



HY463 - Συστήματα Ανάκτησης Πληροφοριών Information Retrieval (IR) Systems

Προχωρημένες Λειτουργίες Επερώτησης Advanced Query Operations

Γιάννης Τζιζίκας

Διάλεξη : 11

Ημερομηνία :



Διάρθρωση Διάλεξης

- Κίνητρο
- **Ανάδραση Συνάφειας (Relevance Feedback)**
- **Αναδιατύπωση Επερωτήσεων (Query Reformulation)**
 - Αναβάρυνση Όρων (Term Reweighting)
 - Επέκταση (Διαστολή) Επερώτησης (Query Expansion),
 - Αναδιατύπωση Επερωτήσεων για το Διανυσματικό Μοντέλο
 - Optimal Query, Rocchio Method, Ide Method, DeHi Method
 - Η έννοια του Optimal (or Best) Query
 - Αξιολόγηση
- **Ψευδο-ανάδραση συνάφειας (Pseudo relevance feedback)**
- **Επέκταση Επερωτήσεων**
 - Αυτόματη Τοπική (Επιτόπια) Ανάλυση (Automatic Local Analysis)
 - Καθολική Ανάλυση
 - Επέκταση Επερώτησης βάσει Θησαυρού (Thesaurus-based Query Expansion)
 - Αυτόματη Καθολική Ανάλυση (Automatic Global Analysis)
 - Στατιστικοί Θησαυροί (Statistical Thesaurus)
 - Κατασκευή Θησαυρών
- **//Γενετικοί Αλγόριθμοι**



Κίνητρο

- Έχει παρατηρηθεί ότι οι χρήστες των ΣΑΠ δαπανούν πολύ χρόνο αναδιατυπώνοντας την αρχική τους επερώτηση προκειμένου να βρουν ικανοποιητικά έγγραφα
- Πιθανές αιτίες
 - ο χρήστης δεν γνωρίζει το περιεχόμενο των υποκείμενων εγγράφων
 - το λεξιλόγιο του χρήστη μπορεί να διαφέρει από αυτό της συλλογής
 - η αρχική επερώτηση μπορεί να είναι πιο γενική ή πιο ειδική από αυτή που θα έπρεπε (καταλήγοντας είτε σε πάρα πολλά ή σε πολύ λίγα έγγραφα)
- Η αρχική επερώτηση μπορεί να θεωρηθεί ως η πρώτη προσπάθεια έκφρασης της πληροφοριακής ανάγκης του χρήστη
- Ανάγκη για τεχνικές αντιμετώπισης αυτού του προβλήματος



Τρόποι Αντιμετώπισης

- (1) Βελτίωση της αρχικής επερώτησης**
- (2) Χρήση Προφίλ Χρήστη**
- (3) Βελτίωση παράστασης κειμένων**
- (4) Βελτίωση αλγορίθμου (μοντέλου) ανάκτησης**

Παρατηρήσεις

- Τα (2) ,(3),(4) έχουν πιο μόνιμο αποτέλεσμα (επηρεάζουν την απάντηση και των επόμενων επερωτήσεων)
- Εδώ θα εστιάσουμε στο (1)



Τεχνικές Βελτίωσης της Αρχικής Επερώτησης

Κατηγορίες:

- (α) τεχνικές που απαιτούν **είσοδο από τον χρήστη**
- (β) τεχνικές που **δεν απαιτούν** είσοδο
 - (β1) που βασίζονται στα **κορυφαία έγγραφα** που ανακτήθηκαν
 - (β2) που βασίζονται σε **όλα τα έγγραφα** της συλλογής



