



CS-463 Information Retrieval Systems

Μοντέλα Ανάκτησης (Retrieval Models)

Part A

Yannis Tzitzikas

University of Crete

CS-463, Spring 05

Lecture : 3

Date : 1-3-2005



Διάρθρωση

PART (A)

- Ανάκτηση και Φιλτράρισμα
- Εισαγωγή στα Μοντέλα Αντήρησης
- Κατηγορίες Μοντέλων
- Exact vs Best Match
- Τα κλασσικά μοντέλα ανάκτησης
 - Το Boolean Μοντέλο
 - Στατιστικά Μοντέλα - Βάρυνση Όρων
 - Το Διανυσματικό Μοντέλο
 - Το Πιθανοκρατικό Μοντέλο

PART (B): Εναλλακτικά μοντέλα

- (I) Συνολοθεωρητικά μοντέλα
 - Fuzzy Retrieval Model
 - Extended Boolean Model
- (II) Αλγεβρικά Μοντέλα
 - Latent Semantic Indexing
 - Neural Network Model

PART (C):

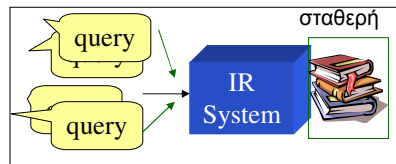
- (III) Πιθανοκρατικά Μοντέλα
 - Inference Network Model
 - Belief Network Model



Ανάκτηση και Φιλτράρισμα

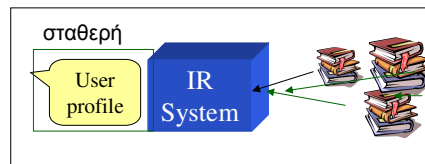
• Ανάκτηση επί σκοπού (ad hoc retrieval):

- Σταθερή συλλογή εγγράφων, μεταβαλλόμενες επερωτήσεις

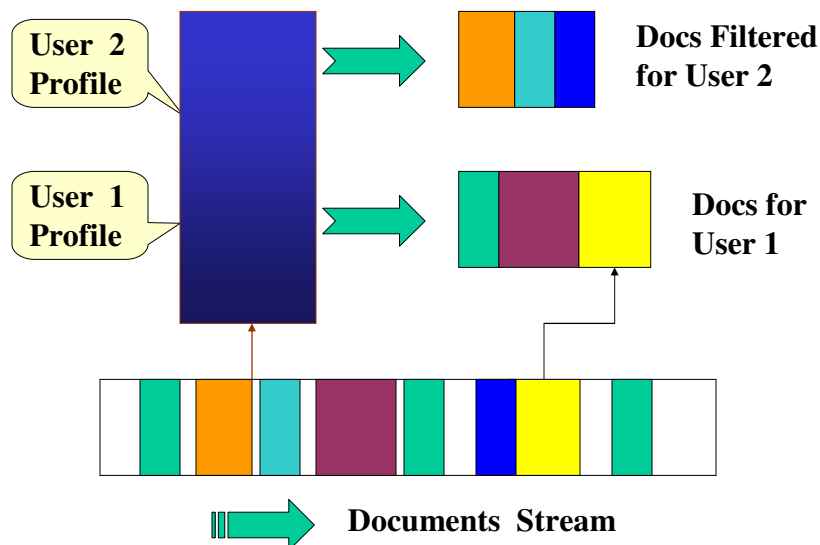


• Φιλτράρισμα(Filtering):

- Σταθερή επερώτηση, **ροή** νέων κειμένων
- **Προφίλ Χρήστη** = Επερώτηση που εκφράζει πιο μόνιμες προτιμήσεις
- έμφαση στη δημιουργία/ενημέρωση του προφίλ
- Routing: Όπως το φιλτράρισμα μόνο που εδώ το σύστημα δίδει διατεταγμένες λίστες



Φιλτράρισμα





Πως βλέπουμε ένα έγγραφο;

- Επιλογές
 - Όροι ευρετηρίασης (Index terms)
 - Πλήρες κείμενο (Full text)
 - Πλήρες κείμενο + Δομή (π.χ. hypertext, XML)



