

Εισαγωγή I

Ανάλυση Συστάδων - Εισαγωγικές Έννοιες

- Στην Διαχωριστική Ανάλυση (γνωστή επίσης και ως supervised pattern recognition), είναι εκ των προτέρων γνωστό ότι έχουμε k ομάδες, στις οποίες θέλουμε να κατατάξουμε κάθε νέα παρατήρηση.
- Στην Ανάλυση Συστάδων (Cluster Analysis), το σύνολο των παρατηρήσεων αποτελεί μια (μεγάλη) ομάδα και το στόχος είναι να σηματίσουμε τις κατάλληλες ομάδες (unsupervised pattern recognition)
- Οι ομάδες αυτές θα πρέπει να περιέχουν 'όμοια' αντικείμενα (ή να έχουν παρόμοιες ιδιότητες).

Εισαγωγή II

Ανάλυση Συστάδων - Εισαγωγικές Έννοιες

- Όταν αναφερόμαστε στη δημιουργία συστάδων, εννοούμε τη δημιουργία k ομάδων με βάση X_1, X_2, \dots, X_n τυχαία διανύσματα (το διαθέσιμο dataset), με την i -οστή ομάδα να περιλαμβάνει n_i από τα n διανύσματα, έτσι ώστε

$$(n_1) + (n_2) + \dots + (n_k) = \underline{n}.$$

με τα στοιχεία της ομάδας να είναι 'όμοια' (within homogeneity) αλλά οι ομάδες μεταξύ τους 'ανόμοιες'.

- Επιθυμητή είναι η μικρή μεταβλητότητα εντός των ομάδων (within groups variability) και η μεγάλη μεταβλητότητα μεταξύ των ομάδων (between groups variability).

Εισαγωγή III

Ανάλυση Συστάδων - Εισαγωγικές Έννοιες

- **Προβλήματα:** (α) Πώς θα μετρήσουμε την ‘ομοιότητα’; (β) πώς θα μετρήσουμε την ‘ανομοιότητα’ και την ‘απόσταση’ μεταξύ των ομάδων; (γ) ποια είναι η ελάχιστη απόσταση για το διαχωρισμό δύο ομάδων; (δ) υπάρχει κάποιο κριτήριο ‘συσταδοποίησης’ (clustering), αφού (όπως θα δούμε και στη συνέχεια) διαφορετικά κριτήρια δίνουν διαφορετικές συστάδες με τα ίδια χαρακτηριστικά;
- Ένα ακόμη πρόβλημα είναι ο αριθμός των συστάδων που πρέπει να προσδιοριστούν. Μπορεί είτε να είναι προκαθορισμένος είτε να καθοριστεί μέσω κάποιου κριτηρίου (π.χ. minimum between-clusters distance)
- Τέλος, αν τα δεδομένα (αντικείμενα) θεωρηθούν ως παρατηρήσεις από μια πεπερασμένη μίξη κατανομών, π.χ. $f(\mathbf{x}) = \sum_{h=1}^K q_h f_h(\mathbf{x})$, τότε εκτός από τους συντελεστές q_1, q_2, \dots, q_K πρέπει να προσδιοριστεί και το K .

Παραδείγματα I

- Έστω η συνήθης τράπουλα (με τα 52 φύλλα). Τότε μπορούμε να ορίσουμε:
 - ▶ 13 συστάδες με βάση τα διαφορετικά φύλλα A, 1, 2, ..., J, Q, K
 - ▶ 4 συστάδες με βάση το σχέδιο (κούπα, σπαθί, μπαστούνι, καρό)
 - ▶ 2 συστάδες με βάση το χρώμα (μαύρο και κόκκινο) κλπ
- Έστω ότι έχουμε 11 γλώσσες στις οποίες χρησιμοποιείται το Λατινικό αλφάβητο. Η ομαδοποίηση των γλωσσών μπορεί να γίνει με βάση το αν οι λέξεις που αντιστοιχούν στους αριθμούς 1, 2, ..., 10 ξεκινούν με το ίδιο γράμμα (Παράδειγμα 12.3, Johnson and Wicher 2007).
- Σε αυτή την περίπτωση, η ομαδοποίηση δίνει **Ομάδα 1**: Αγγλικά, Νορβηγικά, Δανέζικα, Ολλανδικά, Γερμανικά. **Ομάδα 2**: Γαλλικά, Ιταλικά, Ισπανικά, Πολωνικά. **Ομάδα 3**: Ουγγρικά. **Ομάδα 4**: Φινλανδέζικα (Εφαρμογή στο Εργαστήριο).

Παραδείγματα II

- Ως μέτρο ομοιότητας έχει χρησιμοποιηθεί η ‘συμφωνία’ (concordance) με βάση το παραπάνω κριτήριο. Ανάλογα όμως κριτήρια μπορούν να χρησιμοποιηθούν και για άλλες λέξεις στα διαφορετικά αλφάβητα, οδηγώντας σε διαφορετική συσταδοποίηση.
- Τα παραπάνω παραδείγματα δείχνουν το πόσο σημαντική είναι η επιλογή του κριτηρίου συσταδοποίησης.

Εφαρμογές Ανάλυσης Συστάδων I

- **Συμπίεση Δεδομένων (Data Compression)**: Ομαδοποίηση των δεδομένων σε k συστάδες. Από κάθε συστάδα επιλέγω έναν αντιπρόσωπο.
- **Δημιουργία Υποθέσεων (Hypothesis Generation)**: Χρήση συστάδων για δημιουργία ερευνητικών υποθέσεων, τις οποίες πρόκειται να ελέγξουμε και να επιβεβαιώσουμε χρησιμοποιώντας διάφορα *training sets*.
- **Έλεγχος Υποθέσεων**
- **Συσταδοποίηση - Διαχωρισμότητα (Classification - Discrimination)**
- **Πρόβλεψη**: Μπορεί να βασιστεί σε συστάδες, λαμβάνοντας υπόψη τα κύρια χαρακτηριστικά των μελών της. Για δεδομένες τιμές ενός ατόμου (π.χ. διάνυσμα τιμών x), μπορούμε να το ταξινομήσουμε σε μια από τις διαθέσιμες συστάδες.

