

3η Σειρά Ασκήσεων  
Ανάθεση: 12 Μαΐου  
Παράδοση: 30 Μαΐου

**Άσκηση 1** (2 βαθμού)

Θεωρείστε ένα έγγραφο με περιεχόμενο:

«**στο σημερινό μάθημα μάθαμε περισσότερα σε σχέση με το προηγούμενο μάθημα**»

Αγνοώντας τους τόνους, σχεδιάστε

- (α) το trie του λεξιλογίου του παραπάνω εγγράφου,
- (β) το δένδρο καταλήξεων του παραπάνω εγγράφου όπου υπερέχουν (index points) τις αρχές των λέξεων, και
- (γ) συμπτύξτε το παραπάνω δένδρο καταλήξεων στη μορφή ενός Patricia tree.

**Άσκηση 2** (2 βαθμ.)

Θεωρείστε το αλφάριθμο  $\{\alpha, \beta, \gamma, \delta, \epsilon, \zeta\}$  και την εξής πρόταση:

“ $\alpha \alpha \beta \beta \gamma \alpha \gamma \alpha \delta \alpha \alpha \alpha \delta \beta \epsilon \beta \epsilon \zeta$ ”.

α) Βάσει αυτής της φράσης ποια είναι η εντροπία του αλφαριθμού;

β) Δώστε τη συμπλεσμένη μορφή της φράσης χρησιμοποιώντας κανονικοποιημένους κώδικες Huffman.

**Άσκηση 3** (1 βαθμός)

Τιολογίστε την Edit Distance μεταξύ των λέξεων paris και alice. Δώστε τον 5x5 πίνακα που περιγράφει τον τρόπο λειτουργίας του σχετικού αλγορίθμου δυναμικού προγραμματισμού.

**Άσκηση 4** (2 βαθμ.)

Θεωρείστε μια συλλογή εγγράφων στην οποία εμφανίζονται 400 διαφορετικές λέξεις και η συχνότητα της πιο συχνά εμφανιζόμενης ισούται με 900.

(α) Εκτιμήστε τον αριθμό εμφανίσεων της 10ης πιο συχνά εμφανιζόμενης λέξης, και της 20ης.

(β) Αποφασίζετε να φτιάξετε ένα ανεστραμμένο ευρετήριο μόνο για τις 200 πιο συχνά εμφανιζόμενες λέξεις. Πόσο μικρότερο να είναι το μέγεθος του σε σχέση με εκείνο για όλες τις λέξεις; Θεωρείστε οτι δεν κάνουμε απαλοιφή λέξεων αποκλεισμού, ούτε στελέχωση.

**Άσκηση 5** (3 βαθμ.)

Θεωρείστε τα ακόλουθα έγγραφα όπου τα γράμματα A-E συμβολίζουν λέξεις.

$$\begin{aligned} d1 &= "A A B" \\ d2 &= "B A A" \\ d3 &= "A B C" \\ d4 &= "C C B" \\ d5 &= "C D A" \\ d6 &= "D E" \\ d7 &= "A A B" \\ d8 &= "E B" \end{aligned}$$

Έστω ότι τα  $d_1, d_2, d_3$  ανήκουν σε ένα σύστημα  $S_1$ , τα  $d_4, d_5$  σε ένα σύστημα  $S_2$ , και τα υπόλοιπα ( $d_6, d_7, d_8$ ) σε ένα σύστημα  $S_3$ . Θέλουμε να φτιάξουμε έναν μεσίτη  $M$  πάνω από αυτά τα συστήματα.

(α) Για την επιλογή πηγής ο  $M$  θέλει να περιγράψει τα περιεχόμενα της κάθε πηγής με ένα διάνυσμα. Δώστε τα διανύσματα πηγών των  $S_1, S_2$  και  $S_3$ .

(β) Έστω ότι ο  $M$  έχει ήδη τα διανύσματα πηγών των  $S_1, S_2, S_3$  και λαμβάνει την επερώτηση  $q = "B\ C"$ . Αν θέλει να προωθήσει την επερώτηση  $q$  σε μία μόνο πηγή, ποια θα επιλέξει;

(γ) Ο  $M$  λαμβάνει μια επερώτηση, την προωθεί σε όλες τις πηγές, και λαμβάνει τα εξής αποτελέσματα από την κάθε μια:

$S_1: \langle d_1, d_2, d_3 \rangle$

$S_2: \langle d_5, d_4 \rangle$

$S_3: \langle d_8, d_7, d_6 \rangle$

Δώστε την ενοποιημένη διάταξη κατά round robin interleaving

(δ) Προκειμένου ο μεσίτης να λαμβάνει από τις πηγές απαντήσεις με συγκρίσιμα σκορ, αποφασίζει να κάνει αποτίμηση επερωτήσεων σε δυο φάσεις ώστε οι πηγές να λαμβάνουν τα καθολικά στατιστικά που χρειάζονται για το σωστό υπολογισμό των σκορ. Δώστε το idf του κάθε όρου στην καθολική συλλογή εγγράφων.

(ε) Ο μεσίτης βρίσκει άλλο ένα σύστημα  $S_4$  το οποίο έχει την ίδια συλλογή με αυτήν του  $S_1$ , δηλαδή και αυτό παρέχει πρόσβαση στα έγγραφα  $d_1, d_2, d_3$ . Έστω ότι ο  $M$  προωθεί μια επερώτηση  $q$  στα  $S_1$  και  $S_4$  και λαμβάνει τις εξής απαντήσεις:

$S_1: \langle d_2, d_1, d_3 \rangle$

$S_4: \langle d_2, d_3, d_1 \rangle$

Ποιο είναι το κορυφαίο έγγραφο αν ενοποιήσουμε τις διατάξεις: (i) κατά Borda, (ii) κατά Condorcet;

Ο  $M$  αποφασίζει να δίνει στο χρήστη όχι μόνο την ενοποιημένη διάταξη, αλλά και την Kemeny distance μεταξύ των διατάξεων που έλαβε από τα υποσυστήματα (προκειμένου ο χρήστης να παίρνει μια γεύση για το βαθμό συμφωνίας των πηγών). Ποια είναι αυτή η απόσταση στην προκειμένη;

(στ) Τα συστήματα  $S_1, S_2, S_3$  δεν θέλουν πλέον να έχουν ανάγκη τον  $M$  και αποφασίζουν να "ανεξαρτητοποιηθούν" φτιάχνοντας ένα σύστημα ομοτίμων (P2P), συγκεκριμένα ένα δομημένο σύστημα τύπου Chord. Προσελκύουν μάλιστα άλλα δυο συστήματα  $S_5$  και  $S_6$  (τα οποία δεν έχουν καμία συλλογή εγγράφων). Αποφασίζουν να χρησιμοποιήσουν μια συνάρτηση κατακερματισμού  $h$  των 3 bits, και έστω ότι:  $h(IPaddress(S_1)) = 1, h(IPaddress(S_2)) = 2, h(IPaddress(S_3)) = 4, h(IPaddress(S_5)) = 5, h(IPaddress(S_6)) = 6$ . Αποφασίζουν να διανείμουν το ανεστραμμένο ευρετήριο ψεωρώντας κάθε όρο σαν κλειδί και έστω ότι  $h(A) = 2, h(B) = 3, h(C) = 4, h(D) = 4, h(E) = 4$ .

Δώστε (i) τους πίνακες δρομολόγησης των κόμβων  $S_1$  και  $S_3$  και (ii) πως θα κατανεμηθεί το ανεστραμμένο ευρετήριο στους κόμβους του δικτύου (δείξτε τι ακριβώς θα έχει κάθε κόμβος).-