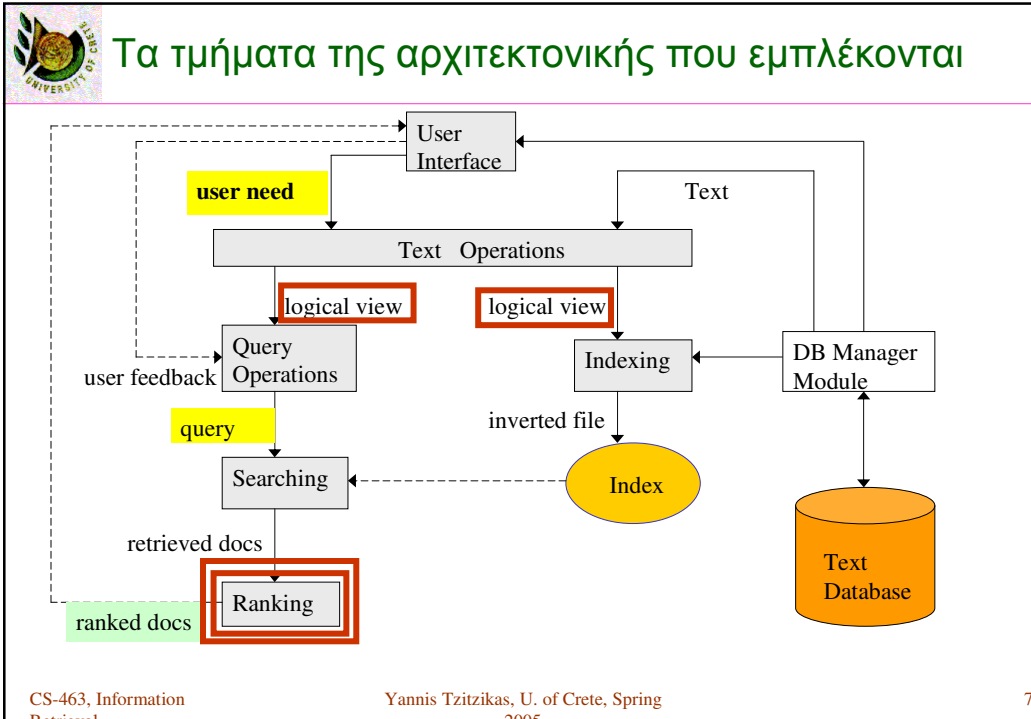
Μοντέλα Ανάκτησης

- Ένα μοντέλο ανάκτησης ορίζει
 - Αναπαράσταση Εγγράφων
 - Αναπαράσταση Επερωτήσεων
 - Καθορίζει και ποσοτικοποιεί την έννοια της συνάφειας
 - η βαθμός συνάφειας μπορεί να είναι δίτιμος (πχ. {1,0}), ή συνεχής(πχ [0,1])

Έστω **D** η συλλογή εγγράφων και **Q** το σύνολο όλων των πληροφοριακών αναγκών που μπορεί να έχει ένας χρήστης.

Μπορούμε να δούμε ένα **μοντέλο ανάκτησης πληροφορίας** ως μια τετράδα $[F, D, Q, R]$ όπου:

- **F**: πλαίσιο μοντελοποίησης εγγράφων, επερωτήσεων και των σχέσεων μεταξύ τους
- **D**: παράσταση εγγράφων $D = \{ F(d) \mid d \in D \}$
- **Q**: παράσταση επερωτήσεων $Q = \{ F(q) \mid q \in Q \}$
- **R**: συνάρτηση κατάταξης που αποδίδει μία τιμή σε κάθε ζεύγος $(d, q) \in D \times Q$
 - δίτιμη: $R: D \times Q \rightarrow \{\text{True/False}\}$
 - συνεχής $R: D \times Q \rightarrow [0,1]$



- ## Κατηγορίες Μοντέλων Ανάκτησης (I)
- **Κλασσικά Μοντέλα (3)**
 - Boolean Model
 - Διανυσματικό (Vector Space)
 - Πιθανοκρατικό (Probabilistic)
 - **Συνολοθεωρητικά (set theoretic)(2)**
 - Εκτεταμένο Boolean (Extended Boolean Model)
 - Fuzzy Model (Ασαφές Μοντέλο)
 - **Διανυσματικά (στατιστικά/αλγεβρικά) (3)**
 - Γενικευμένο Διανυσματικό (Generalized Vector Space Model)
 - Latent Semantic Indexing (Λανθάνων/Αδηλος/Υποβόσκων σημασιολογικός ευρετηριασμός)
 - Μοντέλο Νευρωνικού Δικτύου (Neural Network Model)
- CS-463, Information Retrieval
Yannis Tzitzikas, U. of Crete, Spring 2005

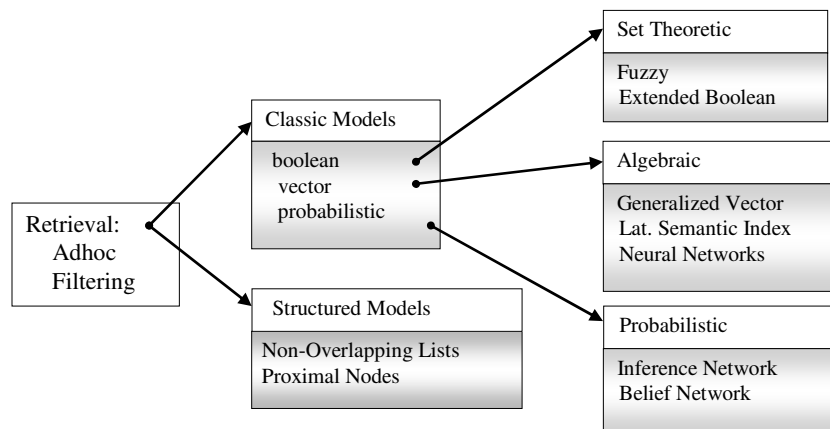


Κατηγορίες Μοντέλων Ανάκτησης (II)

- Πιθανοκρατικά (Probabilistic) (2)
 - Inference Network Model (Μοντέλο Δικτύου Επαγωγών)
 - Belief Network Model (Μοντέλο Δικτύου Πεποιθήσεων)
- Μοντέλα Βασισμένα στη Λογική
- Μοντέλα Δομημένου Κειμένου (Structured Text Retrieval Models)
 - Non-Overlapping Lists
 - Proximal Nodes
 - < Μοντέλα Ανάκτησης XML >