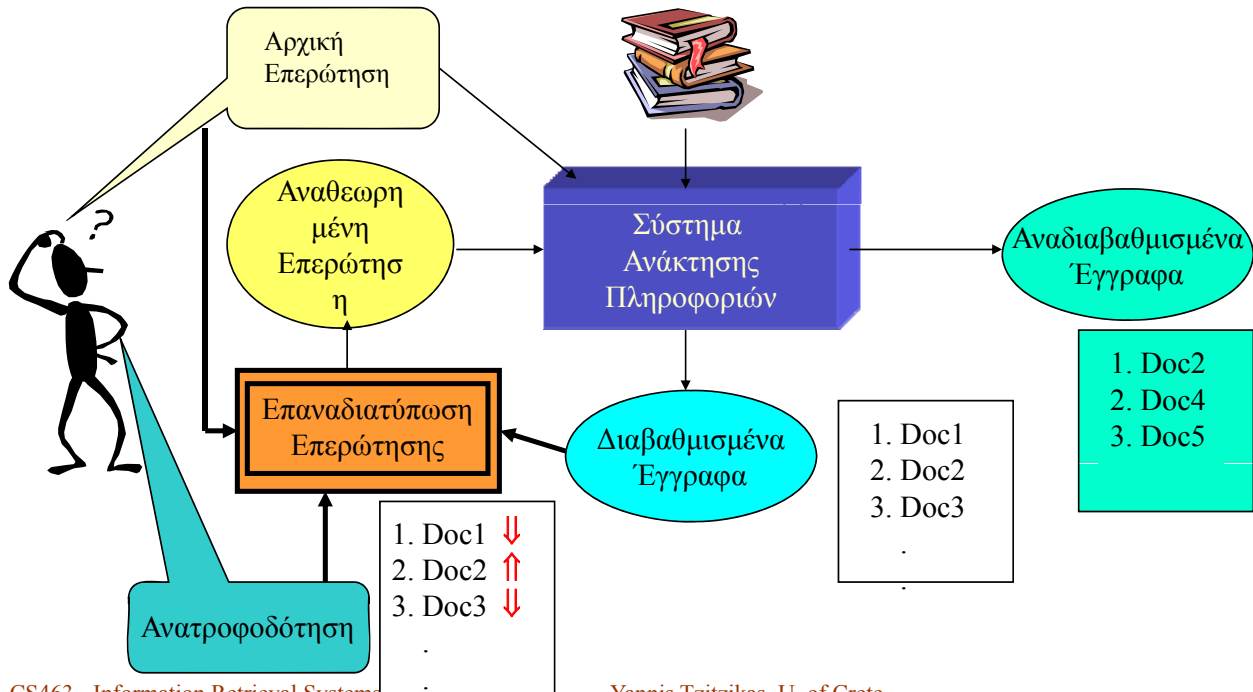
Ανάδραση Συνάφειας (Relevance Feedback): Η βασική ιδέα

Βήματα:

- 1/ Μετά την παρουσίαση των αποτελεσμάτων, επιτρέπουμε στο χρήστη **να κρίνει (θετικά ή αρνητικά) την συνάφεια** ενός ή περισσότερων εγγράφων της απάντησης
- 2/ Αξιοποιούμε αυτήν την πληροφορία για να **αναδιατυπώσουμε** την επερώτηση
- 3/ Κατόπιν δίδουμε στο χρήστη την απάντηση της αναδιατυπωμένης επερώτησης
- 4/ Πήγαινε στο βήμα 1/



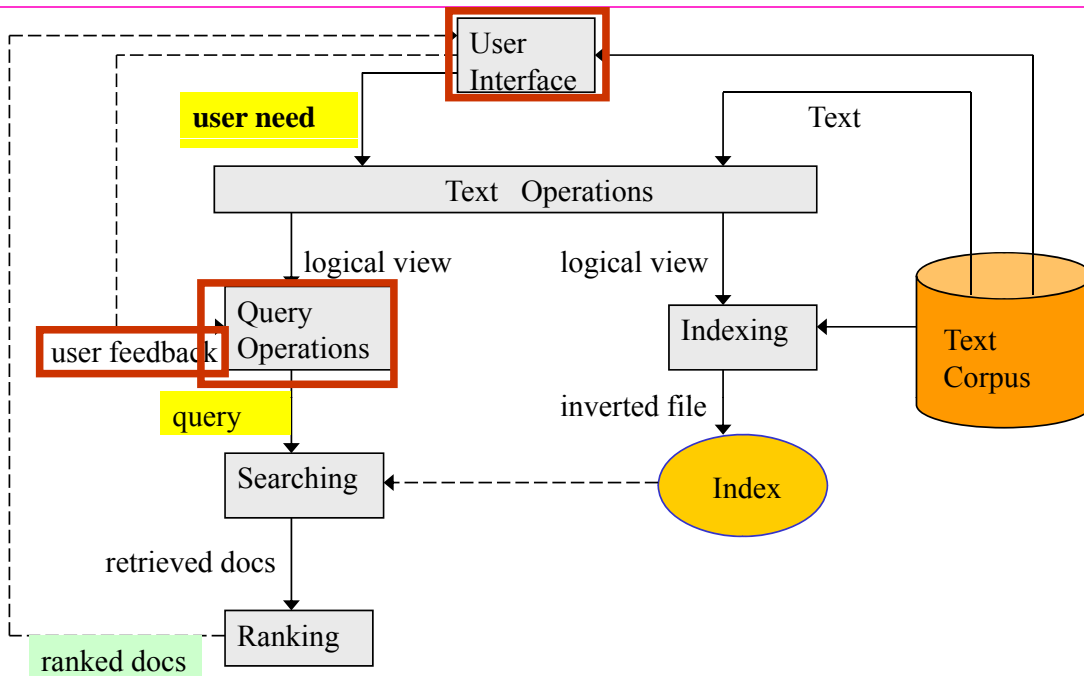
Ανάδραση Συνάφειας (Relevance Feedback): Η βασική ιδέα



