

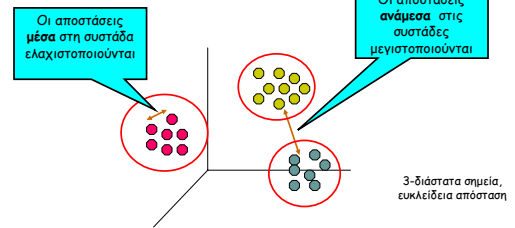
Συσταδοποίηση II

Μέρος των διαφανειών είναι από το P.-N. Tan, M. Steinbach, V. Kumar, «Introduction to Data Mining», Addison Wesley, 2006



Τι είναι συσταδοποίηση

Εύρεση συστάδων αντικειμένων έτσι ώστε τα αντικείμενα σε κάθε ομάδα να είναι όμοια (ή να σχετίζονται) και διαφορετικά (ή μη σχετιζόμενα) από τα αντικείμενα των άλλων ομάδων



Εξώφυλλο Δεδομένων: Ακ. Έτος 2007-2008

ΣΥΣΤΑΔΟΠΟΙΗΣΗ II

2

Γενικές Απαιτήσεις

- Scalability - στον αριθμό σημείων και διαστάσεων
- Να υποστηρίζει διαφορετικούς τύπους δεδομένων
- Να υποστηρίζει συστάδες με διαφορετικά σχήματα (συνήθως, «σφαίρες»)
- Να είναι εύκολο να δώσουμε τιμές στις παραμέτρους εισόδου (αριθμό συστάδων, μέγεθος κλπ)
- Να μην εξαρτάται από τη σειρά επεξεργασίας των σημείων εισόδου

Εξώφυλλο Δεδομένων: Ακ. Έτος 2007-2008

ΣΥΣΤΑΔΟΠΟΙΗΣΗ II

3

Γενικές Απαιτήσεις

- Δυναμικά μεταβαλλόμενα δεδομένα
 - Αλλαγή συστάδων με το πέρασμα του χρόνου
- Απόδοση (scaling)
 - Disk-resident vs Main memory

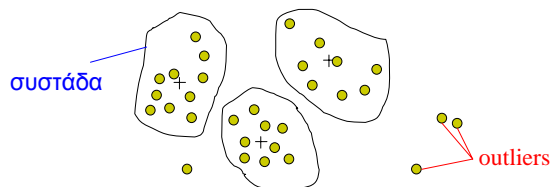
Εξώφυλλο Δεδομένων: Ακ. Έτος 2007-2008

ΣΥΣΤΑΔΟΠΟΙΗΣΗ II

4

Γενικές Απαιτήσεις

Αντιμετώπιση θορύβου και outliers



Εξώφυλλο Δεδομένων: Ακ. Έτος 2007-2008

ΣΥΣΤΑΔΟΠΟΙΗΣΗ II

5

Είδη συσταδοποίησης

Μια συσταδοποίηση είναι ένα σύνολο από συστάδες:

Διαχωριστική Συσταδοποίηση (Partitional Clustering)

Ένας διαμερισμός των αντικειμένων σε μη επικαλυπτόμενα - non-overlapping - υποσύνολα (συστάδες) τέτοιος ώστε κάθε αντικείμενο ανήκει σε ακριβώς ένα υποσύνολο

Ιεραρχική Συσταδοποίηση (Hierarchical clustering)

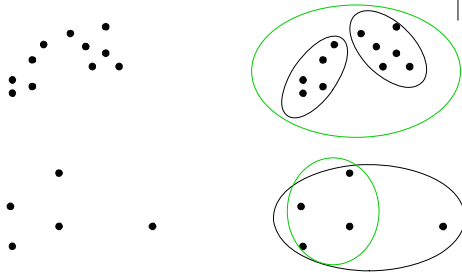
Ένα σύνολο από *εμφωλευμένες* (nested) συστάδες. Επιτρέπουμε σε μια συστάδα να έχει υποσυστάδες οργανωμένες σε ένα ιεραρχικό δέντρο

Εξώφυλλο Δεδομένων: Ακ. Έτος 2007-2008

ΣΥΣΤΑΔΟΠΟΙΗΣΗ II

6

Διαχωριστική και Ιεραρχική Συσταδοποίηση



Αρχικά Σημεία

Άλλες διακρίσεις μεταξύ συνόλων συστάδων

Επικαλυπτόμενο ή όχι

Ένα σημείο ανήκει σε περισσότερες από μια συστάδες (πχ οριακά σημεία)

Ασαφή συσταδοποίηση

Στην ασαφή συσταδοποίηση ένα σημείο ανήκει σε κάθε συστάδα με κάποιο βάρος μεταξύ του 0 και του 1

Συχνά τα βάρη για κάθε σημείο έχουν άθροισμα 1

Η πιθανοτική συσταδοποίηση έχει παρόμοια χαρακτηριστικά

Μερική - Πλήρης

Σε ορισμένες περιπτώσεις θέλουμε να ομαδοποιήσουμε μόνο κάποια από τα δεδομένα (άλλα θόρυβος, ή μη ενδιαφέρουσα πληροφορία)

Ετερογενή - Ομογενή

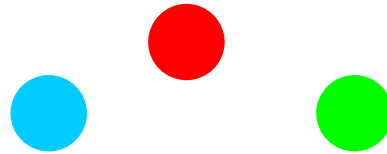
Συστάδες με πολύ διαφορετικά μεγέθη, σχήματα και πυκνότητες (densities)

Είδη Συστάδων

- Καλώς διαχωρισμένες συστάδες
- Συστάδες βασισμένες σε κέντρο
- Συνεχής (contiguous) συστάδες
- Συστάδες βασισμένες σε πυκνότητα
- Βασισμένα σε ιδιότητες ή έννοιες
- Περιγράφονται από μια αντικειμενική συνάρτηση (Objective Function)

Καλώς Διαχωρισμένες Συστάδες

Μια συστάδα είναι ένα σύνολο από σημεία τέτοιο ώστε κάθε σημείο μιας ομάδας είναι **κοντινότερο σε (ή πιο όμοιο με) όλα τα άλλα σημεία της ομάδας** από ότι σε οποιοδήποτε άλλο σημείο που δεν ανήκει στη συστάδα.



3 καλώς-διαχωρισμένες συστάδες

Συχνά υπάρχει η έννοια του κατωφλιού (threshold)

Όχι απαραίτητα κυκλικοί (οποιοδήποτε σχήμα)

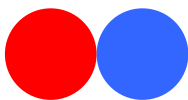
Συστάδες βασισμένες σε κέντρο ή πρότυπο

Μια συστάδα είναι ένα σύνολο από αντικείμενα τέτοιο ώστε ένα αντικείμενο στην ομάδα είναι **κοντινότερο σε (ή πιο όμοιο με) το «κέντρο» ή πρότυπο** της ομάδας από ότι από το κέντρο οποιασδήποτε άλλης ομάδας.

Το κέντρο της ομάδας είναι συχνά

▪ **centroid**, ο μέσος όρος των σημείων της συστάδας, ή

▪ **a medoid**, το πιο «αντιπροσωπευτικό» σημείο της συστάδας (πχ όταν κατηγορικά γνωρίσματα)



4 συστάδες βασισμένες σε κέντρο



Τείνουν στο να είναι κυκλικοί

Συνεχής Συστάδες

Συνεχής Συστάδες (Contiguous Cluster) (Κοντινότερος γείτονα ή μεταβατικά)

Μια συστάδα είναι ένα σύνολο σημείων τέτοιο ώστε κάθε σημείο είναι **πιο κοντά σε ένα ή περισσότερα σημεία της συστάδας από ό,τι σε οποιοδήποτε σημείο εκτός συστάδας**

Συχνά σε περιπτώσεις συστάδων με μη κανονικό σχήμα ή με αλληλοπλεκόμενα σχήματα - ή όταν έχουμε γραφήματα και θέλουμε να βρούμε συνεκτικά υπογραφήματα

Πρόβλημα με θόρυβο

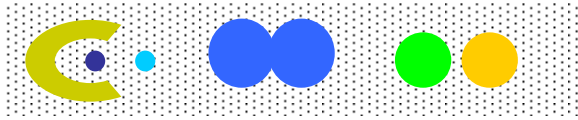


8 συνεχείς συστάδες

Συστάδες βασισμένες στην πυκνότητα

Μια συστάδα είναι μια **πυκνή περιοχή** από σημεία την οποία χωρίζουν από άλλες περιοχές μεγάλης πυκνότητας περιοχές χαμηλής πυκνότητας

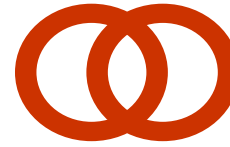
Συχνά σε περιπτώσεις συστάδων με μη κανονικό σχήμα ή με αλληλοπλεκόμενα σχήματα ή όταν θόρυβος ή outliers



6 συστάδες βασισμένες στην πυκνότητα

Εννοιολογική συσταδοποίηση

Συστάδες με κοινή ιδιότητα ή εννοιολογικές συστάδες.