Μια Ταξινόμηση των Μοντέλων Ανάκτησης





Exact vs. Best Match Retrieval Models

- Exact-match
 - μια επερώτηση καθορίζει **αυστηρά κριτήρια ανάκτησης**
 - κάθε έγγραφο **είτε ταιριάζει είτε όχι** με μία επερώτηση
 - το αποτέλεσμα είναι ένα **σύνολο** κειμένων
- Best-match
 - μια επερώτηση **δεν περιγράφει αυστηρά** κριτήρια ανάκτησης
 - **κάθε** έγγραφο ταιριάζει σε μια επερώτηση **σε ένα βαθμό**
 - το αποτέλεσμα είναι μια **διατεταγμένη λίστα** εγγράφων
 - με ένα κατώφλι μπορούμε να ελέγξουμε το μέγεθος της απάντησης
- «Μικτές προσεγγίσεις»
 - E.g., some type of ranking of result set (best of both worlds)
 - E.g., best-match query language that incorporates exact-match operators



Information Retrieval Models **Boolean Retrieval Model**



Boolean Retrieval Model

- Έγγραφο = σύνολο λέξεων κλειδιών (keywords)
- Επερώτηση: Boolean έκφραση λέξεων κλειδιών (AND, OR, NOT, παρενθέσεις)
 - πχ επερώτησης
 - ((Crete AND Greece) OR (Oia AND Santorini)) AND Hotel AND-NOT Hilton
 - ((Crete & Greece) | (Oia & Santorini)) & Hotel & ! Hilton
- Απάντηση= σύνολο εγγράφων
 - απουσία διάταξης



Boolean Retrieval Model: Formally

- $K=\{k_1, \dots, k_t\}$: σύνολο όλων των λέξεων ευρετηρίασης
- Κάθε έγγραφο d_j παριστάνεται με το διάνυσμα $d_j=(w_{1,j}, \dots, w_{t,j})$ όπου:
 - $w_{i,j} = 1$ αν η λέξη k_i εμφανίζεται στο κείμενο d_j (αλλιώς $w_{i,j} = 0$)
- Μια επερώτηση q είναι μια λογική έκφραση στο K , πχ:
 - $q = \text{"k1 and (k2 or not k3)"} \text{ δηλαδή } q = \text{"k1 } \wedge \text{ (k2 } \vee \neg \text{ k3)"} \text{"}$
 - $q_{DNF} = \text{"(k1 } \wedge \text{ k2 } \wedge \text{ k3) } \vee \text{(k1 } \wedge \text{ k2 } \wedge \neg \text{ k3) } \vee \text{(k1 } \wedge \neg \text{ k2 } \wedge \neg \text{ k3)"} \text{"}$
 - $q_{DNF} = \text{"(1,1,1) } \vee \text{(1,1,0) } \vee \text{(1,0,0)"} \text{"}$
- $R(d,q)$
 - **True** αν υπάρχει συζευκτική συνιστώσα του q με λέξεις των οποίων τα βάρη είναι τα ίδια με αυτά των αντίστοιχων λέξεων του εγγράφου d
 - **False**, αλλιώς



Boolean Retrieval Model: Formally

- Πιο απλά
 - ένα κείμενο d είναι μια **σύζευξη** όρων, όπου όρος μια λέξη σε θετική ή αρνητική μορφή
 - μια επερώτηση q είναι μια οποιαδήποτε λογική έκφραση
 - $R(d,q)=\text{True}$ iff $d \models q$
 - δηλαδή αν κάθε ερμηνεία που αληθεύει το d αληθεύει και το q



Παράσταση εγγράφων κατά το Boolean Model

$$\begin{pmatrix} & k_1 & k_2 & \dots & k_t \\ d_1 & w_{11} & w_{21} & \dots & w_{t1} \\ d_2 & w_{12} & w_{22} & \dots & w_{t2} \\ \vdots & \vdots & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ d_n & w_{1n} & w_{2n} & \dots & w_{tn} \end{pmatrix} \quad w_{i,j} = \{0,1\}$$



Ακριβές Ταίριασμα (Exact Match): Συν & Πλην

- Πλεονεκτήματα
 - Προβλέψιμο, εύκολα εξηγήσιμο
 - Αποτελεσματικό όταν γνωρίζεις ακριβώς τι ψάχνεις και τι περιέχει η συλλογή
 - Αποδοτική υλοποίηση
- Αδυναμίες
 - Η διατύπωση των ερωτήσεων είναι δύσκολη για πολλούς χρήστες
 - Ικανοποιητική ακρίβεια (precision) συχνά σημαίνει απαράδεκτη ανάκληση (recall)
 - Τα μοντέλα κατάταξης (ranking models) έχουν αποδειχτεί καλύτερα στην πράξη



Τα προβλήματα του Boolean Model

- Άκαμπτο: AND σημαίνει όλα, OR σημαίνει οποιοδήποτε
- Δυσκολίες
 - η έκφραση σύνθετων πληροφοριακών αναγκών
 - ο έλεγχος του μεγέθους της απάντησης
 - All matched documents will be returned
 - η διάταξη των αποτελεσμάτων
 - All matched documents logically satisfy the query
 - η υποστήριξη ανάδρασης συνάφειας
 - If a document is identified by the user as relevant or irrelevant, how should the query be modified ?