7



Τμήματα της Αρχιτεκτονικής που Εμπλέκονται

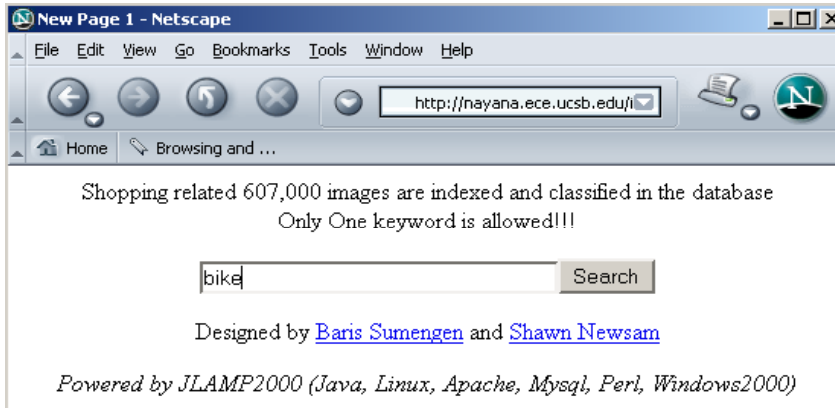


8



Παράδειγμα ανατροφοδότησης συνάφειας σε σύστημα ανάκτησης εικόνων

q=bike



(<http://nayana.ece.ucsb.edu/imsearch/imsearch.html>)



Παράδειγμα ανατροφοδότησης συνάφειας σε σύστημα ανάκτησης εικόνων

Answer("bike")=

The screenshot shows a grid of image search results. At the top, there are navigation buttons: "Browse", "Search", "Prev", "Next", and "Random". The results are arranged in two rows of six images each. Below each image is a small table of numbers representing relevance scores.

(144473, 16458) 0.0 0.0 0.0	(144457, 252140) 0.0 0.0 0.0	(144456, 262857) 0.0 0.0 0.0	(144456, 262863) 0.0 0.0 0.0	(144457, 252134) 0.0 0.0 0.0	(144483, 265154) 0.0 0.0 0.0
(144483, 264644) 0.0 0.0 0.0	(144483, 265153) 0.0 0.0 0.0	(144518, 257752) 0.0 0.0 0.0	(144538, 525937) 0.0 0.0 0.0	(144456, 249611) 0.0 0.0 0.0	(144456, 250064) 0.0 0.0 0.0



Παράδειγμα ανατροφοδότησης συνάφειας σε σύστημα ανάκτησης εικόνων

Μαρκάρισμα των Συναφών (η Επιθυμητών) από τον Χρήστη

Browse Search Prev Next Random

(144473, 16459)	(144457, 252140)	(144456, 263952)	(144456, 263963)	(144457, 252134)	(144483, 265154)
0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.0
(144483, 264644)	(144483, 265153)	(144518, 257752)	(144456, 263963)	(144456, 249611)	(144456, 250064)
0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0	0.0	0.0



Παράδειγμα ανατροφοδότησης συνάφειας σε σύστημα ανάκτησης εικόνων

Απάντηση της αναδιατυπωμένης απάντησης =

Browse Search Prev Next Random

(144538, 523493)	(144538, 523835)	(144538, 523529)	(144456, 253569)	(144456, 253568)	(144538, 523799)
0.54182	0.56319296	0.584279	0.64501	0.650275	0.66709197
0.231944	0.267304	0.280881	0.351395	0.411745	0.358033
0.309876	0.295889	0.303398	0.293615	0.23853	0.309059
(144473, 16249)	(144456, 249634)	(144456, 253693)	(144473, 16328)	(144483, 265264)	(144478, 512410)
0.6721	0.676901	0.676901	0.700339	0.70170796	0.70297
0.393922	0.4639	0.47645	0.309002	0.36176	0.469111
0.278178	0.211118	0.200451	0.391337	0.339948	0.233859



Αναδιατύπωση επερώτησης βάσει Ανάδρασης Συνάφειας (Relevance Feedback: Query Reformulation)

Τρόποι αναδιατύπωσης της επερώτησης μετά την ανάδραση:

- **Αναβάρυνση των Όρων (Term Reweighting):**
 - Αύξηση των βαρών των όρων που εμφανίζονται στα συναφή/επιθυμητά έγγραφα και μείωση των βαρών των όρων που εμφανίζονται στα μη-συναφή/επιθυμητά έγγραφα.
- **Επέκταση επερώτησης (Query Expansion):**
 - Προσθήκη νέων όρων στην επερώτηση (π.χ. από γνωστά συναφή έγγραφα)
- Υπάρχουν πολλοί αλγόριθμοι για επαναδιατύπωση επερώτησης