Εφαρμογές Ανάλυσης Συστάδων II

- **Παράδειγμα:** Ασθενείς που έχουν προσβληθεί από την ίδια ασθένεια, μπορούν να δώσουν συστάδες με βάση την αντίδρασή τους σε συγκεκριμένη θεραπεία. Στη συνέχεια, για κάθε νέο ασθενή, αναγνωρίζοντας σε ποια συστάδα ανήκει (διάγνωση) μπορεί να προταθεί η καταλληλότερη θεραπεία.
- Ανάλογες εφαρμογές μπορεί κανείς να βρει στις περιοχές της Εκπαίδευσης, του Management, της Ψυχολογίας, των Τηλεπικοινωνιών και της Αρχαιολογίας.

↓ Taylor made for each case.

Αλγόριθμοι Συσταδοποίησης I

- **Στόχος:** Η διαμέριση του διαθέσιμου dataset X , το οποίο αποτελείται από n διανύσματα x_1, \dots, x_n διαστάσεων $p \times 1$ το καθένα, σε k μη-κενές συστάδες μεγέθους $n_1 > 0, \dots, n_k > 0$.
- Οι δυνατοί τρόποι για να γίνει αυτό αντιστοιχούν στο πλήθος των ακέραιων λύσεων της εξίσωσης $n_1 + \dots + n_k = n$ και ισούνται με

$n = \text{άξια}$
 $k = \text{ομάδες}$

$$S(n, k) = \frac{1}{k!} \sum_{i=0}^n (-1)^{k-i} \binom{n}{i} i^n,$$

ο οποίος είναι γνωστός και ως ο *Αριθμός Stirling* 2ου είδους.

- Ενδεικτικά, $S(\overset{n}{16}, \overset{k}{2}) = \underline{32767}$ (οι τρόποι για να διαμερίσουμε τα 16 χαρτιά σε 2 ομάδες διαφορετικών μεγεθών) και $S(\overset{n}{16}, \overset{k}{3}) = \underline{7141686}$ (οι τρόποι για να διαμερίσουμε τα 16 χαρτιά σε 3 ομάδες διαφορετικών μεγεθών).

Αλγόριθμοι Συσταδοποίησης II

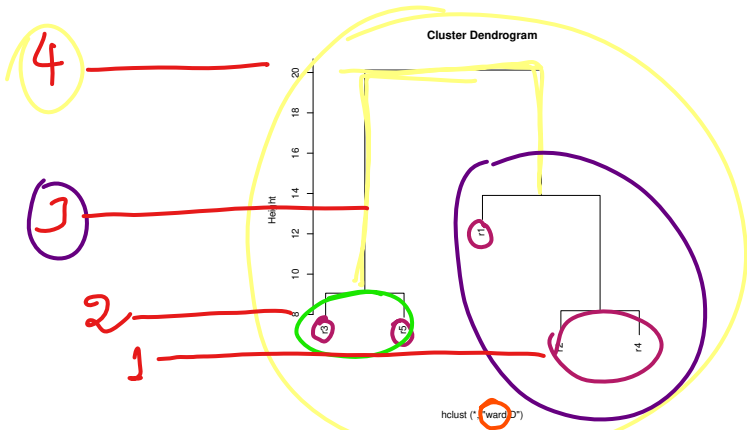
- Ακόμη και οι σύγχρονοι Η/Υ συχνά συναντούν προβλήματα στον προσδιορισμό των δυνατών συστάδων (καθώς το n και το k αυξάνουν).
- Γι'αυτό το λόγο, οι αλγόριθμοί που έχουν αναπτυχθεί συνήθως δίνουν 'καλές' αλλά όχι τις 'βέλτιστες' συστάδες.
- Υπάρχουν 2 βασικές κατηγορίες τέτοιων αλγορίθμων, οι Ιεραρχικοί και οι Μη-Ιεραρχικοί. Στη συνέχεια θα δούμε βασικά χαρακτηριστικά αυτών.

Ιεραρχικοί Αλγόριθμοι Συσταδοποίησης I

- Αποτελούν την κυριότερη κατηγορία και λειτουργούν είτε (α) με διαδοχικές συνενώσεις ομάδων (Agglomerative algorithms) ή (β) με διαδοχικές διαμερίσεις (Divisive algorithms).
- Οι Agglomerative αλγόριθμοι ξεκινούν με τις n διαφορετικές παρατηρήσεις ως μια ομάδα η καθεμία. Οι περισσότερο όμοιες παρατηρήσεις τοποθετούνται σε μια ομάδα και άρα οι ομάδες δημιουργούνται με βάση την ομοιότητα των παρατηρήσεων. Τελικά, όλες οι παρατηρήσεις θα αποτελέσουν (στο τελευταίο βήμα) μια μεγάλη ομάδα.
- Οι Divisive αλγόριθμοι ξεκινούν αντίστροφα, με τις n διαφορετικές παρατηρήσεις να αποτελούν μια μεγάλη ομάδα. Στο επόμενο βήμα οι n παρατηρήσεις χωρίζονται σε δύο ομάδες, και οι παρατηρήσεις που είναι 'πιο μακριά' από τις υπόλοιπες, θα αποτελούν μια ξεχωριστή ομάδα. Ο αλγόριθμος συνεχίζει μέχρι τελικά, όλες οι παρατηρήσεις να αποτελέσουν (στο τελευταίο βήμα) n ξεχωριστές ομάδες.

Ιεραρχικοί Αλγόριθμοι Συσταδοποίησης II

- Οι διαδοχικές διαιρέσεις/συνενώσεις σε κάθε βήμα των παραπάνω αλγορίθμων, αναπαρίστανται με ένα **dendrogram** (δενδρογράμμα).