Η αδυναμία ελέγχου του μεγέθους της απάντησης

- Παράδειγμα:
 - $|\text{Answer}(\text{"Cheap} \wedge \text{Tickets} \wedge \text{Heraklion"})| = 1$
 - $|\text{Answer}(\text{"Cheap} \wedge \text{Tickets})| = 1000$
 - $|\text{Answer}(\text{"Cheap} \wedge \text{Heraklion})| = 1000$
 - $|\text{Answer}(\text{"Tickets} \wedge \text{Heraklion"})| = 1000$
- Άρα είτε παίρνουμε μια απάντηση με ένα έγγραφο είτε ένα σύνολο 1000 στοιχείων. :(



Στατιστικά Μοντέλα

- Έγγραφο: **bag** or words (unordered words with frequencies)
 - Bag = set that allows multiple occurrences of the same element
- Επερώτηση: Σύνολο όρων με προαιρετικά βάρη:
 - Weighted query terms: **$q = \langle \text{database } 0.5, \text{ text } 0.8, \text{ information } 0.2 \rangle$**
 - Unweighted query terms: **$q = \langle \text{database text information} \rangle$**
 - No Boolean conditions specified in the query
- Απάντηση: Διατεταγμένο σύνολο συναφών εγγράφων
 - υπολογίζεται βάσει των συχνοτήτων εμφάνισης των λέξεων στα έγγραφα και στις επερωτήσεις



Κρίσιμα Ερωτήματα

- Πώς να καθορίζουμε τη **σπουδαιότητα** ενός όρου σε ένα έγγραφο και στα πλαίσια ολόκληρης της συλλογής;
- Πώς να καθορίζουμε το **βαθμό ομοιότητας** μεταξύ ενός εγγράφου και μιας επερώτησης;



Information Retrieval Models
Vector Space Model
(Διανυσματικό Μοντέλο)



Διανυσματικό Μοντέλο: Εισαγωγή

- $K=\{k_1, \dots, k_t\}$: σύνολο όλων των λέξεων ευρετηρίασης
- Κάθε έγγραφο d_j παριστάνεται με ένα διάνυσμα $d_j=(w_{1,j}, \dots, w_{t,j})$ όπου $w_{i,j} \in [0, 1]$ (πχ $w_{i,j}=0.3$)
- Μια επερώτηση q παριστάνεται με ένα διάνυσμα $q=(w_{1,q}, \dots, w_{t,q})$ όπου πάλι $w_{i,q} \in [0, 1]$
- $R(d,q)$ εκφράζει το βαθμό ομοιότητας των διανυσμάτων d και q



Βάρη Όρων: Συχνότητα όρου (tf)

- Οι πιο συχνοί όροι σε ένα έγγραφο είναι πιο σημαντικοί (υποδηλώνουν το περιεχόμενο του)
 - $freq_{ij}$ = πλήθος εμφανίσεων του όρου i στο έγγραφο j
- Κανονικοποίηση
 - $tf_{ij} = freq_{ij} / \max_k \{freq_{kj}\}$
 - όπου $\max_k \{freq_{kj}\}$ το μεγαλύτερο πλήθος εμφανίσεων ενός όρου στο έγγραφο j



Βάρη Όρων: Αντίστροφη Συχνότητα Εγγράφων (Inverse Document Frequency)

- Ιδέα: Όροι που εμφανίζονται σε πολλά διαφορετικά έγγραφα έχουν μικρή διακριτική ικανότητα
- df_i = document frequency of term i
 - πλήθος εγγράφων που περιέχουν τον όρο i
- idf_i = inverse document frequency of term i := $\log_2(N/df_i)$
 - (N : συνολικό πλήθος εγγράφων)
- Το idf αποτελεί μέτρο της διακριτικής ικανότητας του όρου
 - ο λογάριθμος ελαφραίνει το βάρος του idf σε σχέση με το tf
- Παράδειγμα:
 - Έστω $N=10$ και $df_{computer}=10$, $df_{aristotle}=2$,
 - Τότε, $N/df_{computer}=10/10=1$, $N/df_{aristotle}=10/2=5$
 - Τότε, $idf_{computer}=\log(1)=0$, $idf_{aristotle}=\log(5)=2.3$



TF-IDF Weighting (βάρυνση TF-IDF)

$$w_{ij} = tf_{ij} idf_i = tf_{ij} \log_2 (N/df_i)$$

- Ένας όρος που εμφανίζεται **συχνά** στο έγγραφο, αλλά **σπάνια** στην υπόλοιπη συλλογή, λαμβάνει **υψηλό** βάρος
- Αν και έχουν προταθεί πολλοί άλλοι τρόποι βάρυνσης, το $tf-idf$ δουλεύει πολύ καλά στην πράξη



Παράδειγμα υπολογισμού TF-IDF

- Έστω ένα έγγραφο που περιέχει όρους με τις εξής συχνότητες:
 - A(3), B(2), C(1), πχ. d="A B A B C A"
- Υποθέστε ότι η συλλογή περιέχει 10.000 έγγραφα και οι συχνότητες κειμένου (document frequencies) αυτών των όρων είναι:
 - A(50), B(1300), C(250)

Τότε:

- A: $tf=3/3$; $idf = \log(10000/50)= 5.3$; $tf-idf=5.3$
- B: $tf=2/3$; $idf = \log(10000/1300)= 2$; $tf-idf=1.3$
- C: $tf=1/3$; $idf = \log(10000/250)= 3.7$; $tf-idf=1.2$



Διάνυσμα Επερώτησης

- Τα διανύσματα των επερωτήσεων θεωρούνται ως έγγραφα και επίσης βαρύνονται με tf-idf
- Εναλλακτικά, ο χρήστης μπορεί να δώσει τα βάρη των όρων της επερώτησης



