

ΑΣΚΗΣΗ

Ημερομηνία Παράδοσης: τέλος εξεταστικής περιόδου

Η άσκηση είναι **ατομική, υποχρεωτική και απαλλακτική** για την εξέταση του μαθήματος (αφορά και φοιτητές παλαιότερων ετών). Από την εκπόνηση της άσκησης εξαιρούνται οι φοιτητές που θα εξεταστούν προφορικά.

Η αναφορά που θα συντάξετε θα παραδοθεί σε ηλεκτρονική μορφή. **Προσοχή**, στην αναφορά που θα παραδώσετε να αναγράφεται το όνομά σας, το επίθετό σας, το ΑΜ σας και το έτος σπουδών σας.

Κάθε απάντηση θα είναι **αιτιολογημένη**. Η αιτιολόγηση μπορεί να γίνει είτε περιγράφοντας τη μεθοδολογία που ακολουθήσατε για την ολοκλήρωση κάθε βήματος είτε υπό τη μορφή παραδειγμάτων. Σε περίπτωση αντιγραφής μηδενίζονται και οι δύο ασκήσεις.

Ασκήσεις που θα παραδοθούν εκπρόθεσμα δεν γίνονται δεκτές!

ΕΚΦΩΝΗΣΗ

Έστω μια συλλογή κειμένων C που περιέχει 5 κείμενα:

$$C \{d_1, d_2, d_3, d_4, d_5\}$$

Όπου για κάθε κείμενο d ισχύουν τα ακόλουθα:

$$d_1 \{w_a, w_b, w_c, w_d, w_e\}$$

$$d_2 \{w_a, w_f, w_c, w_j, w_k, w_m\}$$

$$d_3 \{w_e, w_k, w_n, w_c, w_d, w_p, w_t\}$$

$$d_4 \{w_b, w_f, w_k, w_u, w_c, w_m, w_d\}$$

$$d_5 \{w_a, w_c, w_f, w_k, w_y, w_z\}$$

όπου όλοι οι όροι εκτός από τους w_b, w_c, w_d είναι ΟΥΣΙΑΣΤΙΚΑ.

Επιπλέον, έστω ότι:

$$D_1 \{fw_a=5, fw_b=8, fw_c=7, fw_d=10, fw_e=5\}, \text{ όπου } f = \text{frequency}$$

$$D_2 \{fw_a=3, fw_f=9, fw_c=10, fw_j=8, fw_k=6, fw_m=4\}, \text{ όπου } f = \text{frequency}$$

$D_3 \{fw_e=10, fw_k=2, fw_n=8, fw_c=3, fw_d=7, fw_p=10, fw_t=10\}$, όπου $f = frequency$

$D_4 \{fw_b=5, fw_f=8, fw_k=17, fw_u=3, fw_c=6, fw_m=1, fw_d=8\}$, όπου $f = frequency$

$D_5 \{fw_a=5, fw_c=8, fw_f=10, fw_k=3, fw_y=6, fw_z=8\}$, όπου $f = frequency$

Για τα παραπάνω δεδομένα θέλουμε να δημιουργήσουμε ένα ευρετήριο όρων προκειμένου να μπορούμε να ανακτήσουμε πληροφορία από τα δεδομένα της συλλογής.

(i) Παρουσιάστε το ευρετήριο όπου θα αναγράφονται οι όροι που θα περιέχονται σε αυτό, οι δείκτες προς τα κείμενα που τους περιέχουν και για κάθε κείμενο ο βαθμός ευρετηρίασής τους (w score).

(ii) Έστω ότι ισχύουν τα ακόλουθα για το d_5

$d_5 \{fw_a=5, fw_c=8, fw_f=10, fw_k=3, fw_y=6, fw_z=8\}$, όπου $f = frequency$

Ποιος είναι ο πιο σημαντικός όρος ευρετηρίασης για το d_5 ; Να αιτιολογήσετε την απάντησή σας και να δώσετε την αριθμητική τιμή που δείχνει το βαθμό σπουδαιότητας του όρου που επιλέξατε για το d_5

(iii) Ποιος ο συνολικός βαθμός σπουδαιότητας για το ευρετήριο του όρου που επιλέξατε στον προηγούμενο ερώτημα; (υπολογίστε τη σπουδαιότητα του όρου με τη μετρική $TF*IDF$ για το d_5).

(iv) Έστω το ερώτημα q που περιέχει τους όρους $q = w_a, w_f, w_z$. Ποια κείμενα θα επιστρέφονταν ως απάντηση για το ερώτημα και με ποια σειρά; Αιτιολογήστε την απάντησή σας.

(v) Αν για το ερώτημα $q = w_a$ επιστραφούν τα κείμενα D_5 και D_1 με αυτή τη σειρά, ποια είναι η απόδοση της ανάκτησης (Ακρίβεια / Ανάκληση).

Καλή επιτυχία!