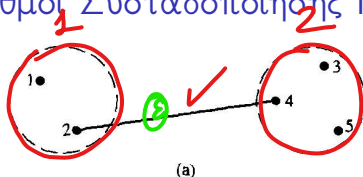
Σχήμα: Παράδειγμα Δενδρογράμματος

Ιεραρχικοί Αλγόριθμοι Συσταδοποίησης III

- Οι μέθοδοι σύνδεσης (Linkage Methods) είναι κατάλληλοι για τη δημιουργία συστάδων, ως μέρος των Ιεραρχικών Αλγορίθμων Συσταδοποίησης.
- Η απλή σύνδεση (single linkage) αναφέρεται στη δημιουργία ομάδων σύμφωνα με μικρότερη απόσταση μεταξύ των μελών.
- Η πλήρης σύνδεση (complete linkage) αναφέρεται στη δημιουργία ομάδων σύμφωνα με τη μεγαλύτερη απόσταση μεταξύ των μελών.
- Τέλος, η μέση σύνδεση (average linkage) αναφέρεται στη δημιουργία ομάδων σύμφωνα με τη μέση απόσταση μεταξύ των μελών. Άλλες μέθοδοι είναι οι Centroid, Ward's method, Median (Παράδειγμα στο εργαστήριο).
- Χαρακτηριστικό είναι το παρακάτω σχήμα.

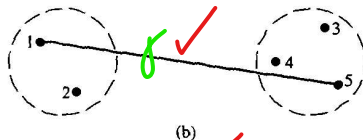
Ιεραρχικοί Αλγόριθμοι Συσταδοποίησης IV

Single
Linkage



min απόσταση

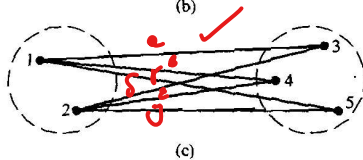
Complete
Linkage



max απόσταση

$D =$ πίνακας απόστασεων

$$\begin{matrix} 1 & \begin{pmatrix} 1 & 3 \\ 0 & 4 \end{pmatrix} \\ 2 & \begin{pmatrix} 6 & 0 \end{pmatrix} \end{matrix}$$



average απόσταση
 μέσος όρος των d_{ij}

Σχήμα: Δημιουργία Συστάδων με Ιεραρχικές Μεθόδους, (a) Single Linkage (b) Complete Linkage, (c) Average Linkage (Πηγή: Johnson and Wichern 2007)

Ιεραρχικοί Αλγόριθμοι Συσταδοποίησης V

- Για την εφαρμογή των Ιεραρχικών Αλγορίθμων Συσταδοποίησης, επιλέγουμε κάποιο μέτρο απόστασης ή ομοιότητας μεταξύ των αντικειμένων (παρατηρήσεων).
- Προφανώς, η ομοιότητα είναι μια φθίνουσα συνάρτηση της απόστασης $d(x, y)$ μεταξύ των x και y ενώ αντίθετα, η ανομοιότητα είναι μια αύξουσα συνάρτηση.
- Γενικά, στην ανάλυση συστάδων, χρησιμοποιείται κάποιο μέτρο (ή δείκτης) ομοιότητας/ανομοιότητας/απόστασης, με την Ευκλείδεια Απόσταση να είναι μια συνήθης επιλογή.
- Θα πρέπει να σημειωθεί ότι η επιλογή της απόστασης $d(x, y)$ μεταξύ των x και y εξαρτάται από τη φύση των διαθέσιμων παρατηρήσεων και από τους σκοπούς για τους οποίους εφαρμόζεται η ανάλυση σε συστάδες, η οποία (κάποιες φορές) μπορεί να είναι και αυθαίρετη.

Ιεραρχικοί Αλγόριθμοι Συσταδοποίησης

Η Περίπτωση των Κατηγορικών Μεταβλητών

- Έστω $\mathbf{X}_i = (X_{i1}, X_{i2}, \dots, X_{ip})$. $i = 1, 2, \dots, n$, όπου τα X_{ij} , $j = 1, 2, \dots, p$ είναι δίτιμες 0-1 τ.μ. Δηλαδή, $X_{ij} = 1$ σημαίνει παρουσία του j -οστού χαρακτηριστικού (συμπτώματος) στο i -οστό άτομο, διαφορετικά έχουμε απουσία του χαρακτηριστικού (συμπτώματος).
- **Παράδειγμα:** $X_{i1} = 1$ το i -οστό άτομο έχει ξανθό μαλλί, $X_{i2} = 1$ το i -οστό άτομο έχει γαλάζια μάτια κ.ο.κ.
- Τότε, μπορούμε να ορίσουμε δείκτες ομοιότητας με βάση τις συμφωνίες/ασυμφωνίες, όπου $a =$ πλήθος των $(X_{i1} = 1, X_{i2} = 1)$, $b =$ πλήθος των $(X_{i1} = 1, X_{i2} = 0)$, $c =$ πλήθος των $(X_{i1} = 0, X_{i2} = 1)$, $d =$ πλήθος των $(X_{i1} = 0, X_{i2} = 0)$.
- Ένας δείκτης σε αυτή την περίπτωση είναι το

$$(a + d) / (a + b + c + d)$$

Αποστάσεις I

Η απόσταση d μεταξύ $\mathbf{x} = (x_1, \dots, x_p)$ και $\mathbf{y} = (y_1, \dots, y_p)$ ικανοποιεί τις παρακάτω ιδιότητες:

- Μη-αρνητική: $d(\mathbf{x}, \mathbf{y}) \geq 0$ με $d(\mathbf{x}, \mathbf{y}) = 0$ αν-ν $\mathbf{x} = \mathbf{y}$.
- Συμμετρική: $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$
- Τριγωνική Ανισότητα: $d(\mathbf{x}, \mathbf{y}) \leq d(\mathbf{y}, \mathbf{z}) + d(\mathbf{x}, \mathbf{z})$