Διανυσματικό Μοντέλο:

- $K=\{k_1, \dots, k_t\}$: σύνολο όλων των λέξεων ευρετηρίασης
- Κάθε έγγραφο d_j παριστάνεται με ένα διάνυσμα $d_j=(w_{1,j}, \dots, w_{t,j})$ όπου $w_{i,j} = \mathbf{tf}_{ij} \mathbf{idf}_i$
- Μια επερώτηση q παριστάνεται με ένα διάνυσμα $q=(w_{1,q}, \dots, w_{t,q})$ όπου πάλι $w_{i,q} = \mathbf{tf}_{iq} \mathbf{idf}_i$
- $R(d,q) = ?$



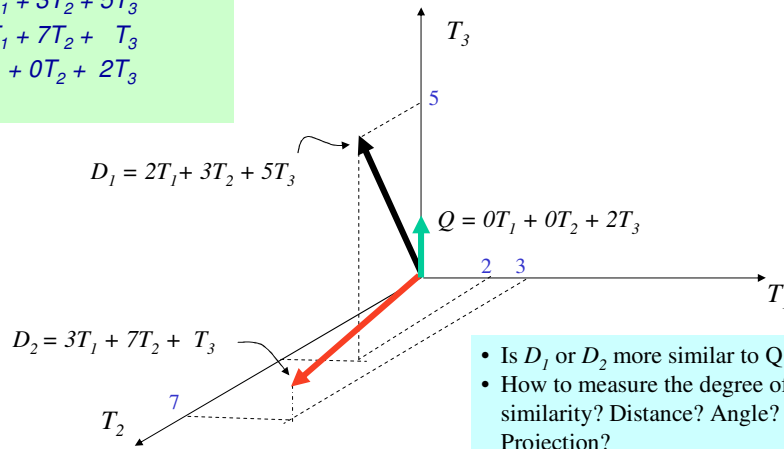
Διανυσματικό Μοντέλο: Μέτρο Ομοιότητας

Example:

$$D_1 = 2T_1 + 3T_2 + 5T_3$$

$$D_2 = 3T_1 + 7T_2 + T_3$$

$$Q = 0T_1 + 0T_2 + 2T_3$$



- Is D_1 or D_2 more similar to Q ?
- How to measure the degree of similarity? Distance? Angle? Projection?



Μέτρο Ομοιότητας: Εσωτερικό Γινόμενο (inner product)

- Η ομοιότητα μεταξύ των διανυσμάτων d και q ορίζεται ως το εσωτερικό τους γινόμενο:

$$sim(d_j, q) = \sum_{i=1}^t w_{ij} \cdot w_{iq}$$

- όπου w_{ij} το βάρος του όρου i στο έγγραφο j και w_{iq} το βάρος του όρου i στην επερώτηση
- Για δυαδικά (0/1) διανύσματα το εσωτερικό γινόμενο είναι ο αριθμός των matched query terms in the document (άρα το μέγεθος της τομής)
- Για βεβαρημένα διανύσματα, είναι το άθροισμα των γινομένων των βαρών των matched terms



Παράδειγμα

Binary:

- $d = 1, 1, 1, 0, 1, 1, 0$
- $q = 1, 0, 1, 0, 0, 1, 1$

$$sim(d, q) = 3$$

Size of vector = size of vocabulary = 7
0 means corresponding term not found in document or query

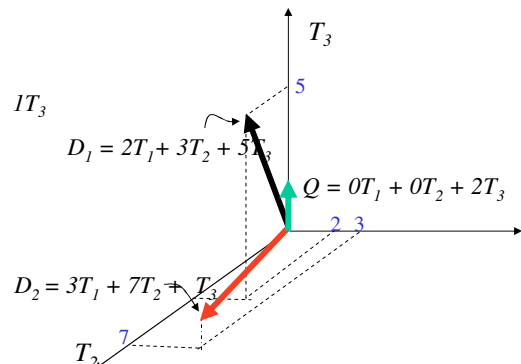
Weighted:

$$D_1 = 2T_1 + 3T_2 + 5T_3 \quad D_2 = 3T_1 + 7T_2 + 1T_3$$

$$Q = 0T_1 + 0T_2 + 2T_3$$

$$sim(D_1, Q) = 2*0 + 3*0 + 5*2 = 10$$

$$sim(D_2, Q) = 3*0 + 7*0 + 1*2 = 2$$





Ιδιότητες του Εσωτερικού Γινομένου

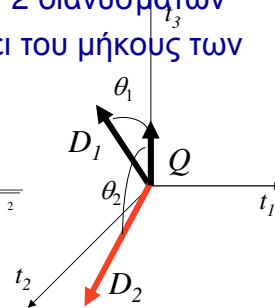
- Το εσωτερικό γινόμενο
 - δεν είναι φραγμένο (unbounded)
 - ευνοεί (μεροληπτεί) μεγάλα έγγραφα με μεγάλο πλήθος διαφορετικών όρων
 - μετρά το πλήθος των όρων που κάνουν match, αλλά αγνοεί αυτούς που δεν κάνουν match



Μέτρο Ομοιότητας Συνημίτονου (Cosine)