Αναδιατύπωση επερώτησης στο Διανυσματικό Χώρο Η έννοια της βέλτιστης επερώτησης

Η βέλτιστη επερώτηση (Optimal Query)

- Ας υποθέσουμε ότι γνωρίζουμε το σύνολο C_r **όλων** των συναφών (με την πληροφοριακή ανάγκη του χρήστη) εγγράφων.
- Η «καλύτερη επερώτηση» (αυτή που κατατάσσει στην κορυφή **όλα** τα συναφή έγγραφα), βάσει του διανυσματικού μοντέλου, θα ήταν:

$$\vec{q}_{opt} = \frac{1}{|C_r|} \sum_{\forall \vec{d}_j \in C_r} \vec{d}_j - \frac{1}{N - |C_r|} \sum_{\forall \vec{d}_j \notin C_r} \vec{d}_j$$

Where N is the total number of documents.

(θα αναλύσουμε περισσότερο αυτό το ζήτημα αργότερα)

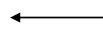
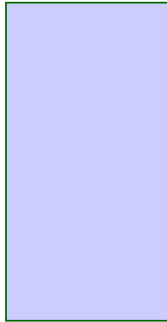
→ Αφού όμως δεν γνωρίζουμε το σύνολο C_r , θα λάβουμε υπόψη την αρχική επερώτηση και την είσοδο/ανατροφοδότηση του χρήστη.



Αναδιατύπωση επερώτησης στο Διανυσματικό Χώρο

Αφού όμως δεν γνωρίζουμε το σύνολο C_r , θα λάβουμε υπόψη την αρχική επερώτηση και την είσοδο του χρήστη.

Answer(q)= Answer (q) + user feedback =



Κόκκινα: ο χρήστης έδωσε αρνητική ανάδραση

Πράσινα: ο χρήστης έδωσε θετική ανάδραση

Μπλε: ο χρήστης δεν έδωσε ανάδραση

Τρόποι αξιοποίησης της ανατροφοδότησης του χρήστη

(I) **Rocchio** Method

(II) **Idc** Method

(III) **DeHi** Method



(I) Standard Rocchio Method

Αφού το σύνολο όλων των συναφών είναι άγνωστο, χρησιμοποίησε τα **γνωστά συναφή** (D_r) και **γνωστά μη-συναφή** (D_n) έγγραφα (από την απάντηση της αρχικής επερώτησης και βάσει της εισόδου από τον χρήστη) και επίσης συμπεριέλαβε την αρχική επερώτηση q .

Αναδιατυπωμένη επερώτηση:

$$\vec{q}_m = \alpha \vec{q} + \frac{\beta}{|D_r|} \sum_{\forall \vec{d}_j \in D_r} \vec{d}_j - \frac{\gamma}{|D_n|} \sum_{\forall \vec{d}_j \in D_n} \vec{d}_j$$

answer(q):

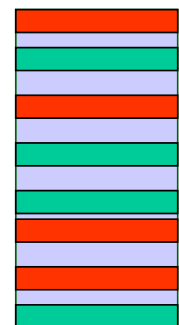
α : Tunable weight for initial query.

β : Tunable weight for relevant documents.

γ : Tunable weight for irrelevant documents.

Usually $\gamma < \beta$ (the relevant docs are more important)

If $\gamma=0$ then we have positive feedback only





(II) IDE Regular Method

Περισσότερη ανάδραση => μεγαλύτερος βαθμός αναδιατύπωσης.

Για αυτό, κατά την IDE Regular μέθοδο δεν κάνουμε κανονικοποίηση (βάσει του ποσού ανάδρασης)

$$\vec{q}_m = \alpha \vec{q} + \beta \sum_{\forall \vec{d}_j \in D_r} \vec{d}_j - \gamma \sum_{\forall \vec{d}_j \in D_n} \vec{d}_j$$

α : Tunable weight for initial query.

β : Tunable weight for relevant documents.

γ : Tunable weight for irrelevant documents.



(III) IDE “Dec Hi” Method

Τάση για απόρριψη **μόνο** των μη-συναφών εγγράφων που έχουν υψηλό σκορ

(Bias towards rejecting **just** the highest ranked of the irrelevant documents:)

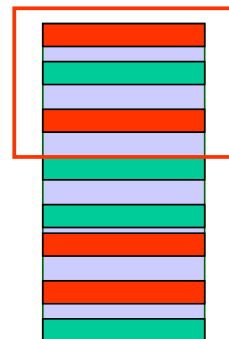
$$\vec{q}_m = \alpha \vec{q} + \beta \sum_{\forall \vec{d}_j \in D_r} \vec{d}_j - \gamma \max_{non-relevant} (\vec{d}_j)$$

α : Tunable weight for initial query.