Κάποιες Αποστάσεις:

- Ευκλείδεια: $d_E = |\mathbf{x} - \mathbf{y}| = \sqrt{\sum_{i=1}^p (x_i - y_i)^2} = \sqrt{(\mathbf{x} - \mathbf{y})'(\mathbf{x} - \mathbf{y})}$
- City-Block ή Manhattan: $d_C = \sum_{i=1}^p |x_i - y_i|$
- Minkowski τάξης s : $d_M = (\sum_{i=1}^p |x_i - y_i|^s)^{1/s}$
- Chebyshev: $d_T = \max\{|x_i - y_i|, i = 1, 2, \dots, p\}$
- Quadratic: Αν \mathbf{Q} είναι ένας θετικά ορισμένος πίνακας, ορίζεται η απόσταση $d_Q = (\mathbf{x} - \mathbf{y})'\mathbf{Q}(\mathbf{x} - \mathbf{y})$.

Αποστάσεις II

I

- Mahalanobis: Αν \mathbf{x}, \mathbf{y} είναι τυχαία διανύσματα με τον ίδιο πίνακα διακυμάνσεων-συνδιακυμάνσεων Σ , τότε η απόσταση Mahalanobis ορίζεται ως $\Delta^2 = (\mathbf{x} - \mathbf{y})' \Sigma^{-1} (\mathbf{x} - \mathbf{y})$. Σημειώνεται ότι η απόσταση αυτή είναι γνωστή και ως 'στατιστική απόσταση' μεταξύ των \mathbf{x}, \mathbf{y} .

Ιεραρχικές Μέθοδοι Ομαδοποίησης I

Συμπληρωματικό Υλικό

- Στην ομαδοποίηση αντικειμένων (ατόμων, σημείων), οι λεγόμενες ιεραρχικές μέθοδοι προχωρούν είτε με μια ακολουθία διαδοχικών συγχωνεύσεων είτε με ακολουθία διαδοχικών διαιρέσεων (διαμερίσεων).
- Οι συγχωνευτικές (ή ανιούσες, agglomerative) ιεραρχικές μέθοδοι ξεκινούν με τα n άτομα ως 'ομάδες' (clusters).
- Πρώτα τοποθετούνται σε ομάδες τα πιο όμοια μεταξύ τους άτομα και οι αρχικές αυτές ομάδες συγχνεύονται σύμφωνα με κάποιο βαθμό ομοιότητας (similarity).
- Τελικά με την ομοιότητα να μικραίνει, όλες οι υποομάδες (sub-clusters ή sub-groups) ενώνονται σε μια ομάδα, με στοιχεία όλα τα n άτομα.
- Οι διαιρετικές (ή διαμεριστικές, divisive) ιεραρχικές μέθοδοι, προχωρούν ανάποδα.

Ιεραρχικές Μέθοδοι Ομαδοποίησης II

Συμπληρωματικό Υλικό

- Η αρχική (ενιαία) cluster των n σημείων διαιρείται σε 2 υπο-ομάδες ώστε τα άτομα της μια υποομάδας να 'απέχουν' των ατόμων της άλλης.
- Οι υποομάδες αυτές υποδιαιρούνται περαιτέρω σε ανόμοιες υποομάδες έως ότου φτάσουμε σε n υποομάδες, δηλαδή κάθε άτομο να αποτελεί και ένα cluster.
- Όπως είδαμε και προηγουμένως, η διαδικασία αυτή μπορεί να παρασταθεί γραφικά με ένα δενδρόγραμμα, το οποίο δείχνει τις διαδοχικές συγχωνεύσεις ή διαιρέσεις.

Συγχωνευτικές Ιεραρχικές Μέθοδοι Ομαδοποίησης I

Linkage Methods

- Οι μέθοδοι Linkage χρησιμεύουν για την ομαδοποίηση είτε ατόμων (είτε μεταβλητών, αν και συστήνεται η αποφυγή της), κάτι το οποίο δεν ισχύει για τις συγχωνευτικές μεθόδους ιεραρχικής ταξινόμησης.
- Έχουμε ήδη ορίσει (α) την απλή σύνδεση ή σύνδεση του εγγύτερου γείτονα (Simple Linkage, Nearest Neighbor), (β) την πλήρη σύνδεση ή σύνδεση απώτερου γείτονα (Complete Linkage, Furthest Neighbor), (γ) τη σύνδεση της μέσης απόστασης (Average Linkage).
- Τα βήματα του αλγορίθμου περιγράφονται παρακάτω:
 - ▶ **B1:** Ξεκίνα με n ομάδες, όσες και τα αντικείμενα x_1, x_2, \dots, x_n και έστω $\mathbf{D} = (d_{ij})$ ο πίνακας των αποστάσεων (ή ομοιοτήτων) των n αντικειμένων.
 - ▶ **Σχόλιο:** Το d_{ij} εκφράζει την απόσταση του i από το j , $i \neq j$. Προφανώς, ο \mathbf{D} είναι συμμετρικός με $d_{ij} = d_{ji}$ και $d_{ii} = 0$.
 - ▶ **B2:** Βρες στον πίνακα \mathbf{D} τις όμοιες ομάδες, έστω αυτές U, V (εκεί δηλαδή που είναι η μικρότερη τιμή d_{ij}).