- Μετρά το συνημίτονο της γωνίας μεταξύ των 2 διανυσμάτων
- Εσωτερικό γινόμενο κανονικοποιημένο βάσει του μήκους των διανυσμάτων

$$\text{CosSim}(d_j, q) = \frac{\vec{d}_j \cdot \vec{q}}{|\vec{d}_j| \cdot |\vec{q}|} = \frac{\sum_{i=1}^t (w_{ij} \cdot w_{iq})}{\sqrt{\sum_{i=1}^t w_{ij}^2} \cdot \sqrt{\sum_{i=1}^t w_{iq}^2}}$$



$$D_1 = 2T_1 + 3T_2 + 5T_3 \quad \text{CosSim}(D_1, Q) = 10 / \sqrt{(4+9+25)(0+0+4)} = 0.81$$

$$D_2 = 3T_1 + 7T_2 + 1T_3 \quad \text{CosSim}(D_2, Q) = 2 / \sqrt{(9+49+1)(0+0+4)} = 0.13$$

$$Q = 0T_1 + 0T_2 + 2T_3$$

D_1 is 6 times better than D_2 using cosine similarity but only 5 times better using inner product.



Διανυσματικό Μοντέλο: Παρατηρήσεις

- **Πλεονεκτήματα**
 - Λαμβάνει υπόψη τις **τοπικές** (tf) και **καθολικές** (idf) συχνότητες όρων
 - Παρέχει **μερικό ταίριασμα** (partial matching) και **διατεταγμένα αποτελέσματα**
 - Τείνει να δουλεύει καλά στην πράξη, παρά τις αδυναμίες του
 - Αποδοτική υλοποίηση για μεγάλες συλλογές εγγράφων
- **Αδυναμίες**
 - Απουσία Σημασιολογίας (π.χ. σημασίας λέξεων)
 - Απουσία Συντακτικής Πληροφορίας (π.χ. δομή φράσης, σειρά λέξεων, εγγύτητα λέξεων)
 - Υπόθεση Ανεξαρτησίας Όρων (π.χ. άγνοια συνωνύμων)
 - Έλλειψη ελέγχου ala Boolean model (π.χ. δεν μπορούμε να απαιτήσουμε την παρουσία ενός όρου στο έγγραφο)
 - Given a two-term query $q = "A B"$, may prefer a document containing A frequently but not B, over a document that contains both A and B but both less frequently



Περίληψη του Διανυσματικού Μοντέλου

- $K = \{k_1, \dots, k_t\}$: σύνολο όλων των λέξεων ευρετηρίασης
- Κάθε **έγγραφο** d_j παριστάνεται με το διάνυσμα $d_j = (w_{1,j}, \dots, w_{t,j})$ όπου $w_{ij} = tf_{ij} idf_i = tf_{ij} \log_2 (N / df_i)$
- Μια **επερώτηση** q παριστάνεται με το διάνυσμα $q = (w_{1,q}, \dots, w_{t,q})$ όπου $w_{iq} = tf_{iq} idf_i = tf_{iq} \log_2 (N / df_i)$

$$R(d_j, q) = \text{CosSim}(d_j, q) = \frac{\bar{d}_j \cdot \bar{q}}{|\bar{d}_j| \cdot |\bar{q}|} = \frac{\sum_{i=1}^t (w_{ij} \cdot w_{iq})}{\sqrt{\sum_{i=1}^t w_{ij}^2} \cdot \sqrt{\sum_{i=1}^t w_{iq}^2}}$$



Απλοϊκή Υλοποίηση

- 1) Φτιάξε το *tf-idf* διάνυσμα για κάθε έγγραφο d_j της συλλογής (έστω V το λεξιλόγιο)
- 2) Φτιάξε το *tf-idf* διάνυσμα q της επερώτησης
- 3) Για κάθε έγγραφο d_j του D
Υπολόγισε το σκορ $s_j = \text{cosSim}(d_j, q)$
- 4) Διέταξε τα έγγραφα σε φθίνουσα σειρά
- 5) Παρουσίασε τα έγγραφα στο χρήστη

Χρονική πολυπλοκότητα: $O(|V| \cdot |D|)$

Κακό για μεγάλο V & D !

$|V| = 10,000$; $|D| = 100,000$; $|V| \cdot |D| = 1,000,000,000$



Γρηγορότερη Υλοποίηση

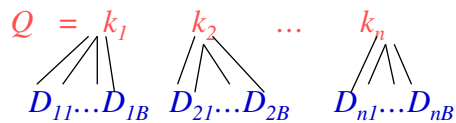
- Ένας όρος που δεν εμφανίζεται και στην επερώτηση και στο έγγραφο **δεν επηρεάζει** το βαθμό ομοιότητας συνημίτονου
 - Το γινόμενο των βαρών είναι 0 και άρα δεν συνεισφέρει στο εσωτερικό γινόμενο
- Συνήθως η επερώτηση είναι μικρή, άρα το διάνυσμα της είναι εξαιρετικά αραιό
- => Μπορούμε να χρησιμοποιήσουμε ένα ευρετήριο ώστε να υπολογίσουμε το βαθμό ομοιότητας μόνο εκείνων των εγγράφων που περιέχουν τουλάχιστον έναν όρο της επερώτησης.

