

2η Σειρά Ασκήσεων (Μοντέλα Ανάκτησης)

Ανάθεση: 18 Μαρτίου

Παράδοση: 1 Απριλίου (2 εβδομάδες)

Άσκηση 1 (2 βαθμοί)

(α) Δώστε τη διανυσματική παράσταση των εγγράφων d_1, \dots, d_5 με βάρη TF-IDF. Θεωρείστε ότι η θέση της κάθε λέξης στα διανύσματα γίνεται κατά αλφαβητική σειρά.

(β) Δώστε την απάντηση που θα έχει η κάθε επερώτηση q_1, \dots, q_4 βάσει του διανυσματικού μοντέλου.

(γ) Σχεδιάστε τη μορφή που θα έχει το ανεστραμμένο ευρετήριο για τη συλλογή D .

Documents D	Queries Q
d_1 : "a b"	q_1 : "a b"
d_2 : "a b a b"	q_2 : "a"
d_3 : "a b a b c"	q_3 : "c"
d_4 : "a b c"	q_4 : "a c"
d_5 : "a a c"	

Άσκηση 2 (2 βαθμοί)

Έστω ότι έχουμε ένα μοντέλο ανάκτησης το οποίο βλέπει τα έγγραφα και τις επερωτήσεις ως σύνολα όρων. Συγκρίνετε τις ακόλουθες συναρτήσεις διαβάθμισης

$$R_1(d, q) = \frac{|d \cap q|}{|q|}$$

$$R_2(d, q) = \frac{|d \cap q|}{|d|}$$

$$R_3(d, q) = \frac{|d \cap q|}{|d \cup q|}$$

$$R_4(d, q) = \frac{|d \cap q|}{|d| + |q|}$$

$$R_5(d, q) = |d \cap q|$$

Μπορείτε να στηρίξετε την απάντησή σας παραθέτοντας παραδείγματα.

Άσκηση 3 (2 βαθμοί)

Θεωρείστε το ακόλουθο ευρετήριο το οποίο έχει τη μορφή: term: doc1: ⟨ position1, position2, . . . ⟩; doc2: ⟨ position1, position2, . . . ⟩; κλπ.

angels : 2 : ⟨36, 174, 252, 651⟩; 4 : ⟨12, 22, 102, 432⟩; 7 : ⟨17⟩;
fools : 2 : ⟨1, 17, 74, 222⟩; 4 : ⟨8, 78, 108, 458⟩; 7 : ⟨3, 13, 23, 193⟩;
fear : 2 : ⟨87, 704, 722, 901⟩; 4 : ⟨13, 43, 113, 433⟩; 7 : ⟨18, 328, 528⟩;
in : 2 : ⟨3, 37, 76, 444, 851⟩; 4 : ⟨10, 20, 110, 470, 500⟩; 7 : ⟨5, 15, 25, 195⟩;
rush : 2 : ⟨2, 66, 194, 321, 702⟩; 4 : ⟨9, 69, 149, 429, 569⟩; 7 : ⟨4, 14, 404⟩;
to : 2 : ⟨47, 86, 234, 999⟩; 4 : ⟨14, 24, 774, 944⟩; 7 : ⟨199, 319, 599, 709⟩;
tread : 2 : ⟨57, 94, 333⟩; 4 : ⟨15, 35, 155⟩; 7 : ⟨20, 320⟩;
where : 2 : ⟨67, 124, 393, 1001⟩; 4 : ⟨11, 41, 101, 421, 431⟩; 7 : ⟨16, 36, 736⟩;

Ποιά έγγραφα (αν υπάρχουν) ικανοποιούν τα επόμενα ερωτήματα:

q_1 : “fools rush in”

q_2 : “fools rush in” AND “angels fear to tread”

Τα εισαγωγικά σηματοδοτούν phrase queries.

Άσκηση 4 (4 βαθμοί)

Θεωρείστε μια μηχανή αναζήτησης η οποία στηρίζεται σε ανεστραμμένο ευρετήριο και υποστηρίζει το Λογικό Μοντέλο (Boolean Model). Για κάθε λέξη το ευρετήριο μας δίνει μια μέθοδο που μας επιστέφει όλα τα αναγνωριστικά των εγγράφων που την περιέχουν σε αύξουσα σειρά. Θεωρείστε τις εξής επερωτήσεις:

q_1 : μέλι AND ουρανός

q_2 : μέλι AND ουρανός AND πορτοκάλι

q_3 : (πορτοκάλι OR δέντρα) AND (μέλι OR ουρανός) AND (όραση OR μάτια)

Επίσης θεωρείστε τα εξής μεγέθη των λιστών εμφάνισης (posting list sizes):

Όρος	Μέγεθος Λίστας
δέντρα	316812
μάτια	213312
μέλι	107913
όραση	87009
ουρανός	271658
πορτοκάλι	46653

(α) Δώστε ένα γρήγορο αλγόριθμο για την αποτίμηση του q_1

(β) Δώστε ένα γρήγορο αλγόριθμο για την αποτίμηση του q_2 (εστιάστε στη σειρά με την οποία συμφέρει να γίνει η αποτίμηση)

(γ) Δώστε ένα γρήγορο αλγόριθμο για την αποτίμηση του q_3 (εστιάστε στη σειρά με την οποία συμφέρει να γίνει η αποτίμηση)

Σημείωση: Για τον υπολογισμό των παραπάνω μπορείτε να φτιάξετε σχετικό πρόγραμμα (χωρίς αυτό να είναι απαραίτητο).