Συγχωνευτικές Ιεραρχικές Μέθοδοι Ομαδοποίησης II

Linkage Methods

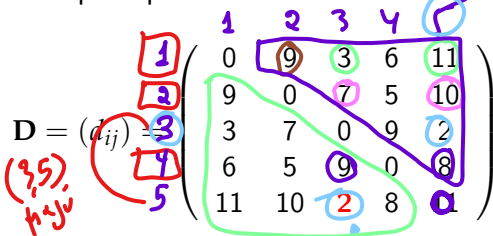
- ▶ **B3:** Συγχώνευσε τις U και V και έστω (UV) η νέα ομάδα. 'Επικαιροποίησε' (update) τον D ως εξής: (1) απαλείφοντας τις γραμμές και στήλες που αντιστοιχούν στις ομάδες U και V και (2) προσθέτοντας μια γραμμή και μια στήλη, στις οποίες θα βρίσκεται πλέον η απόσταση της νέας ομάδας (UV) από τις υπόλοιπες ομάδες.
- ▶ Επανάλαβε τα B2 και B3 για $n - 1$ φορές ώστε τα n άτομα να βρεθούν σε μια ομάδα. Σημειώνουμε ποιές ομάδες συγχωνεύονται και σε ποιες αποστάσεις (ή ομοιότητες) έγιναν οι συγχωνεύσεις.

Παράδειγμα 1

Single Linkage - Nearest Neighbor



- Έστω $n = 5$ αντικείμενα με πίνακα αποστάσεων $D = (d_{ij})$



- Θεωρούμε τα άτομα 1, 2, 3, 4, 5 ως 5 διαφορετικές ομάδες (συστάδες) και η ομαδοποίηση αρχίζει με τη συγχώνευση των 2 πλησιέστερων ατόμων. Προφανώς

$$\min\{d_{ij}, 1 \leq i, j \leq 5\} = d_{53} = 2$$

min αριθμό του D
από εσωτερικό του
3 με το 5

και έτσι οι (5) και (3) ενώνονται και σχηματίζουν τη συστάδα (3,5).

Παράδειγμα II

Single Linkage - Nearest Neighbor

- Στο επόμενο βήμα, ο νέος πίνακας D θα περιέχει τις (εγγύτερες) αποστάσεις της $(3, 5)$ από τα αντικείμενα 1, 2 και 4:

$$d_{(35),1} = \min\{d_{31}, d_{51}\} = \min\{3, 11\} = 3,$$

$$d_{(35),2} = \min\{d_{32}, d_{52}\} = \min\{7, 10\} = 7,$$

$$d_{(35),4} = \min\{d_{34}, d_{54}\} = \min\{9, 8\} = 8.$$

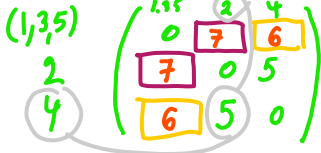
- Διαγράφουμε τις γραμμές και στήλες που αντιστοιχούν στα 3 και 5 και στη συνέχεια προσθέτουμε μια γραμμή και στήλη που αντιστοιχεί στη συστάδα $(3, 5)$, οπότε και προκύπτει ο νέος πίνακας αποστάσεων

$$D = (d_{ij}) = \begin{matrix} & \begin{matrix} (3,5) & 1 & 2 & 4 \end{matrix} \\ \begin{matrix} (3,5) \\ 1 \\ 2 \\ 4 \end{matrix} & \begin{pmatrix} 0 & 3 & 7 & 8 \\ 3 & 0 & 9 & 6 \\ 7 & 9 & 0 & 5 \\ 8 & 6 & 5 & 0 \end{pmatrix} \end{matrix}$$

Συμπερασματικά

Παράδειγμα III

Single Linkage - Nearest Neighbor

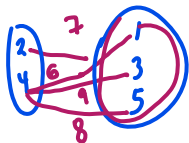


- Πλέον η μικρότερη απόσταση είναι η $d_{(35),1} = 3$ και συνεπώς συγχωνεύουμε τις (3,5) και (1) για το σχηματισμό της (1,3,5).
- Με τρόπο ανάλογο, έχουμε

$$d_{(135),2} = \min\{d_{1,2}, d_{(35),2}\} = \min\{9, 7\} = 7,$$

$$d_{(135),4} = \min\{d_{1,4}, d_{(35),4}\} = \min\{6, 8\} = 6,$$

- Ο νέος πίνακας αποστάσεων είναι



$$D = (d_{ij}) = \begin{matrix} & \begin{matrix} 1,3,5 \\ 2 \\ 4 \end{matrix} \\ \begin{matrix} 1,3,5 \\ 2 \\ 4 \end{matrix} & \begin{pmatrix} 0 & 7 & 6 \\ 7 & 0 & 5 \\ 6 & 5 & 0 \end{pmatrix} \end{matrix}$$

$$\begin{matrix} (1,3,5) & (2,4) \\ (1,3,5) & (2,4) \end{matrix} \begin{pmatrix} 0 & 6 \\ 6 & 0 \end{pmatrix}$$

Παράδειγμα IV

Single Linkage - Nearest Neighbor

(3,5)

- Η νέα απόσταση του εγγύτερου γείτονα είναι η $d_{4,2} = 5$ και άρα τα αντικείμενα 4 και 2 ενώνονται σε μια ομάδα (2,4). Πλέον έχουν δημιουργηθεί 2 ομάδες, η (135) και η (24). Η απόσταση εγγύτερου γείτονα είναι η

$$d_{(135),(24)} = \min\{d_{(135),2}, d_{(135),4}\} = \min\{7, 6\} = 6.$$

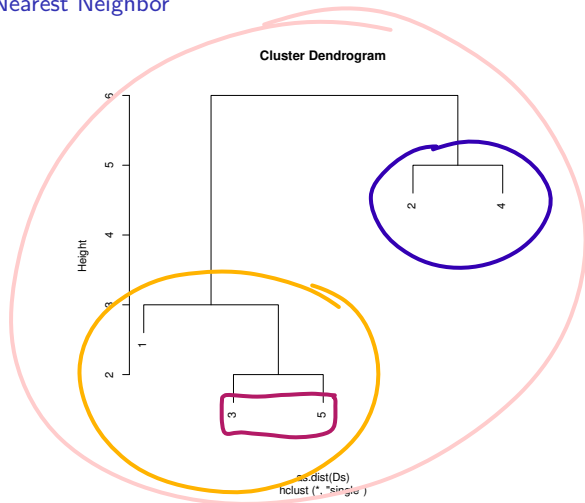
- Ο τελικός 2×2 πίνακας αποστάσεων περιέχει μία μόνο απόσταση, την $d_{(135),(24)} = 6$ και μορφή

$$\mathbf{D} = (d_{ij}) = \begin{pmatrix} 0 & 6 \\ 6 & 0 \end{pmatrix}$$

- Επομένως, οι (135) και (24) συγχωνεύονται σε μια μοναδική ομάδα των 5 ατόμων. Παρακάτω δίνεται και το αντίστοιχο δενδρόγραμμα.

