

Στο έγγραφο αυτό συγκεντρώνουμε τις ποιό συχνές παρατηρήσεις που κάναμε στην πρόοδο της εξόρυξης γνώσης από δεδομένα.

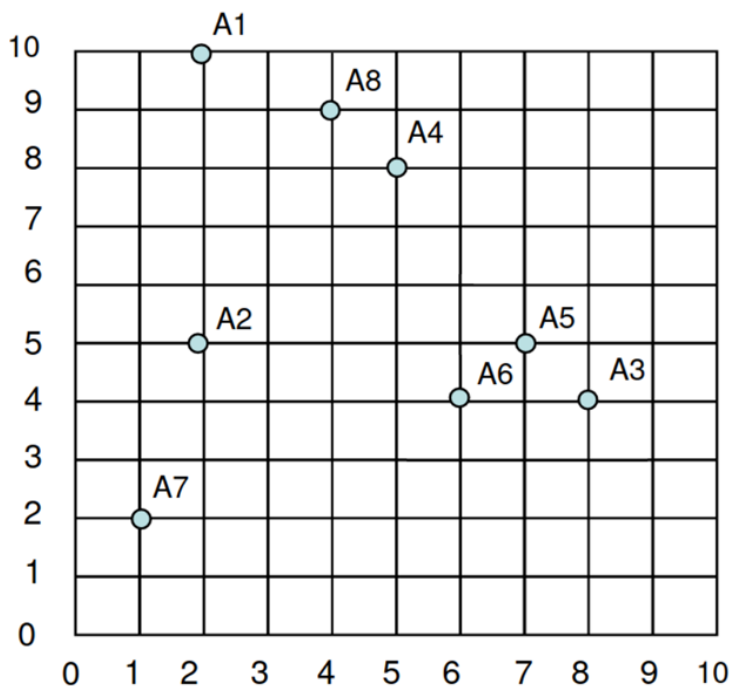
Σημειώστε βέβαια ότι σε όλες τις ερωτήσεις οι μονάδες δεν ήταν “όλες ή μηδέν”, αλλά προσμετρήθηκαν κλιμακωτά, ανάλογα την απάντηση.

Ποιο ή ποιά από τα ακόλουθα είναι παραδείγματα εξαγωγής χαρακτηριστικών;

- A. Κατασκευή διανύσματος “bag of words” από ένα email
- B. Εφαρμογή PCA σε μεγάλα δεδομένα υψηλής διαστάσεων
- Γ. Αφαίρεση stopwords από μια πρόταση
- Δ. Καμία από τις υπόλοιπες απαντήσεις

A, B και Γ.

Δίνονται τα 8 σημεία (A1 ως A8) στον δισδιάστατο χώρο, τα οποία τα συσταδοποιούμε με τη χρήση του αλγορίθμου DBSCAN με $\text{MinPts}=2$ και $\text{Eps}=\sqrt{10}$. Να γράψετε ποιες συστάδες σχηματίζονται, αν υπάρχουν οριακά σημεία (border points), ποια είναι και σε ποια συστάδα είναι οριακά καθώς και αν υπάρχουν σημεία θορύβου (noise points) και ποια είναι αυτά.



Εφόσον $\text{MinPts} = 2$ και $\text{Eps} = 3,16$, αρχικά χρειάζεται να χαρακτηρίσουμε τα σημεία μας σε σχέση με τον 2ο πλησιέστερο γείτονα τους και ακτίνας $= 3,16$. Τα A5,A6,A3 είναι εντός μιας συστάδας γιατί η απόσταση πχ του A5 από το A6 είναι $\sqrt{(7-6)^2+(5-4)^2} = \sqrt{2} < \text{eps}$. Αντίστοιχα προκύπτει και για το A5-A3 και εφόσον αυτά τα τρία είναι δίπλα το ένα στο άλλο για καθένα από αυτά αποτελούν 2ο πλησιέστερο γείτονα τα άλλα δυο. Η δεύτερη συστάδα αποτελείται από τα A8,A4 και A1. Η τελευταία συστάδα είναι η A2,A7.

Με βάση τα παραπάνω, τα σημεία A5, A2 και A8 αποτελούν σημεία πυρήνα και το A7 αποτελεί οριακό σημείο γιατί απέχει απόσταση από το A2 ακριβώς απόσταση ρίζας(10).

Έχουμε έξι δείγματα δύο χαρακτηριστικών το καθένα: [4 1], [6 6], [9 5], [1 2], [7 3], [5 4]. Οι ετικέτες τους είναι [1 0 1 0 1 0]. Για την ταξινόμηση φτιάχνουμε ένα δέντρο αποφάσεων βάθους δύο. Οι διχοτομήσεις προκύπτουν με κανόνες που θέτουν ένα χαρακτηριστικό μεγαλύτερο ίσο από ένα κατώφλι. α) ποια είναι η εντροπία της ρίζας του δέντρου; β) ποιος είναι ο κανόνας της δεύτερης διχοτόμησης;

α) $-0.5 \log_2 0.5 - 0.5 \log_2 0.5 = 1$ ή απλά αν οι ετικέτες είναι σε δυαδικό πρόβλημα οι ετικέτες είναι μισές-μισές η εντροπία είναι μέγιστη δηλαδή 1.

β) Κατατάσσουμε τα δείγματα και τις ετικέτες με βάση το πρώτο χαρακτηριστικό x_1

1 2 0

4 1 1

5 4 0

6 6 0

7 3 1

9 5 1

Κάνουμε το ίδιο για το δεύτερο χαρακτηριστικό x_2

4 1 1

1 2 0

7 3 1

5 4 0

9 5 1

6 6 0

Ο καλύτερος πρώτος κανόνας είναι $x_1 \geq 7$.