2 αλληλοκαλυπτόμενοι κύκλοι

Συστάδες βασισμένες σε μια Αντικειμενική Συνάρτηση

Εύρεση συστάδων που ελαχιστοποιούν ή μεγιστοποιούν μια **αντικειμενική συνάρτηση**

Απαρίθμηση όλων των δυνατών τρόπων χωρισμού των σημείων σε συστάδες και υπολογισμού του «πόσο καλό» ("goodness") είναι κάθε πιθανό σύνολο από συστάδες χρησιμοποιώντας τη δοθείσα αντικειμενική συνάρτηση (NP-hard)

Οι στόχοι (objectives) μπορεί να είναι ολικό (global) ή τοπικό (local)
Οι ιεραρχικοί συνήθως τοπικού
Οι διαχωριστικοί ολικές

Αλγόριθμοι Συσταδοποίησης

Θα δούμε ανάμεσα σε άλλους τους:

- **K-means και παραλλαγές**
- **Ιεραρχική Συσταδοποίηση**
- **Συσταδοποίηση με βάση την Πυκνότητα (DBSCAN)**
- **BIRCH (δεδομένα στο δίσκο!)**

K-means

K-means: Γενικά

Διαχωριστικός αλγόριθμος

(βασισμένος σε πρότυπο) Κάθε συστάδα συσχετίζεται με ένα **κεντρικό σημείο (centroid)**

Κάθε σημείο ανατίθεται στη συστάδα με το κοντινότερο κεντρικό σημείο

Ο αριθμός των ομάδων, K, είναι είσοδος στον αλγόριθμο

K-means: Βασικός Αλγόριθμος



Βασικός αλγόριθμος

- 1: Επιλογή K σημείων ως τα αρχικά κεντρικά σημεία
- 2: **Repeat**
- 3: Ανάθεση όλων των αρχικών σημείων στο κοντινότερο τους από τα K κεντρικά σημεία
- 4: Επανα-υπολογισμός του κεντρικού σημείου κάθε συστάδας
- 5: **Until** τα κεντρικά σημεία να μην αλλάζουν

K-means: Βασικός Αλγόριθμος



Παρατηρήσεις

1. Τα αρχικά κεντρικά σημεία συνήθως επιλέγονται τυχαία

Οι συστάδες που παράγονται διαφέρουν από το ένα τρέξιμο του αλγορίθμου στο άλλο

K-means: Βασικός Αλγόριθμος



Παρατηρήσεις (συνέχεια)

2. Η εγγύτητα των σημείων υπολογίζεται με βάση κάποια απόσταση που εξαρτάται από το είδος των σημείων, στα παραδείγματα θα θεωρήσουμε την Ευκλείδεια απόσταση

▪ Επειδή η απόσταση υπολογίζεται συχνά πρέπει να είναι σχετικά απλή

3. Το κεντρικό σημείο είναι (συνήθως) το μέσο (mean) των σημείων της συστάδας (το οποίο μπορεί να μην είναι ένα από τα δεδομένα εισόδου)

Περίληψη Δεδομένων



Μια παρένθεση

- **Αριθμητικό Μέσο - Mean (αλγεβρική μέτρηση)**
(sample vs. population):

- Αριθμητικό μέσο με βάρος (Weighted arithmetic mean)
- Trimmed mean: κόβουμε τις ακραίες τιμές (πχ τα μεγαλύτερα και μικρότερα (p/2)%

- **Μέσο (median):**

- Μεσαία τιμή αν μονός αριθμός, ο μέσος όρος των δυο μεσαίων τιμών, αλλιώς
Καλύτερα όταν skewed δεδομένα

Distributed measure (κατανομισμένη μέτρηση): μπορούν να υπολογιστούν αν χωρίσουμε τα αρχικά δεδομένα σε μικρότερα υποσύνολα, υπολογίσουμε την τιμή σε κάθε υποσύνολο και τις συγχωνεύουμε πχ sum(), count(), max(), min()

Algebraic measure (αλγεβρική μέτρηση): μπορεί να υπολογιστεί αν εφαρμόσουμε μια αλγεβρική (πολυωνυμική) συνάρτηση σε μία ή περισσότερες κατανομισμένες μετρήσεις (πχ avg()= sum()/count())

Holistic measure (ολιστική μέτρηση) πρέπει να υπολογιστεί στο σύνολο των δεδομένων

Γενική Τάση

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\bar{x} = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$$

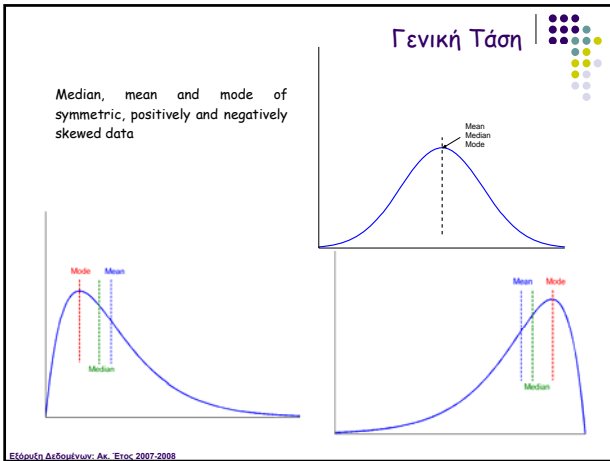


Γενική Τάση

- **Mode**
 - Η τιμή που εμφανίζεται πιο συχνά στα δεδομένα
 - Unimodal, bimodal, trimodal
 - Εμπειρικός τύπος:

$$mean - mode = 3 \times (mean - median)$$

- **Midrange**
 - $(\min() + \max()) / 2$



Διασπορά

Variance (σ^2)

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \left[\sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \right]$$

Standard deviation (σ)

Εύρεση Δεδομένων: Ακ. Έτος 2007-2008

Περίληψη Δεδομένων

Κλείνει η παρένθεση

Εύρεση Δεδομένων: Ακ. Έτος 2007-2008

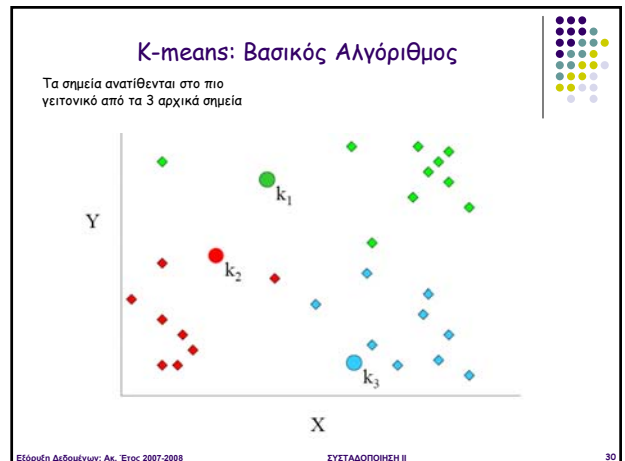
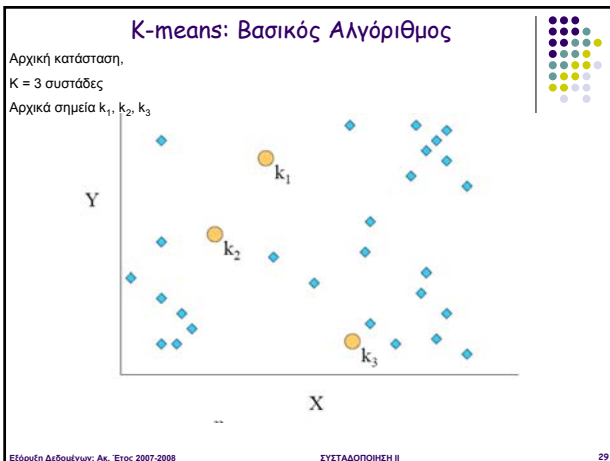
K-means: Βασικός Αλγόριθμος