Παράδειγμα V

Single Linkage - Nearest Neighbor



Σχήμα: Δενδρόγραμμα για το παραπάνω παράδειγμα

Παράδειγμα 1

Complete Linkage - Furthest Neighbor

- Ας θεωρήσουμε το ίδιο πλήθος αντικειμένων όπως στο προηγούμενο παράδειγμα και ας εφαρμόσουμε τη μέθοδο της Πλήρους Ένωσης - Απώτατου Γείτονα για τη δημιουργία συστάδων.

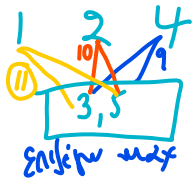
- Όπως και στην Απλή Ένωση, βρίσκουμε ότι

$$\min\{d_{ij}, 1 \leq i, j \leq 5\} = d_{53} = 2$$

| | | | | |
|----|----|---|---|----|
| 0 | 9 | 3 | 5 | 11 |
| 9 | 0 | 7 | 5 | 10 |
| 3 | 7 | 0 | 9 | 2 |
| 6 | 5 | 9 | 0 | 8 |
| 11 | 10 | 2 | 8 | 0 |

και έτσι οι (5) και (3) ενώνονται και σχηματίζουν τη συστάδα (3,5).

- Στο επόμενο βήμα, ο νέος πίνακας **D** θα περιέχει τις (απώτερες) αποστάσεις της (3,5) από τα αντικείμενα 1, 2 και 4:



$$d_{(35),1} = \max\{d_{31}, d_{51}\} = \max\{3, 11\} = 11$$

$$d_{(35),2} = \max\{d_{32}, d_{52}\} = \max\{7, 10\} = 10,$$

$$d_{(35),4} = \max\{d_{34}, d_{54}\} = \max\{9, 8\} = 9.$$

| | | | |
|----|----|----|---|
| 35 | 1 | 2 | 4 |
| 0 | 11 | 10 | 9 |
| 11 | 0 | 9 | 6 |
| 10 | 9 | 0 | 5 |
| 9 | 6 | 5 | 0 |

Παράδειγμα II

Complete Linkage - Furthest Neighbor

- Διαγράφουμε τις γραμμές και στήλες που αντιστοιχούν στα 3 και 5 και στη συνέχεια προσθέτουμε μια γραμμή και στήλη που αντιστοιχεί στη συστάδα (3,5), οπότε και προκύπτει ο νέος πίνακας αποστάσεων

$$D = (d_{ij}) = \begin{matrix} & \begin{matrix} 3,5 & 1 & 2 & 4 \end{matrix} \\ \begin{matrix} 3,5 \\ 1 \\ 2 \\ 4 \end{matrix} & \begin{pmatrix} 0 & 11 & 10 & 9 \\ 11 & 0 & 9 & 6 \\ 10 & 9 & 0 & 5 \\ 9 & 6 & 5 & 0 \end{pmatrix} \end{matrix}$$

min distance
Θα ευθεία το

- Πλέον η μικρότερη απόσταση είναι η $d_{(3,5),1} = 5$ και συνεπώς συγχωνεύουμε τις (2) και (4) στην ομάδα (2,4).

- Με τρόπο ανάλογο, έχουμε

$$\begin{matrix} & \begin{matrix} 3,5 & 2,4 & 1 \end{matrix} \\ \begin{matrix} 3,5 \\ 2,4 \\ 1 \end{matrix} & \begin{pmatrix} 0 & 10 & 11 \\ 10 & 0 & 9 \\ 11 & 9 & 0 \end{pmatrix} \end{matrix}$$

min distance
= 9

$$d_{(2,4),(3,5)} = \max\{d_{2,(3,5)}, d_{4,(3,5)}\} = \max\{10, 9\} = 10,$$

$$d_{(2,4),1} = \max\{d_{2,1}, d_{4,1}\} = \max\{9, 6\} = 9.$$



Παράδειγμα III

Complete Linkage - Furthest Neighbor

- Ο νέος πίνακας αποστάσεων είναι

$$D = (d_{ij}) = \begin{pmatrix} 0 & 10 & 11 \\ 10 & 0 & 9 \\ 11 & 9 & 0 \end{pmatrix}$$

- Η νέα μικρότερη απόσταση είναι η $d_{(2,4),1} = 9$ και άρα τα αντικείμενα 2, 4 και 1 ενώνονται σε μια ομάδα (1, 2, 4). Πλέον έχουν δημιουργηθεί 2 ομάδες, η (35) και η (124). Η απόσταση απώτερου γείτονα είναι η

$$d_{(35),(124)} = \max\{d_{1,(35)}, d_{(24),(35)}\} = \max\{11, 10\} = 11.$$

- Ο τελικός 2×2 πίνακας αποστάσεων περιέχει μία μόνο απόσταση, την $d_{(35),(124)} = 11$ και μορφή

$$D = (d_{ij}) = \begin{pmatrix} 0 & 11 \\ 11 & 0 \end{pmatrix}$$

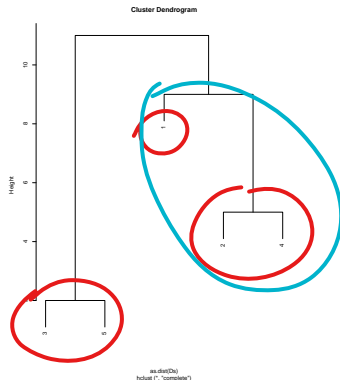
1 ενώνεται με {34}



Παράδειγμα IV

Complete Linkage - Furthest Neighbor

- Επομένως, οι (35) και (124) συγχωνεύονται σε μια μοναδική ομάδα των 5 ατόμων. Παρακάτω δίνεται και το αντίστοιχο δενδρόγραμμα.