β : Tunable weight for relevant documents.

γ : Tunable weight for irrelevant document.

answer(q):





Σύγκριση μεθόδων (I) (II) (III)

$$\vec{q}_m = \alpha \vec{q} + \frac{\beta}{|D_r|} \sum_{\forall \vec{d}_j \in D_r} \vec{d}_j - \frac{\gamma}{|D_n|} \sum_{\forall \vec{d}_j \in D_n} \vec{d}_j$$

$$\vec{q}_m = \alpha \vec{q} + \beta \sum_{\forall \vec{d}_j \in D_r} \vec{d}_j - \gamma \sum_{\forall \vec{d}_j \in D_n} \vec{d}_j$$

$$\vec{q}_m = \alpha \vec{q} + \beta \sum_{\forall \vec{d}_j \in D_r} \vec{d}_j - \gamma \max_{non-relevant} (\vec{d}_j)$$

- Γενικά, τα πειραματικά δεδομένα δεν δίνουν καθαρό προβάδισμα σε κάποια τεχνική.
- Όλες οι τεχνικές βελτιώνουν την απόδοση (recall & precision)
- Συνήθως $\alpha=\beta=\gamma=1$



Αξιολόγηση Αποτελεσματικότητας Τεχνικών Ανάδρασης Συνάφειας

Remarks

- By construction, reformulated query will rank **explicitly-marked relevant** documents higher and **explicitly-marked irrelevant** documents lower.
- When evaluating such methods, a method should not get credit for improvement on **these** documents, since it was told their relevance.
- In machine learning, this error is called “testing on the training data.”
- Evaluation should focus on generalizing to **other** un-rated documents.

Fair Process for Evaluating the Effectiveness of Relevance Feedback

- **Remove** from the corpus any document for which feedback was provided.
- Measure recall/precision performance on the remaining **residual collection**.
- *Compared to complete corpus, specific recall/precision numbers may decrease since relevant documents were removed.*
- Measure recall/precision after relevance feedback (on the residual collection)
- Relative performance on the residual collection provides fair data on the effectiveness of relevance feedback



Relevance Feedback Evaluation

TABLE 4. Evaluation of typical relevance feedback methods for five collections (weighted documents, weighted queries).

Relevance Feedback Method	Rank of Method and Avg Precision	CACM 3204 docs 64 queries	CISI 1460 docs 112 docs	CRAN 1397 docs 225 queries	INSPEC 12684 docs 84 queries	MED 1033 docs 30 queries	Average
Initial Run (reduced collection)		.1459	.1184	.1156	.1368	.3346	
Idle (dec hi)							
expand by	Rank	+49%	+44%	+92%	+32%	+79%	+59%
Rocchio (standard $\beta = .75, \alpha = .25$)							
expand by all terms	Rank	2	39	8	14	17	16
	Precision	.2552	.1404	.2955	.1821	.5630	
expand by most common terms	Improvement	+75%	+19%	+156%	+33%	+68%	+70%
	Rank	3	12	12	10	24	
	Precision	.2491	.1623	.2534	.1861	.5279	
	Improvement	+71%	+37%	+119%	+36%	+55%	+64%
Probabilistic (adjusted revised derivation)							

Simulated interactive retrieval consistently outperforms non-interactive retrieval (70% here).



Relevance Feedback Evaluation: Case Study

Example of evaluation of interactive information retrieval [Koenemann & Belkin 1996]

Goal of study: show that relevance feedback improves retrieval effectiveness

Details

- 64 novice searchers (43 female, 21 male, native English)
- TREC test bed (Wall Street Journal subset)
- Two search topics
 - Automobile Recalls
 - Tobacco Advertising and the Young
- Relevance judgements from TREC and experimenter
- System was INQUERY (vector space with some bells and whistles)
- Subjects had a tutorial session to learn the system
- Their goal was to keep modifying the query until they have developed one that gets high precision
- Reweighting of terms similar to but different from Rocchio



Rutgers INQUERY

Reset All UNDO LAST RUN QUERY Show Search Topic Text Show Tutorial EXIT RU INQUERY

Enter (next) query term below and hit <RETURN> Clear All Marks You marked 0 documents

Current Query Has 4 term(s):
 automobil* manufactur*
 car*
 defect*
 recal*

1. GM Plans to Recall 62,000 1988-89 Cars With Quad 4 Engines
 2. GM, Ford Recall Vehicles to Repair Defective Parts ---- By Neal Templin S
 3. Isuzu Motors, Honda Commence Car Recalls ---- A Wall Street Journal News I
 4. Ford and GM Recall Series Of Pickup Trucks, Coupes
 5. General Motors Corp. Recalls 196,000 Cars For Defective Brakes

