

4^η Σειρά ασκήσεων
(Συμπύεση, Ομαδοποίηση, Ευρετηρίαση Πολυμέσων, Κατανεμημένη Ανάκτηση)
Ημερομηνία Ανάθεσης: 5/5/2006
Ημερομηνία Παράδοσης: 5/6/2006

Άσκηση 1 (1.5 βαθμ.)

Θεωρείστε τα εξής λόγια του Σενέκα:

«Ο γνωστικός γνωρίζει τον αμαθή γιατί υπήρξε και ο ίδιος αμαθής. Ο αμαθής δεν γνωρίζει τον γνωστικό γιατί δεν υπήρξε ποτέ γνωστικός».

α) Θεωρώντας την κάθε λέξη ως σύμβολο του αλφαβήτου, ποια είναι η εντροπία του αλφαβήτου;

β) Δώστε τη συμπιεσμένη μορφή του κειμένου χρησιμοποιώντας κανονικοποιημένους κώδικες Huffman.

Άσκηση 2 (1.5 βαθμ.)

Θεωρείστε 5 έγγραφα A, B, C, D, E και έστω ότι οι αποστάσεις μεταξύ τους είναι αυτές του παρακάτω πίνακα. Δώστε το δενδρικό διάγραμμα που προκύπτει εφαρμόζοντας ιεραρχική ομαδοποίηση εγγράφων τύπου: (α) SingleLink, (β) CompleteLink, και (γ) Average Link.

A					
B	6				
C	4.2	5			
D	4	2	5		
E	2	8	5	6	
	A	B	C	D	E

Άσκηση 3 (1 βαθμ.)

Έστω ότι έχουμε 5 εικόνες A, B, C, D, E των οποίων οι αποστάσεις είναι αυτές που δίνονται στον πίνακα της Άσκησης 2. Προκειμένου να μπορούμε να απαντήσουμε επερωτήσεις γρήγορα θέλουμε να φτιάξουμε ένα μετρικό ευρετήριο, συγκεκριμένα ένα Vantage-Point-Tree (VTP). Σχεδιάστε το VTP που προκύπτει:

α) αν επιλέξουμε την εικόνα A ως κεντρική (pivot),

β) αν επιλέξουμε την εικόνα C ως κεντρική (pivot).

Άσκηση 4 (60 βαθμοί)

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

d1 = «A B Γ », d2 = «B E B»

d3 = «Δ B », d4 = «Γ E Γ»

d5 = «Δ Γ E Γ», d6 = «Γ E»

d7 = «B Δ B», d8 = «E B»

Έστω ότι τα d1,d5, d6 ανήκουν σε ένα σύστημα S1, τα d2,d4 σε ένα σύστημα S2, και τα υπόλοιπα (d3,d7,d8) σε ένα σύστημα S3. 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 S1, S2 και S3.

(β) Έστω ότι ο M έχει ήδη τα διανύσματα πηγών των S1, S2, S3 και λαμβάνει την επερώτηση $q = \text{"A Γ"}$. 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:

S1: $\langle d1, d6, d5 \rangle$

S2: $\langle d4, d2 \rangle$

S3: $\langle d7, d8, d3 \rangle$

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 τον σωστό υπολογισμό των σκορ. Δώστε το idf του κάθε όρου στην καθολική συλλογή εγγράφων.

(ε) Ο μεσίτης βρίσκει άλλο ένα σύστημα S4 το οποίο έχει την ίδια συλλογή με αυτήν του S1, δηλαδή και αυτό παρέχει πρόσβαση στα έγγραφα d1, d5, d6. Έστω ότι ο M προωθεί μια επερώτηση q στα S1 και S4 και λαμβάνει τις εξής απαντήσεις:

S1: $\langle d1, d5, d6 \rangle$

S4: $\langle d6, d1, d5 \rangle$

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 S1, S2, S3 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 S5 και S6 (τα οποία δεν έχουν καμία συλλογή εγγράφων).

Αποφασίζουν να χρησιμοποιήσουν μια συνάρτηση κατακερματισμού h των 3 bits, και έστω ότι

$h(\text{IPaddress}(S1))=1$, $h(\text{IPaddress}(S2))=2$, $h(\text{IPaddress}(S3))=4$,

$h(\text{IPaddress}(S5))=7$, $h(\text{IPaddress}(S6))=5$

Αποφασίζουν να διανείμουν το ανεστραμμένο ευρετήριο θεωρώντας κάθε όρο σαν κλειδί και έστω ότι

$h(A)=2$, $h(B)=3$, $h(\Gamma)=6$

$h(\Delta)=6$, $h(E)=5$

Δώστε (i) τους πίνακες δρομολόγησης των κόμβων S1 και S3 και (ii) πως θα κατανεμηθεί το ανεστραμμένο ευρετήριο στους κόμβους του δικτύου (δείξτε τι ακριβώς θα έχει κάθε κόμβος)