Για τη δεύτερη διχοτόμηση: το υποδέντρο με ετικέτες [1, 1] δεν χρειάζεται να χωριστεί.

Κάνουμε τις δύο κατατάξεις για το υποδέντρο με ετικέτες [0, 1, 0, 0]

1 2 0

4 1 1

5 5 0

6 5 0

και

4 1 1

1 2 0

5 4 0

6 6 0

Ο καλύτερος δεύτερος κανόνας είναι ο $x_2 \geq 2$

Έχουμε έξι δείγματα δύο χαρακτηριστικών: [4 1], [6 6], [9 5], [1 2], [7 3], [5 4]. Οι ετικέτες τους είναι [1 0 1 0 1 0]. Για την ταξινόμηση φτιάχνουμε ένα δέντρο βάθους δύο. Οι διχοτομήσεις προκύπτουν με κανόνες που θέτουν ένα χαρακτηριστικό μεγαλύτερο ίσο από ένα κατώφλι. α) στη ρίζα του δέντρου, ποιοι κανόνες δίνουν μηδενικό κέρδος πληροφορίας; β) ποιος είναι ο κανόνας της πρώτης διχοτόμησης;

α) Κατατάσσουμε τα δείγματα και τις ετικέτες με βάση το πρώτο χαρακτηριστικό x_1

1 2 0
4 1 1
5 4 0
6 6 0
7 3 1
9 5 1

Κάνουμε το ίδιο για το δεύτερο χαρακτηριστικό x_2

4 1 1
1 2 0
7 3 1
5 4 0
9 5 1
6 6 0

Οι κανόνες $x_1 \geq 5$, $x_2 \geq 3$ ή $x_2 \geq 5$ δεν πετυχαίνουν να μειώσουν την εντροπία.

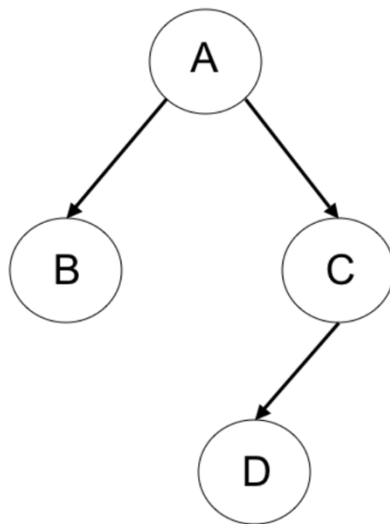
β) Ο καλύτερος πρώτος κανόνας είναι ο $x_1 \geq 7$.

Εκτελούμε τον αλγόριθμο Expectation-Maximization για ένα μοντέλο μίξης δύο γκαουσιανών κατανομών με ίδια βάρη για ένα σύνολο 3 δειγμάτων δεδομένων και σε κάποιο βήμα της επανάληψης έχουμε $p(x_1=-2|\theta_1)=A$, $p(x_1=-2|\theta_2)=B$, $p(x_2=1|\theta_1)=C$, $p(x_2=1|\theta_2)=D$, $p(x_3=3|\theta_1)=E$, $p(x_3=3|\theta_2)=F$. Έστω $A=0.3$, $B=0.7$, $C=0.45$, $D=0.55$, $E=0.2$, και $F=0.8$. Να υπολογίσετε την πιθανότητα το κάθε δείγμα δεδομένων να προέρχεται από κάθε μια από τις δύο κατανομές καθώς και τις μέσες τιμές των δύο κατανομών μετά την ολοκλήρωση αυτού του βήματος (στο διάστημα $[0,1]$, με ακρίβεια 2 δεκαδικών ψηφίων)

$p(j=1|x_1=-2)=0.3$, $p(j=2|x_1=-2)=0.7$, $p(j=1|x_2=1)=0.45$, $p(j=2|x_2=1)=0.55$, $p(j=1|x_3=3)=0.2$,
 $p(j=2|x_3=3)=0.8$
 $\mu_1=0.47$, $\mu_2=0.76$

Ενδεικτικά. Σε όλα τα versions του, το ερώτημα για τον EM μοιάζει να δυσκόλεψε.

**Δίνεται το Μπεϋζιανό δίκτυο πεποίθησης του σχήματος και οι εξείς πιθανότητες:
 $P(A)=0.3$, $P(B|A)=0.7$, $P(B|\sim A)=0.5$, $P(C|A)=0.2$, $P(C|\sim A)=0.6$, $P(D|C)=0.7$, $P(D|\sim C)=0.7$.
Υπολογίστε την πιθανότητα $P(C|B)$**



$P(C|B) = P(C, B) / P(B)$ Ομως $p(B, C) = P(B, C|A) + P(B, C|\sim A) = P(B|A)P(C|A)P(A) + P(B|\sim A)P(C|\sim A)P(\sim A)$
 $P(\sim A) = 0.7 \times 0.2 \times 0.3 + 0.6 \times 0.5 \times 0.7 = 0.252$ $P(B) = P(B|A)P(A) + P(B|\sim A)P(\sim A) = 0.7 \times 0.3 + 0.5 \times 0.7 = 0.56$ Άρα $P(C|B) = 0.45$

Ενδεικτικά. Σε όλα τα versions του, το ερώτημα για τα δίκτυα πεποίθησης μοιάζει να δυσκόλεψε.