Παράδειγμα

2 4 10 12 3 20 30 11 15

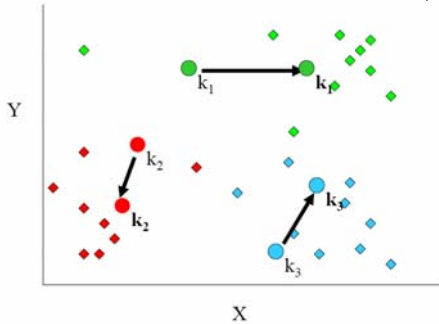
Έστω $k = 2$, και αρχικά επιλεγούμε το 3 και το 4

Εύρεση Δεδομένων: Ακ. Έτος 2007-2008



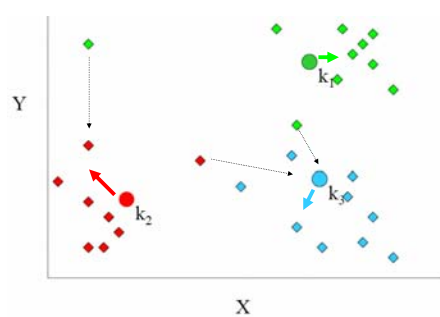
K-means: Βασικός Αλγόριθμος

Επανα-υπολογισμός του κέντρου (κέντρου βάρους) κάθε σημείου

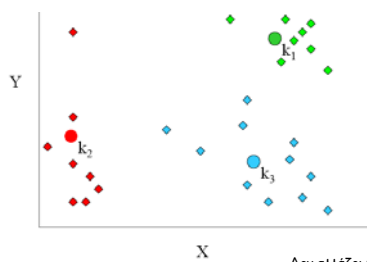


K-means: Βασικός Αλγόριθμος

Νέα ανάθεση των σημείων
Νέα κέντρα βάρους



K-means: Βασικός Αλγόριθμος



K-means: Βασικός Αλγόριθμος

Παρατηρήσεις (συνέχεια)

- Χώρος: αποθηκεύουμε μόνο τα κέντρα
- Η πολυπλοκότητα είναι $O(I \cdot n \cdot K \cdot d)$
 n = αριθμός σημείων,
 K = αριθμός συστάδων,
 I = αριθμός επαναλήψεων,
 d = αριθμός γνωρισμάτων (διάσταση)

K-means: Βασικός Αλγόριθμος

Παρατηρήσεις (συνέχεια)

- Για συνηθισμένα μέτρα ομοιότητας, ο αλγόριθμος **συγκλίνει**. Η σύγκλιση συμβαίνει συνήθως τις αρχικές πρώτες επαναλήψεις
- Συχνά η **τελική συνθήκη** αλλάζει σε
Until σχετικά λίγα σημεία να αλλάζουν συστάδα – ή η απόσταση μεταξύ των νέων κεντρικών σημείων από τα παλιά να είναι μικρή

K-means: Εκτίμηση ποιότητας

Ουσιαστικά, ο αλγόριθμος προσπαθεί επαναληπτικά να «μειώσει» την απόσταση από ένα σημείο της συστάδας

Η πιο συνηθισμένη μέτρηση είναι το *άθροισμα των τετράγωνων του λάθους* (Sum of Squared Error (SSE))

Για κάθε σημείο, το λάθος είναι η απόστασή του από την κοντινότερη συστάδα

Για να πάρουμε το SSE, παίρνουμε το τετράγωνο αυτών των λαθών και τα προσθέτουμε

$$SSE = \sum_{i=1}^K \sum_{x \in C_i} dist^2(m_i, x)$$

Όπου $dist$ Ευκλείδεια απόσταση, x είναι ένα σημείο στη συστάδα C_i και m_i είναι ο αντιπρόσωπος (κεντρικό σημείο) της συστάδας C_i

Μπορούμε να δείξουμε ότι το σημείο που ελαχιστοποιεί το SSE για τη συστάδα είναι ο μέσος όρος $c_i = 1/m_i \sum_{x \in C_i} x$

Δοθέντων δύο συστάδων, μπορούμε να επιλέξουμε αυτήν με το μικρότερο λάθος

K-means: Εκτίμηση ποιότητας

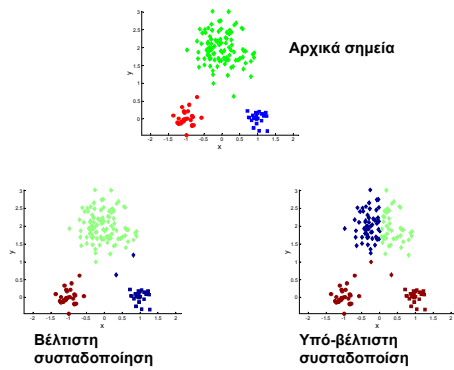
Ένας τρόπος να βελτιώσουμε τη συσταδοποίηση (ελάττωση του SSE) είναι να μεγαλώσουμε το K

Αλλά γενικά μια καλή συσταδοποίηση με μικρό K μπορεί να έχει μικρότερο SSE από μια κακή συσταδοποίηση με μεγάλο K