3) Για κάθε έγγραφο d_j του D
Υπολόγισε το σκορ $s_j = \text{cosSim}(d_j, q)$

3') Για κάθε έγγραφο d_j που περιέχει τουλάχιστον έναν όρο του **query**
Υπολόγισε το σκορ $s_j = \text{cosSim}(d_j, q)$



Γρηγορότερη Υλοποίηση (II)



- Ας υποθέσουμε ότι ένας όρος της επερώτησης εμφανίζεται σε B έγγραφα
- Τότε η χρονική πολυπλοκότητα είναι $O(|Q| B)$
- Συνήθως είναι πολύ μικρότερη του απλοϊκού τρόπου (που είχε πολυπλοκότητα $O(|V||D|)$), διότι:
 - $|Q| \ll |V|$ και
 - $B \ll |D|$.



Διάρθρωση Διάλεξης

PART (A)

- Ανάκτηση και Φιλτράρισμα
- Εισαγωγή στα Μοντέλα Αντιληψης
- Κατηγορίες Μοντέλων
- Exact vs Best Match
- Τα κλασσικά μοντέλα ανάκτησης
 - Το Boolean Μοντέλο
 - Στατιστικά Μοντέλα - Βάρυνση Όρων
 - Το Διανυσματικό Μοντέλο
 - Το Πιθανοκρατικό Μοντέλο

PART (B): Εναλλακτικά μοντέλα

- (I) Συνολοθεωρητικά μοντέλα
 - Fuzzy Retrieval Model
 - Extended Boolean Model
- (II) Αλγεβρικά Μοντέλα
 - Latent Semantic Indexing
 - Neural Network Model

PART (C):

- (III) Πιθανοκρατικά Μοντέλα
 - Inference Network Model
 - Belief Network Model



Information Retrieval Models

Probabilistic Model



The Probabilistic Model

(πιθανοκρατικό μοντέλο)

- Στόχος:
 - σύλληψη του προβλήματος της ΑΠ με χρήση πιθανοτήτων
- Προσέγγιση
 - Υπόθεση: Για κάθε επερώτηση υπάρχει μια **ιδανική** απάντηση
 - Υπολογισμός της απάντησης βάσει των *ιδιοτήτων* (όρων) της ιδανικής απάντησης
 - Κρίσιμο ερώτημα: Ποιες είναι αυτές οι ιδιότητες;
 - Προσέγγιση: Μάντεψε αρχικά και κατόπιν βελτίωσε με επαναλήψεις



Τα βήματα

- 1/ Ένα αρχικό σύνολο εγγράφων ανακτάται με κάποιο τρόπο
- 2/ Ο χρήστης τα παρατηρεί (πχ τα πρώτα 10-20) και μαρκάρει τα συναφή
- 3/ Το ΣΑΠ αξιοποιεί αυτά τα έγγραφα για να εκλεπτύνει την περιγραφή της ιδανικής απάντησης
- 4/ Επαναλαμβάνοντας αυτή τη διαδικασία αναμένουμε ότι η περιγραφή της ιδανικής απάντησης όλο και βελτιώνεται

Παρατηρήσεις

- Πάντα πρέπει να μαντέψουμε αρχικά την περιγραφή της ιδανικής απάντησης
- Η περιγραφή της ιδανικής απάντησης γίνεται πιθανοκρατικά



Πιθανοκρατική Διάταξη

Έστω επερώτηση q και έγγραφο d_j

Το πιθανοκρατικό μοντέλο προσπαθεί να εκτιμήσει την πιθανότητα να βρει ο χρήστης το έγγραφο d_j συναφές με την επερώτηση q .

Το μοντέλο κάνει την υπόθεση ότι αυτή η πιθανότητα εξαρτάται **μόνο** από το d_j και q (και όχι από τα υπόλοιπα έγγραφα)

Η ιδανική απάντηση (σύνολο συναφών εγγράφων) συμβολίζεται με R .

Ερωτήματα:

- πώς να υπολογίσουμε τις πιθανότητες;
- ποιος ο δειγματικός χώρος;



Πιθανοκρατική Διάταξη (II)

- Η πιθανοκρατική διάταξη ορίζεται ως εξής:
 - $\text{sim}(dj, q) = P(dj \text{ relevant-to } q) / P(dj \text{ non-relevant-to } q)$
 - This is the odds of the document dj being relevant
 - Taking the odds minimize the probability of an erroneous judgement
- Ορισμοί:
 - $w_{ij} \in \{0, 1\}$
 - Μία επερώτηση είναι ένα **σύνολο** όρων.
 - $P(R | \text{vec}(dj))$: πιθανότητα το doc να είναι συναφές (να ανήκει στην ιδανική απάντηση R)
 - $P(\neg R | \text{vec}(dj))$: πιθανότητα το doc να **μην** είναι συναφές



Πιθανοκρατική Διάταξη (III)

$$\text{sim}(dj, q) = P(R | \text{vec}(dj)) / P(\neg R | \text{vec}(dj))$$

Bayes Rule:

$$P(H|e) = \frac{P(e|H) P(H)}{P(e)}$$

$$= \frac{[P(\text{vec}(dj) | R) * P(R)]}{[P(\text{vec}(dj) | \neg R) * P(\neg R)]} \quad (\text{Bayes' rule})$$

$$\sim \frac{P(\text{vec}(dj) | R)}{P(\text{vec}(dj) | \neg R)} \quad (\text{υποθέτοντας } P(R)=P(\neg R))$$

όπου:

- $P(R)$: πιθανότητα ένα τυχαίο επιλεγμένο έγγραφο να είναι συναφές
- $P(\neg R)$: πιθανότητα ένα τυχαίο επιλεγμένο έγγραφο να μην είναι συναφές
- $P(\text{vec}(dj) | R)$: πιθανότητα επιλογής του dj από το σύνολο R των συναφών