Total of 6747 documents retrieved Jump to rank: _____

Document # 1 of 6747

GM Plans to Recall
 62,000 1988-89 Cars
 With Quad 4 Engines

WSJ900413-0013
 04/13/90 WALL STREET JOURNAL (J), PAGE B2

DETROIT -- General Motors Corp. said it is recalling 62,000 1988-89 model cars equipped with its high-tech Quad 4 engine to fix defective fuel lines linked to 24 engine fires. GM said the 1988-89 Pontiac Grand Am, Oldsmobile Cutlass Calais and Buick Skylark cars equipped with the 16-valve, four-cylinder Quad 4 engine have fuel lines that could crack or separate from the engines. Although GM has received reports of 24 fires caused by leaks attributable to the faulty fuel lines, a spokesman says the company knows of no injuries resulting from the incidents. GM sold about 312,000 cars equipped with Quad 4 engines in the 1988-89 model years.

In another action, GM said it is recalling about 3,200 of its 1990 Oldsmobile Cutlass Calais and Buick Skylark models to fix fuel-line defects on three engines: the Quad 4, 3.3-liter V-6, and 2.5-liter four cylinder. GM isn't aware of any fires or injuries related to the fuel line problems in this group of cars, the spokesman said.

All repairs will be done free of charge to owners, the company said.

Separately, the U.S. sales arm of Volkswagen AG's Audi subsidiary said it is recalling 1,600 1990-model Audi 80, 90 and Coupe Quattro luxury cars to replace a defective bolt in the assembly that locks the steering when the car is parked. The defective bolt could break, causing the steering wheel to remain locked even after the driver starts the car and begins

Credit: Marti Hearst



Evaluation: Precision vs. RF condition (from Koenemann & Belkin 96)

Criterion: p@30 (precision at 30 documents)

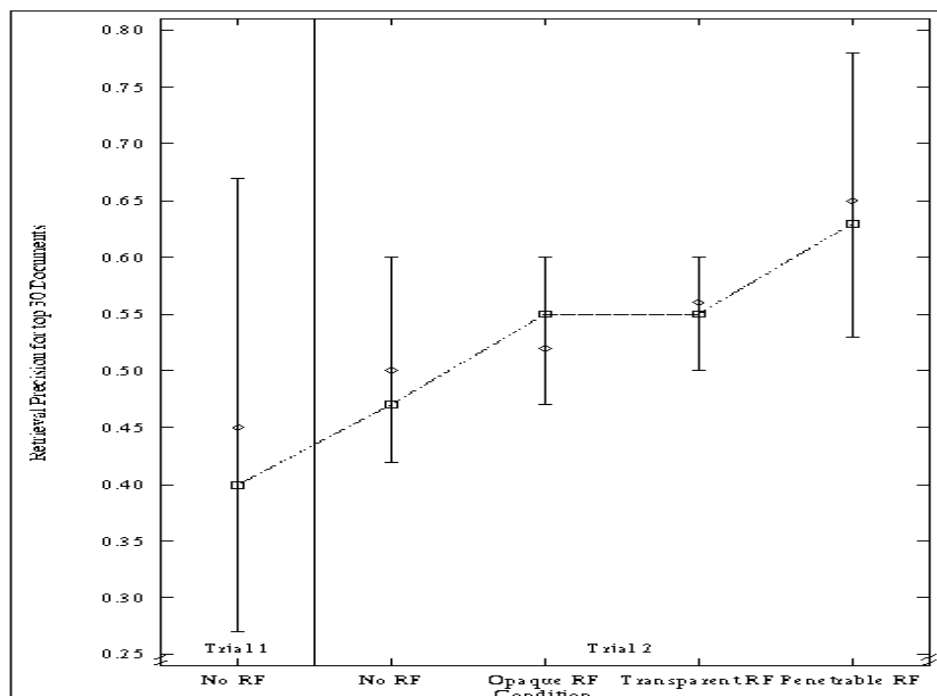
Compare:

- p@30 for users with relevance feedback
- p@30 for users without relevance feedback

Goal: show that users with relevance feedback do better

Results:

- Subjects with relevance feedback had 17-34% better performance
- But: Difference in precision numbers not statistically significant. Search times approximately equal





A1: User has sufficient knowledge for formulating the initial query.

- However:
 - User does not always have sufficient initial knowledge.
 - Examples: Misspellings, Mismatch of searcher's vocabulary vs collection vocabulary.