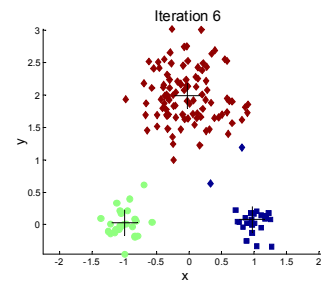
K-means: Βασικός Αλγόριθμος

- Το αποτέλεσμα εξαρτάται από την επιλογή των αρχικών σημείων

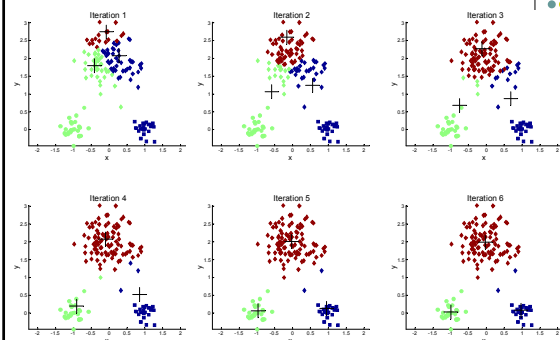
K-means: Παράδειγμα



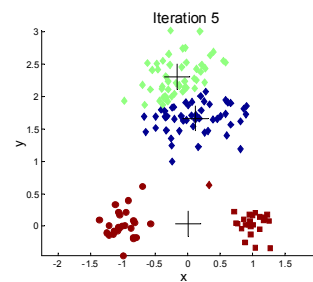
K-means: Επιλογή αρχικών σημείων

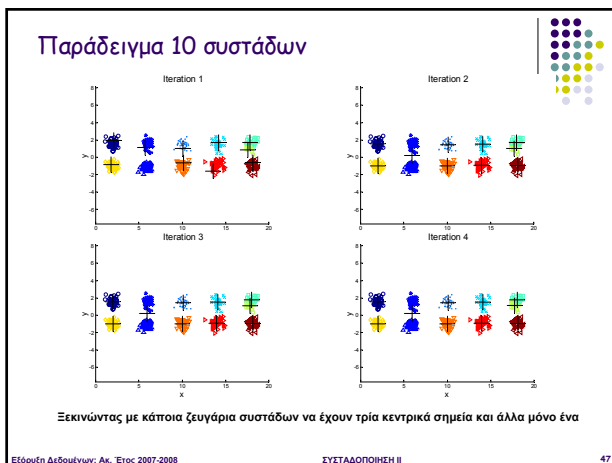
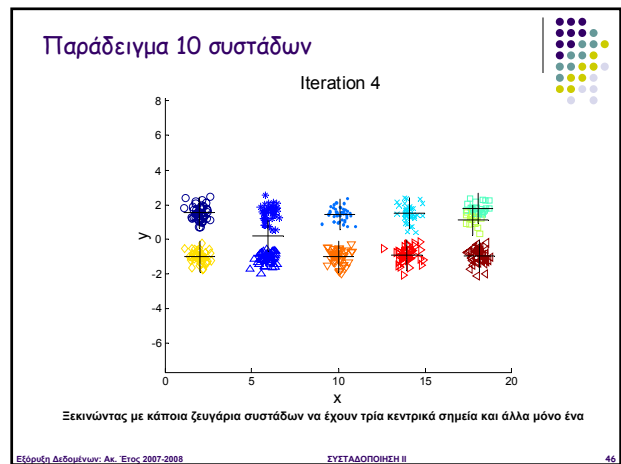
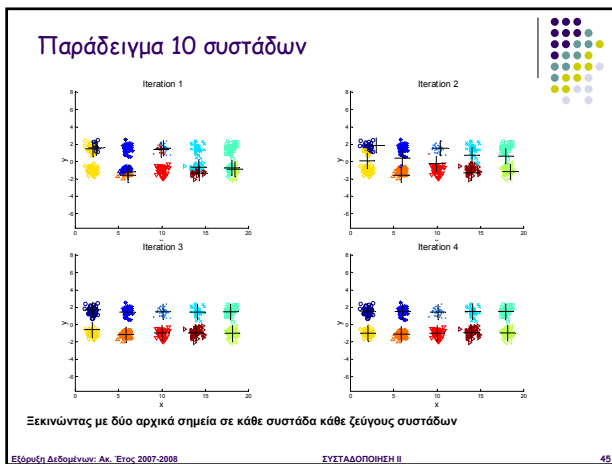
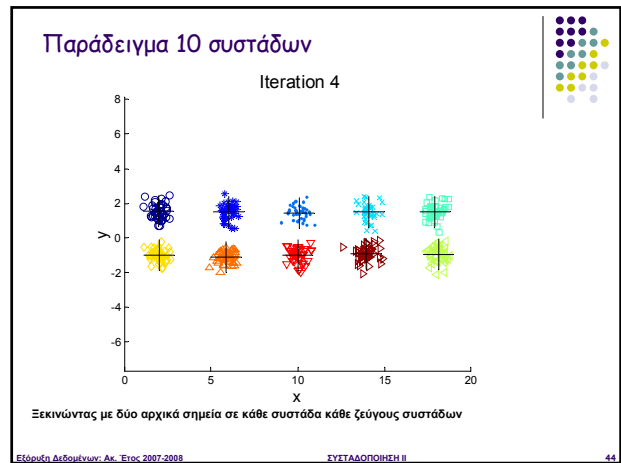
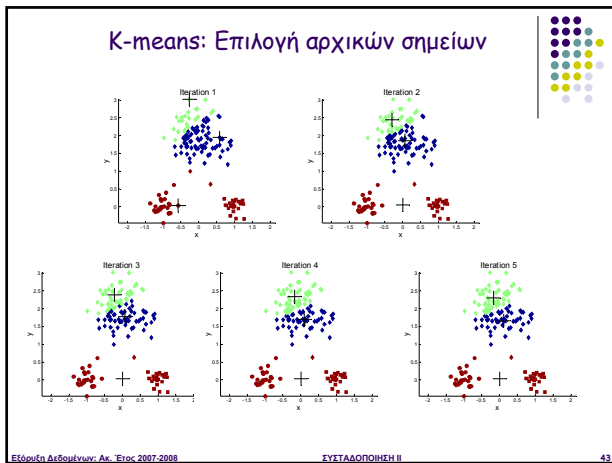


K-means: Επιλογή αρχικών σημείων



K-means: Επιλογή αρχικών σημείων





K-means: Λύσεις για την επιλογή αρχικών σημείων

Πολλά τρεξίματα
Βοηθά, αλλά πολλές περιπτώσεις

Δειγματοληψία και χρήση κάποιας ιεραρχικής τεχνικής

Επιλογή παραπάνω από k αρχικών σημείων και μετά επιλογή k από αυτά τα αρχικά κεντρικά σημεία (πχ τα πιο απομακρυσμένα μεταξύ τους)

Σταδιακή επιλογή
Επιλογή του πρώτου σημείου τυχαία ή ως το μέσο όλων των σημείων
Για καθένα από τα υπόλοιπα αρχικά σημεία επέλεξε αυτό που είναι πιο μακριά από τα μέχρι τώρα επιλεγμένα αρχικά σημεία

- Μπορεί να οδηγήσει στην επιλογή outliers
- Ο υπολογισμός του πιο απομακρυσμένου σημείου είναι δαπανηρός
- Συχνά εφαρμόζεται σε δείγματα

Εξομολογή Διδασκόντων: Ακ. Έτος 2007-2008 ΣΥΓΓΡΑΦΟΠΟΙΗΣΗ II 48

K-means: Άδειες συστάδες

Ο βασικός αλγόριθμος μπορεί να οδηγήσει σε **άδειες αρχικές συστάδες**

Πολλές στρατηγικές

Επιλογή του σημείου που είναι πιο μακριά από όλα τα τωρινά κέντρα = επιλογή του σημείου που συμβάλει περισσότερο στο SSE

Ένα σημείο από τη συστάδα με το υψηλότερο SSE - θα οδηγήσει σε «σπάσιμο» της άρα σε μείωση του λάθους

Αν πολλές *άδειες συστάδες*, τα παραπάνω βήματα μπορεί να επαναληφθούν πολλές φορές

K-means: Σταδιακή ενημέρωση κεντρικών σημείων

Στο βασικό K-means, το κέντρα ενημερώνεται αφού όλο τα σημεία έχουν ανατεθεί στο κέντρο

Μια παραλλαγή είναι να ενημερώνονται τα κέντρα μετά από κάθε ανάθεση (incremental approach)

- Κάθε ανάθεση ενημερώνει 0 ή 2 κέντρα
- Πιο δαπανηρό
- Έχει σημασία η σειρά εισαγωγής/εξέτασης των σημείων
- Δεν υπάρχουν άδειες συστάδες
- Μπορεί να χρησιμοποιηθούν βάρη - αν υπάρχει κάποια τυχαία αντικειμενική συνάρτηση - έλεγχος τι συμφέρει κάθε φορά

Προ και Μετα Επεξεργασία

Ολικό SSE και SSE Συστάδας

Προ-επεξεργασία

Κανονικοποίηση των δεδομένων
Απομάκρυνση outliers

Post-processing

Split-Merge (διατηρώντας το ίδιο K)

Διαχωρισμός (split) συστάδων με το σχετικά μεγαλύτερο SSE
Δημιουργία μια νέας συστάδας: πχ επιλέγοντας το σημείο που είναι πιο μακριά από όλα τα κέντρα ή τυχαία επιλογή σημείου ή επιλογή του σημείου με το μεγαλύτερο SSE

Συνένωση (merge) συστάδων που είναι σχετικά κοντινές (τα κέντρα τους έχουν την μικρότερη απόσταση) ή τις δυο συστάδες που οδηγούν στην μικρότερη αύξηση του SSE

Διαγραφή συστάδας και ανακατανομή των σημείων της σε άλλες συστάδες (αυτό που οδηγεί στην μικρότερη αύξηση του SSE)

K-means με διχοτόμηση (bisecting k-means)

Παραλλαγή που μπορεί να παράγει μια διαχωριστική ή ιεραρχική συσταδοποίηση

- 1: Αρχικοποίηση της λίστας των συστάδων ώστε να περιέχει μια συστάδα που περιέχει όλα τα σημεία
- 2: **Repeat**
- 3: Επιλογή μιας συστάδας από τη λίστα των συστάδων
- 4: **for** $i = 1$ to $number_of_trials$ **do**
- 5: διχοτόμισε την επιλεγμένη συστάδα χρησιμοποιώντας το βασικό k-means
- 6: Πρόσθεσε στη λίστα από τις δυο συστάδες που προέκυψαν από τη διχοτόμηση αυτήν με το μικρότερο SSE
- 5: **Until** η λίστα των συστάδων να έχει K συστάδες

K-means με διχοτόμηση (bisecting k-means)

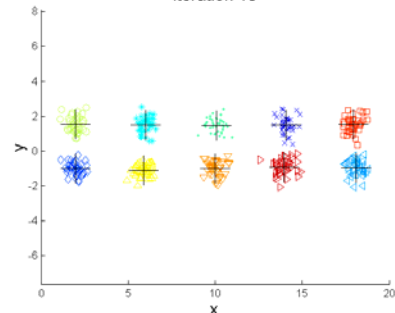
Ποια συστάδα να διασπάσουμε:

- Τη μεγαλύτερη
- Αυτή με το μεγαλύτερο SSE
- Συνδυασμό των παραπάνω

Μπορεί να χρησιμοποιηθεί και ως ιεραρχικός

K-means με διχοτόμηση

Iteration 10



K-means: Περιορισμοί

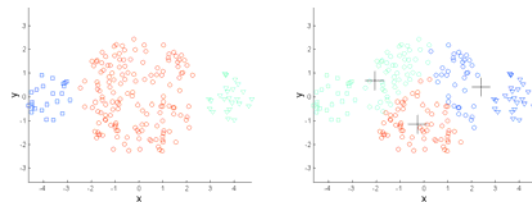
Ο K-means έχει προβλήματα όταν οι συστάδες έχουν διαφορετικά

Διαφορετικά Μεγέθη
Διαφορετικές Πυκνότητες
Non-globular shapes

Έχει προβλήματα όταν τα δεδομένα έχουν outliers



K-means: Περιορισμοί - διαφορετικά μεγέθη

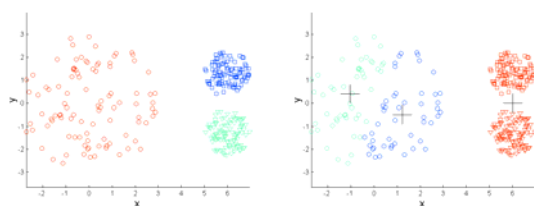


Αρχικά σημεία

K-means (3 συστάδες)

Δεν μπορεί να βρει το μεγάλο κόκκινο, γιατί είναι πολύ μεγαλύτερος από τους άλλους

K-means: Περιορισμοί - διαφορετικές πυκνότητες

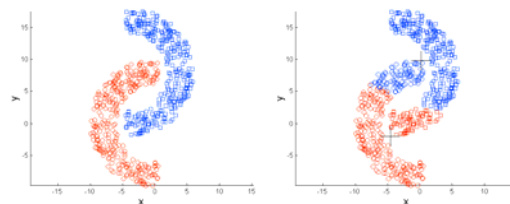


Αρχικά σημεία

K-means (3 συστάδες)

Δεν μπορεί να διαχωρίσει τους δύο μικρούς γιατί είναι πολύ πυκνοί σε σχέση με τον ένα μεγάλο

K-means: Περιορισμοί - μη κυκλικά σχήματα

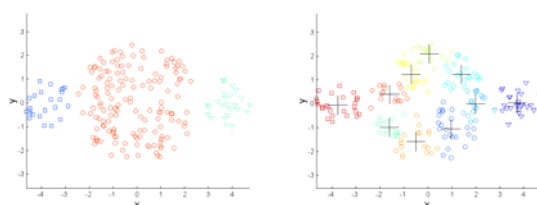


Αρχικά σημεία

K-means (2 συστάδες)

Δεν μπορεί να βρει τις δύο συστάδες γιατί έχουν μη κυκλικά σχήματα

K-means: Περιορισμοί

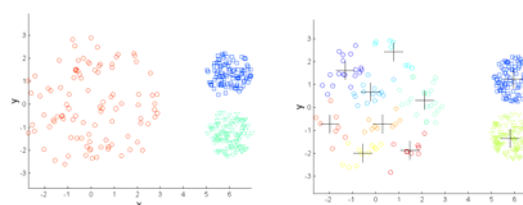


Αρχικά Σημεία

K-means Συστάδες

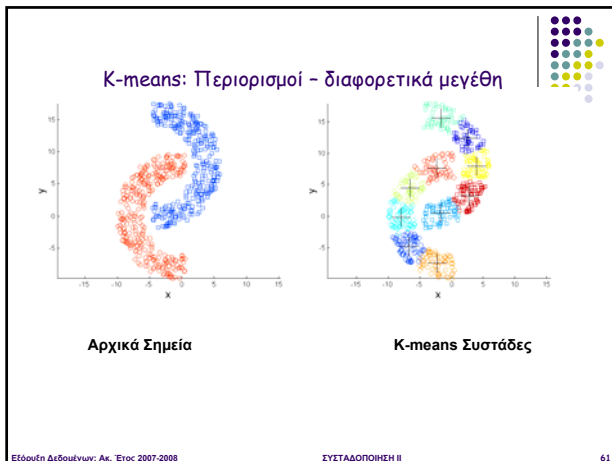
Μια λύση είναι να χρησιμοποιηθούν πολλές συστάδες
Βρίσκει τμήματα των συστάδων, αλλά πρέπει να τα συγκεντρώσουμε

K-means: Περιορισμοί



Αρχικά σημεία

K-means Συστάδες



K-means: Επιλογή αρχικών σημείων

Αν υπάρχουν K «πραγματικές συστάδες» η πιθανότητα να επιλέξουμε ένα κέντρο από κάθε συστάδα είναι μικρή, συγκεκριμένα αν όλες οι συστάδες έχουν το ίδιο μέγεθος n , τότε:

$$P = \frac{\text{number of ways to select one centroid from each cluster}}{\text{number of ways to select } K \text{ centroids}} = \frac{K!n^K}{(Kn)^K} = \frac{K!}{K^K}$$

Για παράδειγμα, αν $K = 10$, η πιθανότητα είναι $= 10!/10^{10} = 0.00036$

Μερικές φορές τα αρχικά σημεία βελτιώνουν τη θέση τους και άλλες φορές όχι

Θα δούμε ένα παράδειγμα με 5 ζευγάρια συστάδων

Εξώφυλλο Διδασκόντων: Ακ. Έτος 2007-2008 ΣΥΓΧΡΟΝΗ ΣΤΑΤΙΣΤΙΚΗ II 62

K-medoid

Συνήθως συνεχής d-διάστατο χώρο

Διαλέγει ένα αντιπροσωπευτικό σημείο από τα δεδομένα και ελαχιστοποιεί την απόσταση από αυτό - Medoid: το πιο κεντρικό σημείο της συστάδας (αντί να χρησιμοποιεί το mean)

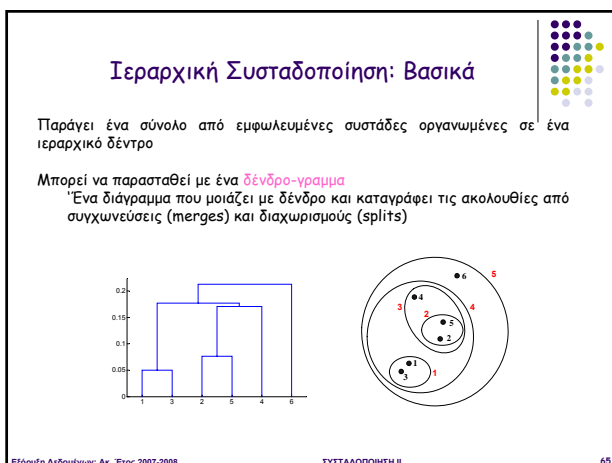
Μειώνει την ευαισθησία σε outliers

Μπορεί να εφαρμοστεί σε δεδομένα οποιουδήποτε τύπου (πχ και για κατηγορικά δεδομένα)

Εξώφυλλο Διδασκόντων: Ακ. Έτος 2007-2008 ΣΥΓΧΡΟΝΗ ΣΤΑΤΙΣΤΙΚΗ II 63

Ιεραρχική Συσταδοποίηση

Εξώφυλλο Διδασκόντων: Ακ. Έτος 2007-2008 ΣΥΓΧΡΟΝΗ ΣΤΑΤΙΣΤΙΚΗ II 64



Ιεραρχική Συσταδοποίηση: Πλεονεκτήματα

- Δε χρειάζεται να υποθέσουμε ένα συγκεκριμένο αριθμό από συστάδες
- Οποιοσδήποτε επιθυμητός αριθμός από συστάδες μπορεί να επιτευχθεί κόβοντας το δενδρόγραμμα στο κατάλληλο επίπεδο
- Μπορεί να αντιστοιχούν σε λογικές ταξινομήσεις

Για παράδειγμα στις βιολογικές επιστήμες (ζωικό βασίλειο, phylogeny reconstruction, ...)