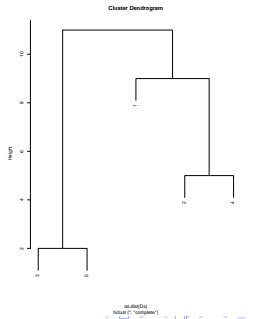
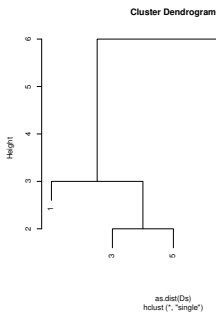


Σχήμα: Δενδρόγραμμα για το παραπάνω παράδειγμα, Furthest Neighbor

Παράδειγμα V

Complete Linkage - Furthest Neighbor

- Συγκρίνοντας τα δύο δενδρογράμματα της απλής και της πλήρους ένωσης, βλέπουμε ότι οι σχηματισμοί των ομάδων (35) και (24) γίνονται σε διαφορετικά στάδια της διαδικασίας συσταδοποίησης.
- Και οι 2 τεχνικές οδηγούν στις ομάδες (24) και (35) με κάποια διαφοροποίηση ως προς την τοποθέτηση του 1, είτε στην (124) (simple linkage) είτε στην (135) (complete linkage).



Μη-Ιεραρχικοί Αλγόριθμοι Συσταδοποίησης I

k-Means Method

ΠΑΛΙ ΚΟΙΤΗΣΤΕ

SSE
SSE
SSE ⊕
SSTO =

SSE =
sum of squares
error = $\sum (y - \hat{y})^2$

- Οι πιο συχνά χρησιμοποιούμενες μέθοδοι δημιουργίας συστάδων βασίζονται στο κριτήριο του τετραγωνικού σφάλματος (square-error criterion) **LSE**
- Έστω X το σύνολο $n \times p$ των δεδομένων, το οποίο έχει διαμεριστεί (με κάποιον τρόπο) στις συστάδες C_1, C_2, \dots, C_k , καθεμία από τις οποίες περιέχει n_i παρατηρήσεις εκ των x_1, x_2, \dots, x_n .
- Έστω

$$\bar{x}_i = \frac{1}{n_i} \sum_{s=1}^{n_i} x_s^{(i)}, \quad i = 1, 2, \dots, k$$

το οποίο αποτελεί το κέντρο της i -οστής συστάδας.

$\bar{x}_1, \bar{x}_2, \dots, \bar{x}_k = k$ κέντρα των k συστάδων

$C_i = i$ -th cluster
(με n_i observations)

$i = 1, 2, 3, \dots, k$

που ανήκουν ως
 $x_1^{(i)}, x_2^{(i)}, \dots, x_{n_i}^{(i)}$

Μη-Ιεραρχικοί Αλγόριθμοι Συσταδοποίησης II

$$n_1 + n_2 + \dots + n_k = n$$

k-Means Method

- Η μεταβλητότητα εντός της C_i (within-cluster variation) ισούται με το άθροισμα των τετραγώνων της Ευκλείδειας απόστασης μεταξύ $x_s^{(i)}$ και \bar{x}_i , δηλ.



$$e_i^2 = \sum_{s=1}^{n_i} (x_s^{(i)} - \bar{x}_i)' (x_s^{(i)} - \bar{x}_i)$$

απόσταση κάθε
παρατήρησης από
το κέντρο της

$i=1,2,3, \dots, k = \# \text{ clusters}$

- Το ολικό σφάλμα συσταδοποίησης (overall clustering error) είναι το

$$\text{Sum of Squared errors} = E_k^2 = \sum_{i=1}^k e_i^2$$

$e_1^2 =$ σφάλμα του σφάλματος
N το cluster 1
 $e_k^2 =$ τέτοιο σφάλμα cluster k

και ο στόχος είναι να βρεθεί μια διαμέριση του συνόλου δεδομένων X η οποία να ελαχιστοποιεί το E_k^2 για δεδομένο k .

- Ο αλγόριθμος για την εύρεση αυτής της διαμέρισης είναι γνωστός ως k-means Algorithm και τα βήματά του είναι τα εξής:

Μη-Ιεραρχικοί Αλγόριθμοι Συσταδοποίησης III

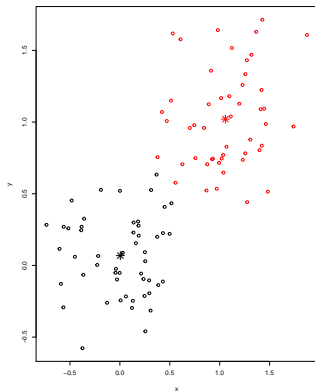
k-Means Method

- ▶ **B1:** Επιλέγω τα k αρχικά κέντρα (seed points) για τις k (μη προσδιορισμένες) ομάδες.
 - ▶ **B2:** Τοποθετούμε καθένα από τα $x_i, i = 1, \dots, n$, αντικείμενα στην πιο κοντινή ομάδα (μικρότερη απόσταση με το κέντρο της ομάδας).
 - ▶ **B3:** Υπολογίζουμε εκ νέου τα κέντρα των νέων ομάδων (μετά την τοποθέτηση των αντικειμένων στο B2)
 - ▶ **B4:** Επαναλαμβάνουμε τα B2 και B3 μέχρι να προκύψει μια βέλτιστη τιμή για το E_k^2 .
 - ▶ **Σχόλιο:** Αφού σχηματιστούν οι ομάδες, μπορούμε να τροποποιήσουμε τον τελικό αριθμό τους είτε συγχωνεύοντας είτε διαιρώντας κάποιες από τις υπάρχουσες, είτε ακόμη και απομακρύνοντας μικρές ή εξωκείμενες ομάδες.
- Ο αλγόριθμος *k*-means είναι υπολογιστικά αποδοτικός και δίνει πολύ καλά αποτελέσματα (συμπαγείς ομάδες, σφαιρικές και καλώς διαχωρισμένες).