A2: Relevance prototypes are “well-behaved”.

- Either: All relevant documents are similar to a single prototype.
- Or: There are different prototypes, but they have significant vocabulary overlap.
- However:
 - There are several relevance prototypes.



- Οι χρήστες συχνά διστάζουν να δώσουν είσοδο
- Η ανάδραση έχει ως αποτέλεσμα μεγάλες επερωτήσεις των οποίων ο υπολογισμός απαιτεί περισσότερο χρόνο
 - σε σύγκριση με τις συνηθισμένες επερωτήσεις που διατυπώνουν οι χρήστες οι οποίες αποτελούνται από 2-3 λέξεις
 - (search engines process lots of queries and allow little time for each one)
- Μερικές φορές η νέα απάντηση περιέχει έγγραφα τα οποία δεν μπορούμε να καταλάβουμε πως προέκυψαν



Ανάδραση Συνάφειας στον Παγκόσμιο Ιστό

The screenshot shows a Google search interface. The search term is 'Information Retrieval'. Below the search bar, there are navigation links: 'Ιστός', 'Εικόνες', 'Υμνοί', and 'Καταλόγος'. The search results show 'Αποτελέσματα 1 - 10 από περίπου 6.270.0'. The first result is 'INFORMATION RETRIEVAL' by CJ Rijsbergen, University of Glasgow. Below the title, there are two links: 'Αποθηκευμένη Σελίδα' and 'Παρόμοιες σελίδες', which is circled in red to highlight the 'similar/related pages' feature.

- Some search engines offer a similar/related pages feature (simplest form of relevance feedback)
 - Πολλές φορές ο υπολογισμός αυτών των όμοιων/σχετικών σελίδων δεν γίνεται βάσει του περιεχομένου αλλά βάσει της δομής του γράφου (θυμηθείτε την ανάλυση συνδέσμων). Ο υπολογισμός είναι αρκετά πιο γρήγορος.
- But some don't because it's hard to explain to average user.
 - “Excite” initially had true relevance feedback, but abandoned it due to lack of use.



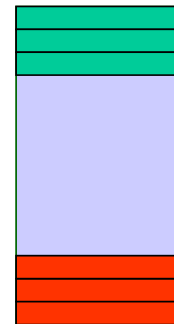
Ψευδοανάδραση Συνάφειας



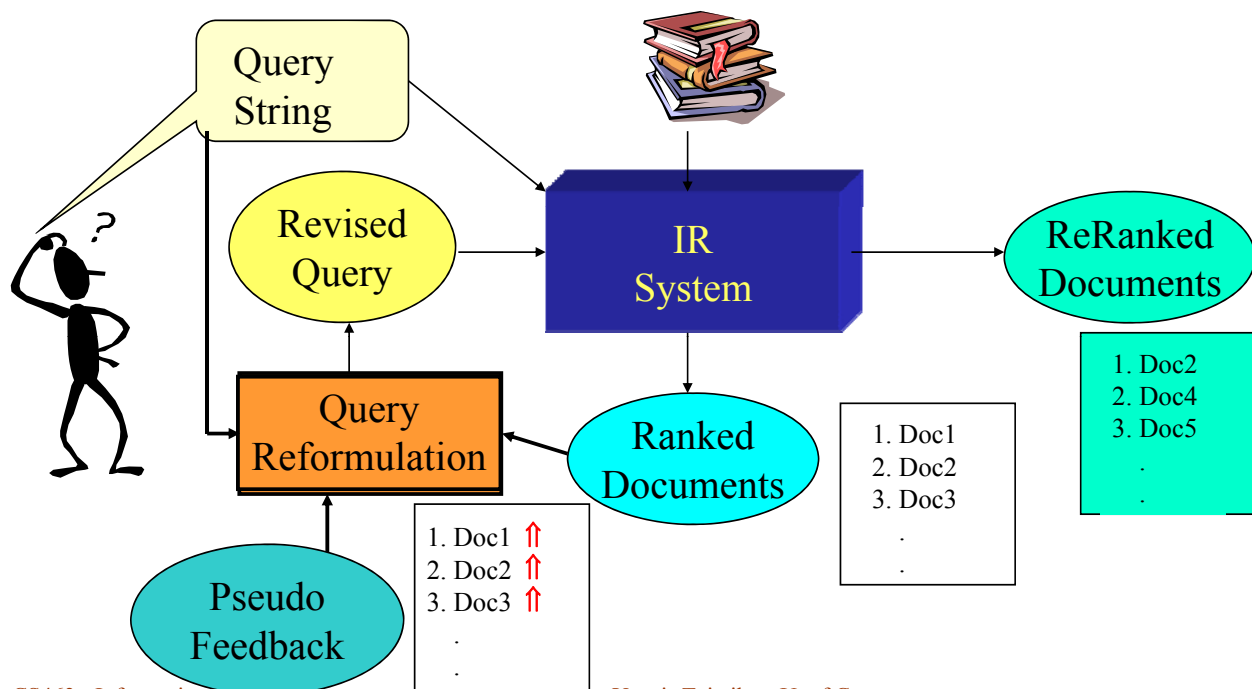
Ψευδοανάδραση Συνάφειας Pseudo Relevance Feedback