Εξώφυλλο Διδασκόντων: Ακ. Έτος 2007-2008 ΣΥΓΧΡΟΝΗ ΣΤΑΤΙΣΤΙΚΗ II 66

Ιεραρχική Συσταδοποίηση

Διο βασικοί τύποι ιεραρχικής συσταδοποίησης

▪ Συσσωρευτικός (Agglomerative):

- Αρχίζει με τα σημεία ως ξεχωριστές συστάδες
- Σε κάθε βήμα, συγχωνεύει το πιο κοντινό ζευγάρι συστάδων μέχρι να μείνει μόνο μία (ή k) συστάδες

▪ Διαιρετικός (Divisive):

- Αρχίζει με μία συστάδα που περιέχει όλα τα σημεία
- Σε κάθε βήμα, διαχωρίζει μία συστάδα, έως κάθε συστάδα να περιέχει μόνο ένα σημείο (ή να δημιουργηθούν k συστάδες)

Ιεραρχική Συσταδοποίηση

Οι παραδοσιακοί αλγόριθμοι

- χρησιμοποιούν έναν πίνακα ομοιότητα ή απόστασης
- διαχωρισμός ή συγχώνευση μιας ομάδας τη φορά

Συσσωρευτική Ιεραρχική Συσταδοποίηση (ΣΙΣ)

Η πιο δημοφιλής τεχνική συσταδοποίησης

Βασικός Αλγόριθμος

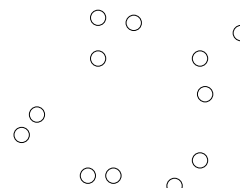
- 1: Υπολογισμός του Πίνακα Γειτνίασης
- 2: Έστω κάθε σημείο αποτελεί και μια συστάδα
- 3: **Repeat**
- 4: Συγχώνευση των δύο κοντινότερων συστάδων
- 5: Ενημέρωση του Πίνακα Γειτνίασης
- 6: **Until** να μείνει μία μόνο συστάδα

Βασική λειτουργία είναι ο υπολογισμός της γειτνίασης δυο συστάδων

Διαφορετικοί αλγόριθμοι με βάση το πως ορίζεται η απόσταση ανάμεσα σε δύο συστάδες

Συσσωρευτική Ιεραρχική Συσταδοποίηση

Αρχικά: Κάθε σημείο και συστάδα και ένας Πίνακας Γειτνίασης (proximity matrix)



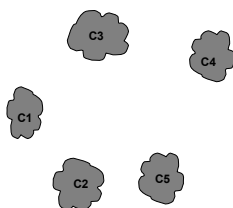
	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
...						

Πίνακας Γειτνίασης



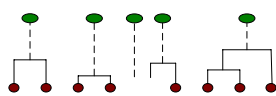
Συσσωρευτική Ιεραρχική Συσταδοποίηση

Μετά από κάποιες συγχωνεύσεις, έχουμε κάποιες συστάδες



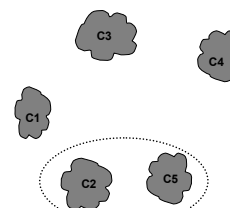
	C1	C2	C3	C4	C5
C1					
C2					
C3					
C4					
C5					

Πίνακας Γειτνίασης



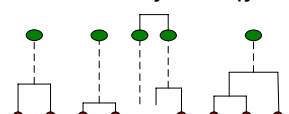
Συσσωρευτική Ιεραρχική Συσταδοποίηση

Θέλουμε να συγχωνεύσουμε τις δύο κοντινότερες συστάδες (C2 και C5) και να ενημερώσουμε τον πίνακα γειτνίασης.



	C1	C2	C3	C4	C5
C1					
C2					
C3					
C4					
C5					

Πίνακας Γειτνίασης



Συσσωρευτική Ιεραρχική Συσταδοποίηση

Μετά τη συγχώνευση η ερώτηση είναι: Πώς ενημερώνουμε τον πίνακα γειννιάσης

	C1	C2 U C5	C3	C4
C1	?			
C2 U C5	?	?	?	?
C3	?			
C4	?			

Πίνακας Γειννιάσης

Εύρωδη Διδασκννν: Ακ. Έτος 2007-2008 ΣΥΣΤΑΔΟΠΟΙΗΣΗ II 73

ΣΙΣ: Ορισμός απόστασης μεταξύ συστάδων

	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						

Πίνακας Γειννιάσης

- MIN
- MAX
- Μέσος όρος της συστάδας
- Η απόσταση μεταξύ των κεντρικών σημείων
- Άλλες μέθοδοι βασισμένες σε μια αντικειμενική συνάρτηση
 - Η μέθοδος του Ward χρησιμοποιεί τετραγωνικά λάθη

Εύρωδη Διδασκννν: Ακ. Έτος 2007-2008 ΣΥΣΤΑΔΟΠΟΙΗΣΗ II 74

ΣΙΣ: Ορισμός απόστασης μεταξύ συστάδων

	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						

Πίνακας Γειννιάσης

- MIN
- MAX
- Μέσος όρος της ομάδας
- Η απόσταση μεταξύ των κεντρικών σημείων
- Άλλες μέθοδοι βασισμένες σε μια αντικειμενική συνάρτηση
 - Η μέθοδος του Ward χρησιμοποιεί τετραγωνικά λάθη

Εύρωδη Διδασκννν: Ακ. Έτος 2007-2008 ΣΥΣΤΑΔΟΠΟΙΗΣΗ II 75

ΣΙΣ: Ορισμός απόστασης μεταξύ συστάδων: MIN

MIN ή μοναδικής ακμής ή απλού συνδέσμου (single link)

Η ομοιότητα μεταξύ δυο συστάδων βασίζεται στα δυο πιο όμοια (πιο γειτονικά) σημεία στις διαφορετικές συστάδες (με όρους γραφημάτων - shortest edge)

Καθορίζεται από ένα ζεύγος τιμών, δηλαδή **μια ακμή (link)** του γραφήματος γειννιάσης.

Ονομάζεται και μέθοδος συσταδοποίησης **κοντινότερου γείτονα**

Εύρωδη Διδασκννν: Ακ. Έτος 2007-2008 ΣΥΣΤΑΔΟΠΟΙΗΣΗ II 76

p1 p2 p3 p4 p5 p10 p11 p12

ΣΙΣ: Ορισμός απόστασης μεταξύ συστάδων: MIN

MIN ή μοναδικής ακμής ή απλού συνδέσμου (single link)

Η ομοιότητα μεταξύ δυο συστάδων βασίζεται στα δυο πιο όμοια (πιο γειτονικά) σημεία στις διαφορετικές συστάδες (με όρους γραφημάτων - shortest edge)

Καθορίζεται από ένα ζεύγος τιμών, δηλαδή **μια ακμή (link)** του γραφήματος γειννιάσης.

	I1	I2	I3	I4	I5
I1	1,00	0,90	0,10	0,65	0,20
I2	0,90	1,00	0,70	0,60	0,50
I3	0,10	0,70	1,00	0,40	0,30
I4	0,65	0,60	0,40	1,00	0,80
I5	0,20	0,50	0,30	0,80	1,00

Προσοχή: ομοιότητα!!

