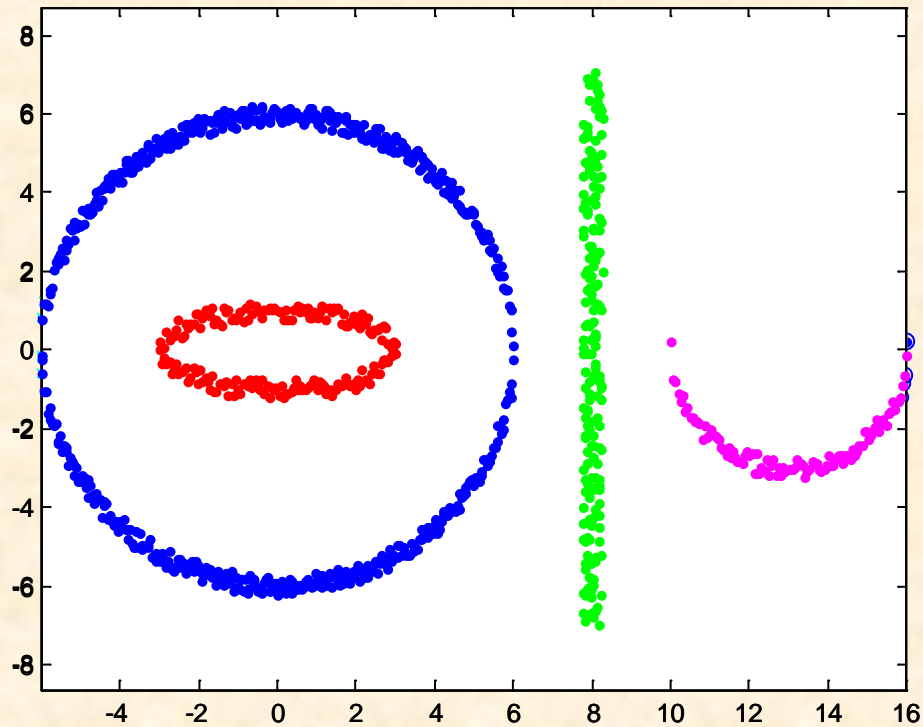
- Όρισε τον **πλήρη γράφο** με **κορυφές** τα **σημεία δεδομένων** και **ακμές** τα **τμήματα** που συνδέουν **κάθε ζεύγος κορυφών**.
- **Στάθμισε** κάθε **ακμή** με την **απόσταση** μεταξύ των **δύο κορυφών** που αυτή συνδέει.
- Όρισε το **MST** του γράφου **και αφαίρεσε** τις **“ασυνήθιστα μεγάλες”** ακμές.
- Οι **υπό-γράφοι** που δημιουργούνται **αντιστοιχούν** στις **ομάδες**.

Παραμετρικοί έναντι μη παραμετρικών αλγορίθμων ομαδοποίησης



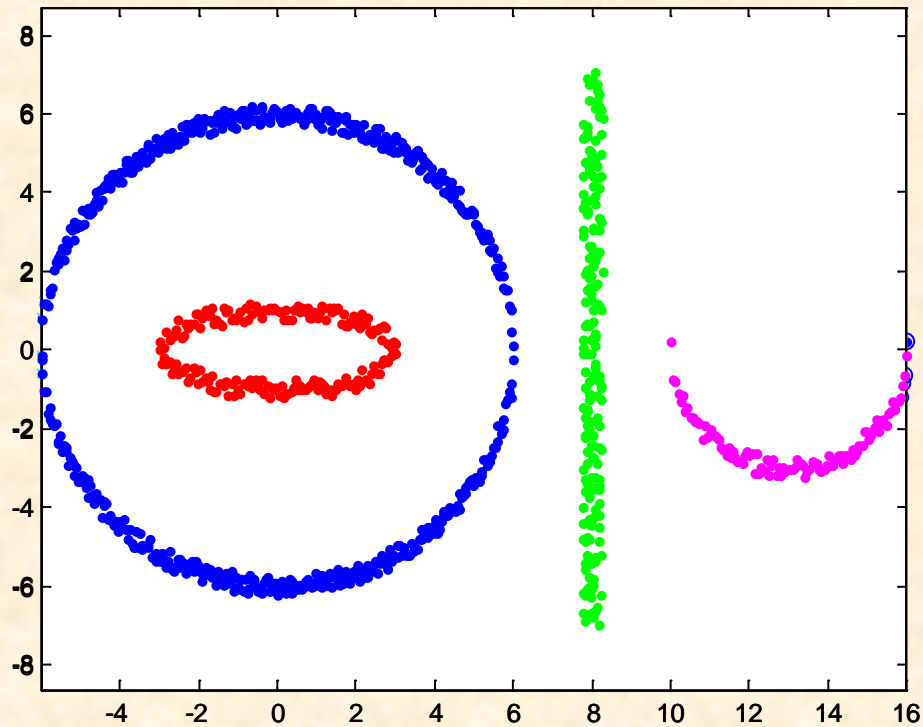
Η περίπτωση ομάδων που δεν έχουν την ίδια δομή και δεν τέμνονται μεταξύ τους.