Μη-Ιεραρχικοί Αλγόριθμοι Συσταδοποίησης IV

k-Means Method

- Επίσης, έχει καλή απόδοση στην ανίχνευση/δημιουργία ομάδων ελλειψοειδούς σχήματος, αν χρησιμοποιηθεί η απόσταση Mahalanobis αντί της Ευκλείδειας.



Μη Ιεραρχικές Μέθοδοι Ταξινόμησης I

Συμπληρωματικό Υλικό

- Οι τεχνικές αυτές δεν απαιτούν τον προσδιορισμό του πίνακα αποστάσεων D ενώ τα δεδομένα x_i , $i = 1, \dots, n$ δεν είναι ανάγκη να αποθηκεύονται κατά την εκτέλεση του αλγορίθμου.
- Κατά συνέπεια, οι μη-ιεραρχικές μέθοδοι ταξινόμησης είναι προτιμότερες για μεγάλα σύνολα δεδομένων.
- **Παράδειγμα:** Έστω $k = 2$ ο προκαθορισμένος αριθμός ομάδων και $n = 4$ άτομα (έστω A, B, C, D), για τα οποία έχουν σημειωθεί οι τιμές των (X_1, X_2)

| Άτομο | x_1 | x_2 |
|-------|-------|-------|
| A | 5 | 3 |
| B | -1 | 1 |
| C | 1 | -2 |
| D | -3 | -2 |

Μη Ιεραρχικές Μέθοδοι Ταξινόμησης II

Συμπληρωματικό Υλικό

- Παίρνουμε αυθαίρετα τη διαμέριση (AB) και (CD) και υπολογίζουμε τα κέντρα (\bar{x}_1, \bar{x}_2) για την κάθε ομάδα.

| Ομάδα | \bar{x}_1 | \bar{x}_2 |
|-------|-------------------------|--------------------------|
| (AB) | $\frac{5+(-1)}{2} = 2$ | $\frac{3+1}{2} = 2$ |
| (CD) | $\frac{1+(-3)}{2} = -1$ | $\frac{-2+(-2)}{2} = -2$ |

- Στο επόμενο βήμα, υπολογίζουμε την Ευκλείδεια απόσταση κάθε ατόμου από τα κέντρα των 2 ομάδων. Τοποθετούμε το κάθε άτομο στην πλησιέστερη ομάδα.
- Όταν ένα άτομο φύγει από την αρχική του ομάδα, τα κέντρα των (νέων) ομάδων 'διορθώνονται' προτού προχωρήσουμε στο επόμενο βήμα.

Μη Ιεραρχικές Μέθοδοι Ταξινόμησης III

Συμπληρωματικό Υλικό

- Υπολογίζουμε τις Ευκλείδειες αποστάσεις (τα τετράγωνα τους), αρχικά για το A:

$$d^2(A, (AB)) = (5 - 2)^2 + (3 - 2)^2 = 10,$$

$$d^2(A, (CD)) = (5 + 1)^2 + (3 + 2)^2 = 61.$$

- Άρα το A είναι πιο κοντά στην (AB) απ'ότι στη (CD), οπότε και παραμένει στην (AB).
- Ομοίως, για το B:

$$d^2(B, (AB)) = (-1 - 2)^2 + (-1 - 2)^2 = 10,$$

$$d^2(B, (CD)) = (-1 + 1)^2 + (1 + 2)^2 = 9.$$

- Επομένως, το B τοποθετείται στη (CD), δίνοντας την ομάδα (BCD). Τα νέα κέντρα είναι:

Μη Ιεραρχικές Μέθοδοι Ταξινόμησης IV

Συμπληρωματικό Τλικό

| Ομάδα | \bar{x}_1 | \bar{x}_2 |
|-------|----------------------------|---------------------------|
| A | 5 | 3 |
| (BCD) | $\frac{-1+1+(-3)}{3} = -1$ | $\frac{1-2+(-2)}{3} = -1$ |

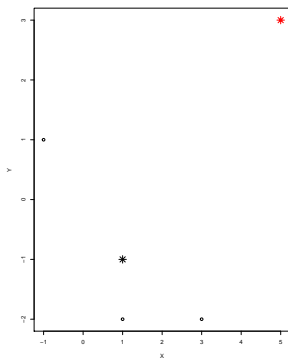
- Πάλι ελέγχουμε κάθε άτομο για ενδεχόμενο μετακίνηση σε άλλη ομάδα. Βρίσκουμε τον πίνακα με τις Ευκλείδειες αποστάσεις από τα κέντρα των (A) και (BCD).

| Ομάδα | A | B | C | D |
|-------|----|----|----|----|
| A | 0 | 40 | 41 | 89 |
| (BCD) | 52 | 4 | 5 | 5 |

- Βλέπουμε ότι κάθε άτομο παραμένει στην ομάδα του και η διαδικασία σταματά με $k = 2$ ομάδες, την (A) και την (BCD).

Μη Ιεραρχικές Μέθοδοι Ταξινόμησης V

Συμπληρωματικό Υλικό



Σχήμα: Διαγραμματική Απεικόνιση των 2 Ομάδων, K-Means example