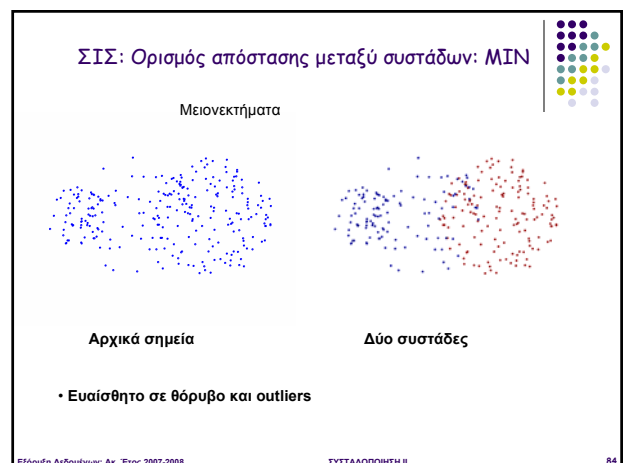
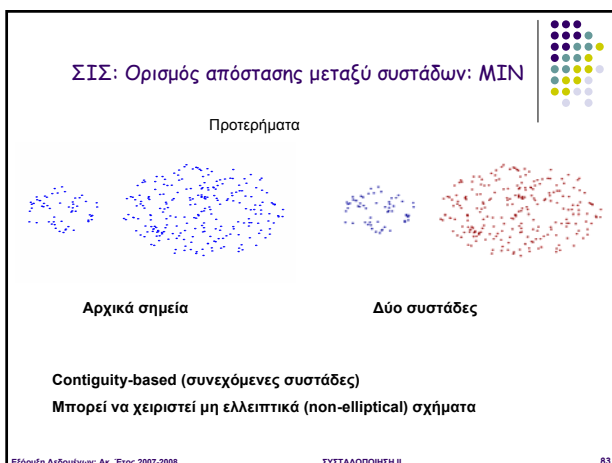
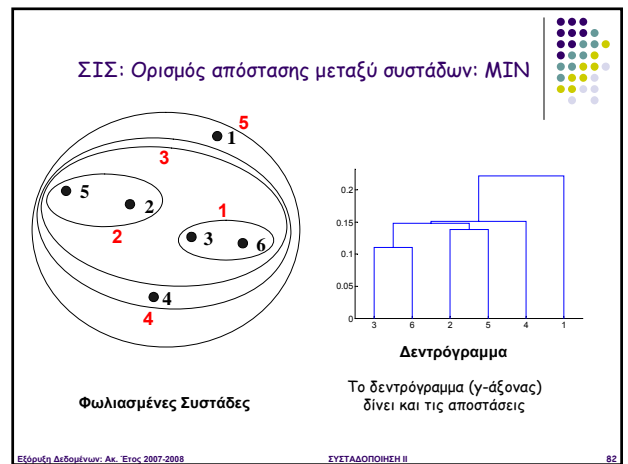
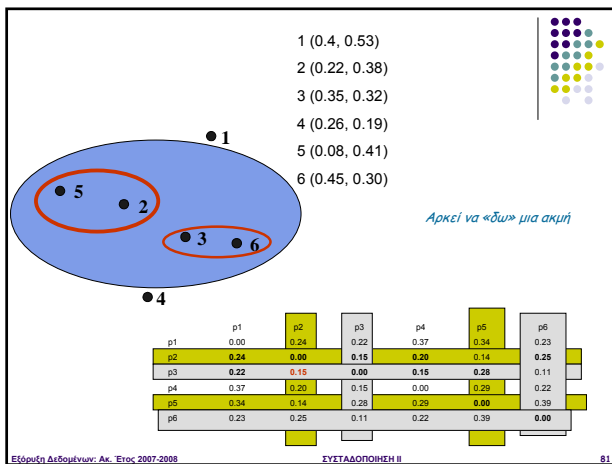
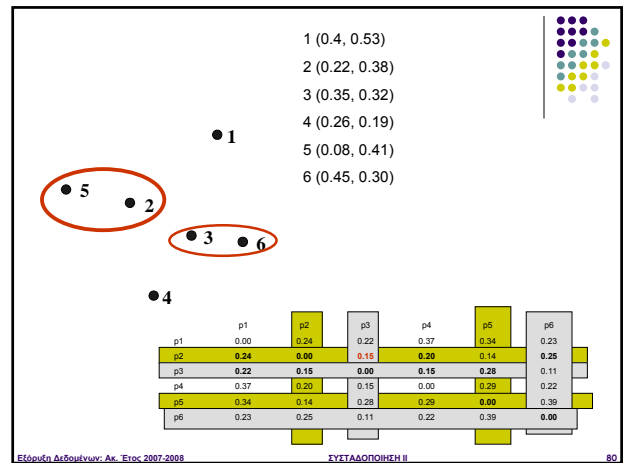
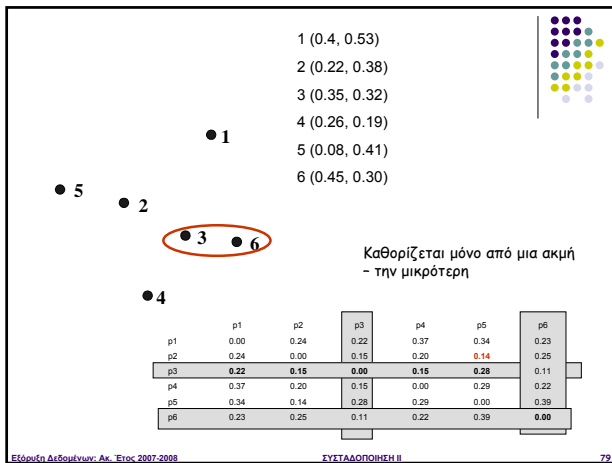
Εύρωδη Διδασκννν: Ακ. Έτος 2007-2008 ΣΥΣΤΑΔΟΠΟΙΗΣΗ II 77

1 (0,4, 0,53)
2 (0,22, 0,38)
3 (0,35, 0,32)
4 (0,26, 0,19)
5 (0,08, 0,41)
6 (0,45, 0,30)

Πίνακας απόστασης

	p1	p2	p3	p4	p5	p6
p1	0,00	0,24	0,22	0,37	0,34	0,23
p2	0,24	0,00	0,15	0,20	0,14	0,25
p3	0,22	0,15	0,00	0,15	0,28	0,11
p4	0,37	0,20	0,15	0,00	0,29	0,22
p5	0,34	0,14	0,28	0,29	0,00	0,39
p6	0,23	0,25	0,11	0,22	0,39	0,00

Εύρωδη Διδασκννν: Ακ. Έτος 2007-2008 ΣΥΣΤΑΔΟΠΟΙΗΣΗ II 78



ΣΙΣ: Ορισμός απόστασης μεταξύ συστάδων

	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						

· Πίνακας Γεινιάσης

- MIN
- MAX**
- Μέσος όρος της ομάδας
- Η απόσταση μεταξύ των κεντρικών σημείων
- Άλλες μέθοδοι βασισμένες σε μια αντικειμενική συνάρτηση
 - Η μέθοδος του Ward χρησιμοποιεί τετραγωνικά λάθη

Εξόρυξη Δεδομένων: Ακ. Έτος 2007-2008 ΣΥΓΓΡΑΦΟΠΟΙΗΣΗ II 85

ΣΙΣ: Ορισμός απόστασης μεταξύ συστάδων: MAX

MAX ή πλήρους συνδεσιμότητας (complete linkage)
 - Αναζητά κλίκες
 Η ομοιότητα μεταξύ δυο συστάδων βασίζεται στα δυο λιγότερο όμοια (πιο μακρινά) σημεία στις διαφορετικές συστάδες (longest edge)

Καθορίζεται από **όλα τα ζεύγη τιμών** στις δύο συστάδες.

	11	12	13	14	15
11	1.00	0.90	0.10	0.65	0.20
12	0.90	1.00	0.70	0.60	0.50
13	0.10	0.70	1.00	0.40	0.30
14	0.65	0.60	0.40	1.00	0.80
15	0.20	0.50	0.30	0.80	1.00