Πιθανοκρατική Διάταξη (IV)

•Ερώτημα:

- Πως μπορούμε να υπολογίσουμε το $P(\text{vec}(\mathbf{d}_j) | R)$;
- Έστω $\text{vec}(\mathbf{d}_j) = \langle 1, 0, 1, 0 \rangle$, ποια η τιμή του $P(\text{vec}(\mathbf{d}_j) | R)$;

•Τρόπος:

- Έστω $P(k_i | R)$: η πιθανότητα εμφάνισης του όρου k_i σε ένα τυχαία επιλεγμένο έγγραφο από το σύνολο των συναφών εγγράφων R
- Με υπόθεση ανεξαρτησίας όρων παίρνουμε:

$$P(\text{vec}(\mathbf{d}_j) | R) =$$

$$P(\langle 1, 0, 1, 0 \rangle | R) =$$

$$P(k_1 | R) * P(\neg k_2 | R) * P(k_3 | R) * P(\neg k_4 | R)$$



Πιθανοκρατική Διάταξη (V)

Αρα:

$$\text{sim}(\mathbf{d}_j, q) \sim \frac{P(\text{vec}(\mathbf{d}_j) | R)}{P(\text{vec}(\mathbf{d}_j) | \neg R)}$$

(assuming term independence)

$$\sim \frac{[\prod_{w_{ij}=1} P(k_i | R)] * [\prod_{w_{ij}=0} P(\neg k_i | R)]}{[\prod_{w_{ij}=1} P(k_i | \neg R)] * [\prod_{w_{ij}=0} P(\neg k_i | \neg R)]}$$



Πιθανοκρατική Διάταξη (VI)

$$\bullet \text{ sim}(d_j, q) \sim \log \frac{[\prod P(k_i | R)] * [\prod P(\neg k_j | R)]}{[\prod P(k_i | \neg R)] * [\prod P(\neg k_j | \neg R)]}$$

$$\sim \sum w_{iq} * w_{ij} * \left(\log \frac{P(k_i | R)}{P(\neg k_i | R)} + \log \frac{P(k_i | \neg R)}{P(\neg k_i | \neg R)} \right)$$

όπου $P(\neg k_i | R) = 1 - P(k_i | R)$
 $P(\neg k_i | \neg R) = 1 - P(k_i | \neg R)$

•Ερώτημα: από πού ήρθε το w_{iq} ?



Η Αρχική Διάταξη

- Ποιές είναι οι τιμές των $P(k_i | R)$ και $P(k_i | \neg R)$?
- Έστω
 - N το πλήθος των εγγράφων, και
 - n_i το πλήθος των εγγράφων που περιέχουν τον όρο k_i
- Κάνουμε μια εκτίμηση βασισμένη στις υποθέσεις:
 - $P(k_i | R) = 0.5$
 - $P(k_i | \neg R) = \frac{n_i}{N}$ (το ποσοστό των εγγράφων του έχουν τον όρο k_i)
 - Χρησιμοποιούμε αυτό το αρχικό μάντεμα για να ανακτήσουμε μια αρχική διάταξη (εγγράφων που περιέχουν τους όρους της επερώτησης)
 - Κατόπιν βελτιώνουμε τη διάταξη



Βελτίωση της Αρχικής Διάταξης (χωρίς επέμβαση του χρήστη)

- Έστω
 - V : το σύνολο των εγγράφων που ανακτήθηκαν αρχικά
 - V_i : το υποσύνολο των εγγράφων του V που περιέχουν τον όρο k_i
- Επανατοίμηση:
 - $P(k_i | R) = \frac{V_i}{V}$ (το ποσοστό των ανακτημένων που περιέχουν το k_i)
 - $P(k_i | \neg R) = \frac{n_i - V_i}{N - V}$ (μη ανακτημένα έγγραφα που περιέχουν τον k_i
(μη ανακτημένα έγγραφα))
- Επανέλαβε αναδρομικά



Βελτίωση της Αρχικής Διάταξης

- Για να μην έχουμε πρόβλημα στην περίπτωση που $V=1$ και $V_i=0$, ορίζουμε:
 - $P(k_i | R) = \frac{V_i + 0.5}{V + 1}$
 - $P(k_i | \neg R) = \frac{n_i - V_i + 0.5}{N - V + 1}$
- Εναλλακτικά:
 - $P(k_i | R) = \frac{V_i + n_i/N}{V + 1}$
 - $P(k_i | \neg R) = \frac{n_i - V_i + n_i/N}{N - V + 1}$



Συν και Πλην του Πιθανοκρατικού Μοντέλου

- **Συν:**
 - Διάταξη εγγράφων ως προς την **πιθανότητα** συνάφειας τους
- **Πλην**
 - ανάγκη μαντέματος των $P(k_i | R)$ για κάθε k_i
 - δεν λαμβάνονται υπόψη τα tf και idf



Σύντομη Σύγκριση των Κλασσικών Μοντέλων

- Το Boolean model δεν υποστηρίζει μερικό ταίριασμα και είναι το πιο αδύνατο από τα κλασσικά μοντέλα
- Μια σειρά πειραμάτων απέδειξε ότι γενικά το **διανυσματικό μοντέλο είναι αποτελεσματικότερο** του πιθανοκρατικού [Salton & Buckley]
- Η ερευνητική κοινότητα φαίνεται να συρρέει στην άποψη