Παραμετρικοί έναντι μη παραμετρικών αλγορίθμων ομαδοποίησης



Στην περίπτωση ομάδων που δεν έχουν την ίδια δομή και δεν τέμνονται μεταξύ τους, οι παραμετρικοί αλγόριθμοι ομαδοποίησης είναι πιο πιθανό να αποτύχουν, σε σχέση με τους μη-παραμετρικούς αλγόριθμους.

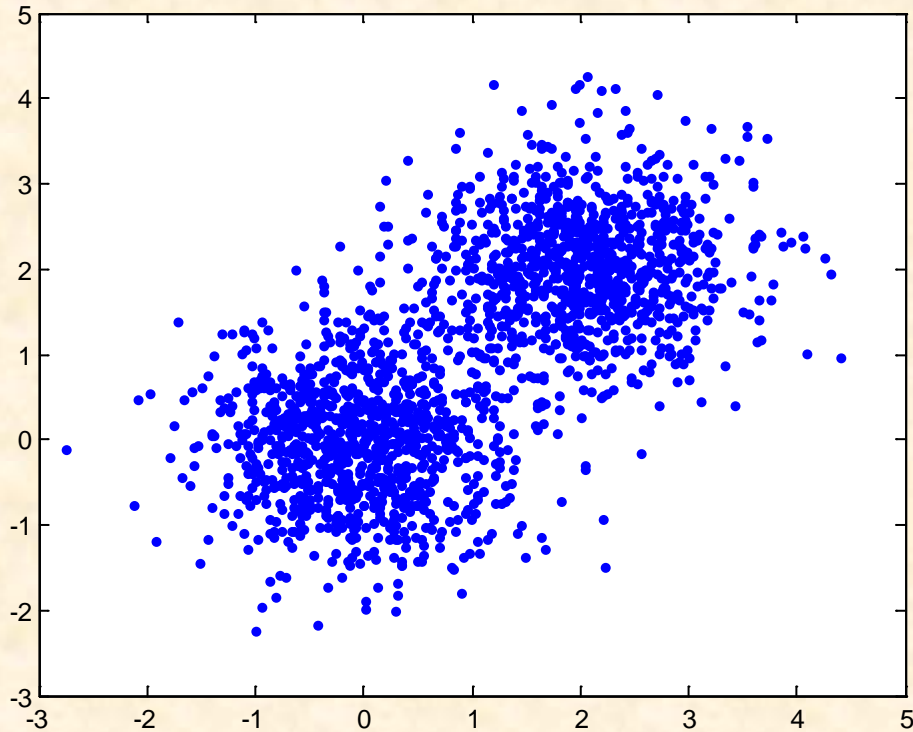
Παραμετρικοί έναντι μη παραμετρικών αλγορίθμων ομαδοποίησης



Μη παραμετρικοί
αλγόριθμοι

Στην περίπτωση ομάδων που δεν έχουν την ίδια δομή και δεν τέμνονται μεταξύ τους, οι παραμετρικοί αλγόριθμοι ομαδοποίησης είναι πιο πιθανό να αποτύχουν, σε σχέση με τους μη-παραμετρικούς αλγόριθμους.