ομοιότητα

Εξόρυξη Δεδομένων: Ακ. Έτος 2007-2008 ΣΥΓΓΡΑΦΟΠΟΙΗΣΗ II 86

	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00

Εξόρυξη Δεδομένων: Ακ. Έτος 2007-2008 ΣΥΓΓΡΑΦΟΠΟΙΗΣΗ II 87

	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00

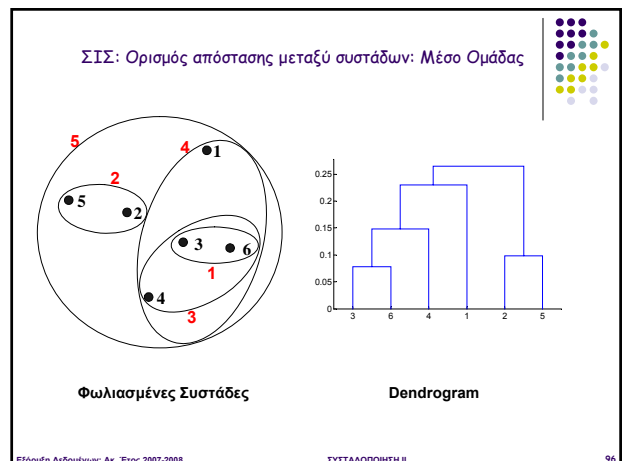
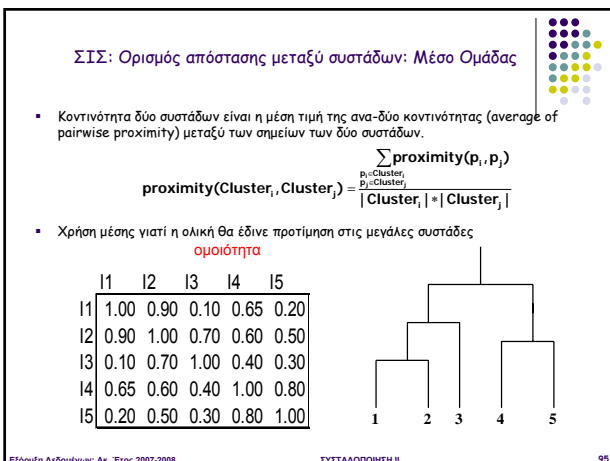
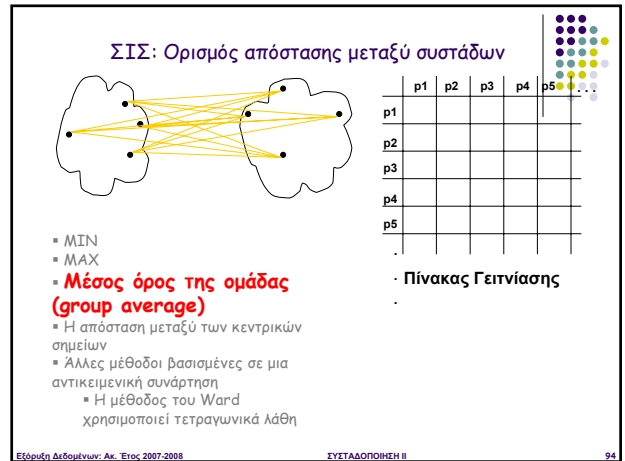
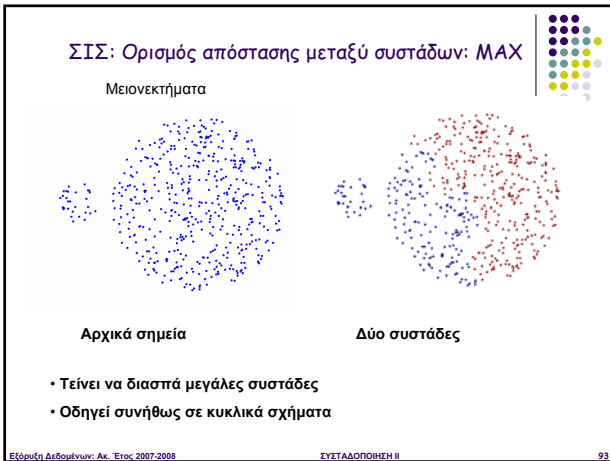
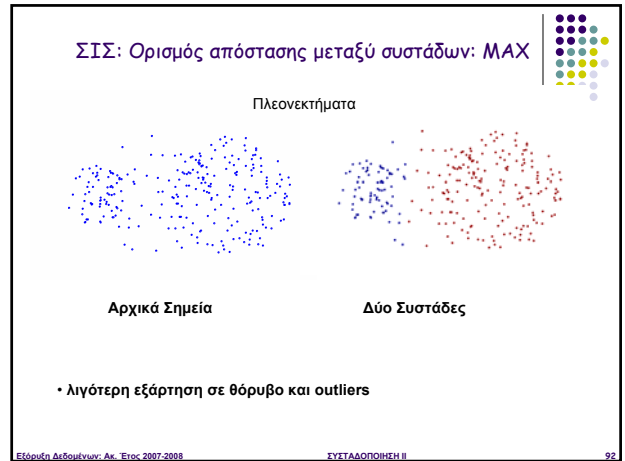
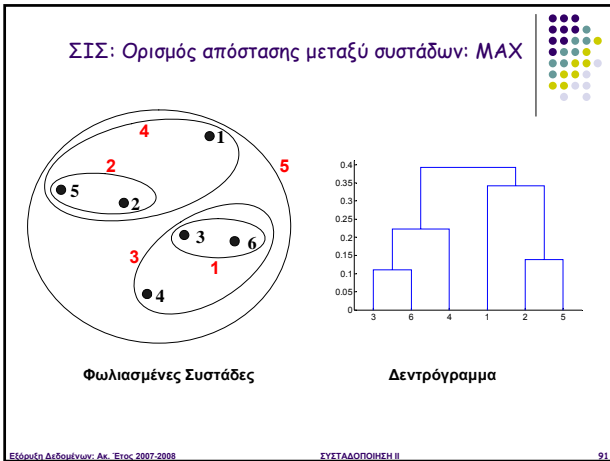
Εξόρυξη Δεδομένων: Ακ. Έτος 2007-2008 ΣΥΓΓΡΑΦΟΠΟΙΗΣΗ II 88

	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00

Εξόρυξη Δεδομένων: Ακ. Έτος 2007-2008 ΣΥΓΓΡΑΦΟΠΟΙΗΣΗ II 89

	p1	p2	p3	p4	p5	p6
p1	0.00	0.24	0.22	0.37	0.34	0.23
p2	0.24	0.00	0.15	0.20	0.14	0.25
p3	0.22	0.15	0.00	0.15	0.28	0.11
p4	0.37	0.20	0.15	0.00	0.29	0.22
p5	0.34	0.14	0.28	0.29	0.00	0.39
p6	0.23	0.25	0.11	0.22	0.39	0.00

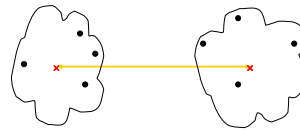
Εξόρυξη Δεδομένων: Ακ. Έτος 2007-2008 ΣΥΓΓΡΑΦΟΠΟΙΗΣΗ II 90



ΣΙΣ: Ορισμός απόστασης μεταξύ συστάδων: Μέσο Ομάδας

- Ανάμεσα σε MIN-MAX
- Πλεονεκτήματα: μικρότερη ευαισθησία σε θόρυβο και outliers
- Μειονεκτήματα: Ευνοεί κυκλικές συστάδες

ΣΙΣ: Ορισμός απόστασης μεταξύ συστάδων



	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						

- MIN
- MAX
- Μέσος όρος της ομάδας
- **Η απόσταση μεταξύ των κεντρικών σημείων**
- Άλλες μέθοδοι βασισμένες σε μια αντικειμενική συνάρτηση
 - Η μέθοδος του Ward χρησιμοποιεί τετραγωνικά λάθη

Πίνακας Γειτνίασης

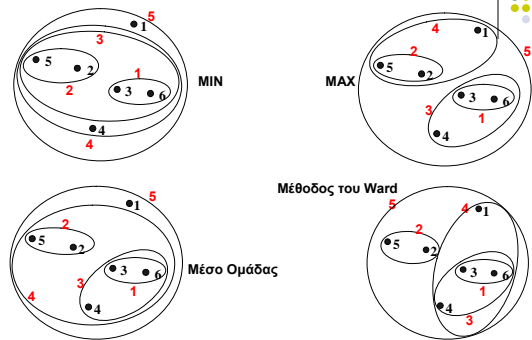
Πρόβλημα: μη μονότονη αύξηση της απόστασης

Δηλαδή, δύο συστάδες που συγχωνεύονται μπορεί να έχουν μικρότερη απόσταση από συστάδες που έχουν συγχωνευτεί σε προηγούμενα βήματα

ΣΙΣ: Ορισμός απόστασης μεταξύ συστάδων: Μέθοδος του Ward

- Βασισμένο στην αύξηση του SSE όταν συγχωνεύονται οι δύο συστάδες
- Ιεραρχικό ανάλογο του k-means
- Μπορεί να χρησιμοποιηθεί για την αρχικοποίηση του k-means

ΣΙΣ: Ορισμός απόστασης μεταξύ συστάδων: Σύγκριση



ΣΙΣ: Πολυπλοκότητα Χρόνου και Χώρου

- $O(m^2)$ χώρος για την αποθήκευση του πίνακα γειτνίασης
 - m αριθμός σημείων.
- $O(m^3)$
 - Ξεκινάμε με m συστάδες και μειώνουμε 1 τη φορά
 - Αν γραμμική αναζήτηση του πίνακα $O(m^2)$
 - Καλύτερος χρόνος αν διατηρούμε κάποια ταξινόμηση των αποστάσεων πχ heap

ΣΙΣ: Περιορισμοί και Προβλήματα

Οι αποφάσεις είναι τελικές - αφού δύο συστάδες συγχωνευτούν αυτό δεν μπορεί να αλλάξει

Δεν ελαχιστοποιούν άμεσα κάποια αντικειμενική συνάρτηση



Μια διαιρετική παραλλαγή του ΜΙΝ βασίζεται σε spanning tree (σκελετικά δέντρα)

1. Χρησιμοποίησε τον πίνακα απόστασης και κατασκεύασε ένα ελάχιστο σκελετικό δέντρο
2. Δημιούργησε μια νέα συστάδα «σπάζοντας» το δέντρο στην ακμή με τη μεγαλύτερη απόσταση (μικρότερη ομοιότητα)