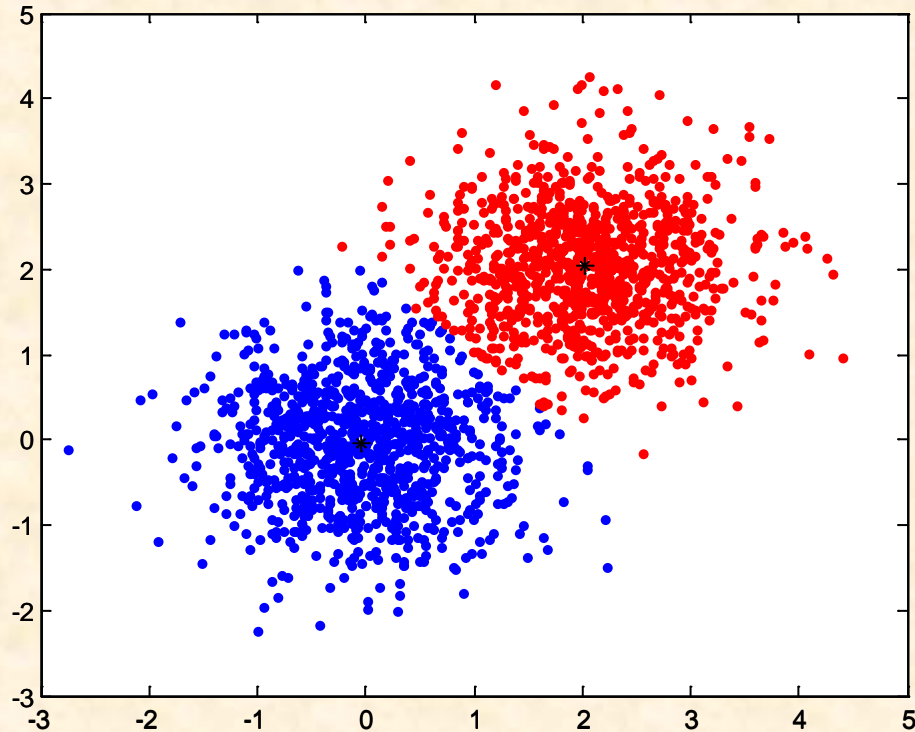
Παραμετρικοί έναντι μη παραμετρικών αλγορίθμων ομαδοποίησης



Μη παραμετρικοί
αλγόριθμοι

Όταν οι ομάδες έχουν την ίδια δομή και παρουσιάζουν επικάλυψη, οι μη παραμετρικοί αλγόριθμοι είναι πιο πιθανό να αποτύχουν σε σχέση με τους παραμετρικούς αλγόριθμους.

Παραμετρικοί έναντι μη παραμετρικών αλγορίθμων ομαδοποίησης



Παραμετρικοί
αλγόριθμοι

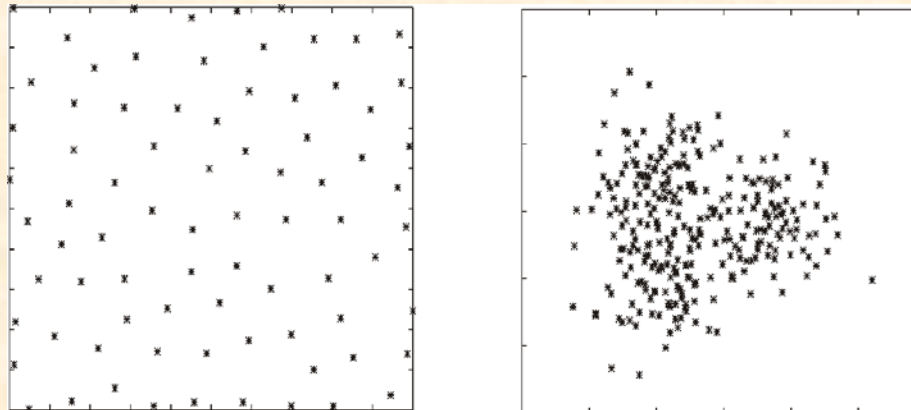
Όταν οι ομάδες έχουν την ίδια δομή και παρουσιάζουν επικάλυψη, οι μη παραμετρικοί αλγόριθμοι είναι πιο πιθανό να αποτύχουν σε σχέση με τους παραμετρικούς αλγόριθμους.

Αξιολόγηση αποτελεσμάτων ομαδοποίησης

- Αυτή πραγματοποιείται από έναν **ειδικό** στο πεδίο της εφαρμογής.
- Εμπεριέχει πάντα την **υποκειμενικότητα** του **ειδικού**.
- Αν τα αποτελέσματα **δεν είναι ικανοποιητικά**, η όλη διαδικασία ομαδοποίησης μπορεί να επαναληφθεί με
 - a) διαφορετικό μέτρο εγγύτητας ή/και
 - b) διαφορετική αναπαράσταση δεδομένων ή/και
 - c) διαφορετικό αλγόριθμο ομαδοποίησης

Μια τελευταία σημείωση: Δεν είναι όλα τα προβλήματα κατάλληλα για επεξεργασία με τεχνικές ομαδοποίησης.

(Για **παράδειγμα**, τα διανύσματα σε δεδομένο πρόβλημα, μπορεί να μην σχηματίζουν καθόλου φυσικές ομάδες).



Δεδομένα που παρουσιάζουν **κανονικότητες** (Regularly spaced Data)

Δεδομένα που σχηματίζουν **μία φυσική ομάδα**.