- Χρήση μεθόδων ανάδρασης αλλά **χωρίς είσοδο από το χρήστη**
- **Υπόθεση** ότι τα **κορυφαία m** από τα ανακτημένα έγγραφα είναι συναφή (και χρήση αυτών για ανάδραση)
 - Μπορούμε επίσης να χρησιμοποιήσουμε τα τελευταία έγγραφα για αρνητική ανάδραση
- Επιτρέπει την επέκταση της επερώτησης με όρους που σχετίζονται με τους όρους της επερώτησης

answer(q):



Ψευδοανάδραση Συνάφειας





Αξιολόγηση Ψευδοανάδρασης

- Βρέθηκε να βελτιώνει την απόδοση στο διαγωνισμό του TREC (ad-hoc retrieval task)
- Δουλεύει ακόμα καλύτερα αν τα κορυφαία έγγραφα πρέπει να ικανοποιούν και μια boolean έκφραση προκειμένου να χρησιμοποιηθούν για ανάδραση
 - (π.χ. να περιέχουν όλους του όρους της επερώτησης)



Αναλύοντας περισσότερο την έννοια της βέλτιστης επερώτησης (optimal query)

Πηγή:

Yannis Tzitzikas and Yannis Theoharis, **Naming Functions for the Vector Space Model**, *29th European Conference on Information Retrieval, Rome 2-5 April 2007*



The Naming Problem

We can view an IR system as a function from set of Queries to set of Answers

$$S: \text{Queries} \longrightarrow \text{Answers}$$

If $q \in \text{Queries}$, $S(q)$ denotes the answer of q .

Classically IR systems are good at “solving” the equation $S(q)=A$ wrt A , i.e.:

$$S(q) = ?$$

The **naming problem** is the problem of solving the equation wrt q , i.e. :

$$S(?) = A$$

We can distinguish two formulations of the naming problem:

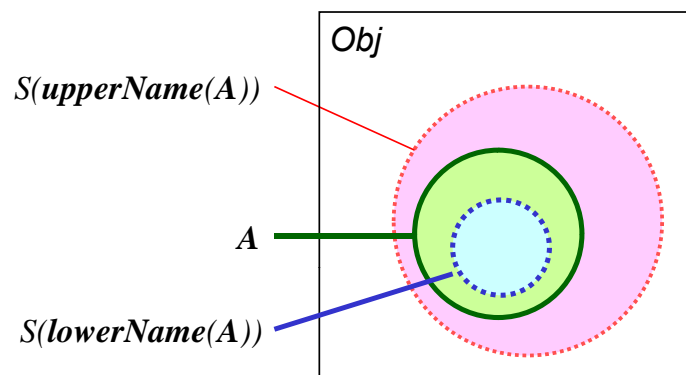
- For **unordered** answers (here A is a subset of Obj)
 - For **ordered** answers (here A is an ordered subset of Obj).
- where Obj is the set of stored objects (e.g. documents, web pages, etc)



The Naming Problem for Unordered Sets

Basic Notions

- **Exact** name
- **Upper** name
 - **Best Upper** Name
- **Lower** name
 - **Best Lower** Name
- **Relaxed** name





The Naming Problem for **Ordered Sets**

Basic Notions

- **Exact name**
- **Upper name**
 - **Best Upper Name**
- **Lower name**
 - **Best Lower Name**
- **Relaxed name**

Example:

$$A = \langle \mathbf{d1}, \mathbf{d2}, \mathbf{d3} \rangle$$

$$S(\mathit{exactName}(A)) = \langle \mathbf{d1}, \mathbf{d2}, \mathbf{d3}, d8, d9, \dots \rangle$$

$$S(\mathit{lowerName}(A)) = \langle \mathbf{d1}, \mathbf{d2}, d5, d7, \dots \rangle$$

$$S(\mathit{upperName}(A)) = \langle \mathbf{d1}, d9, \mathbf{d2}, d5, \mathbf{d3}, d8, \dots \rangle$$



Notations

Let A be an answer. Some notations that will be used

- $A(k)$: the **ordered** set comprising the **first k** elements of A
- $A\{k\}$: the **set** of elements that appear in $A(k)$
- $A|_F$: the **restriction** of A on the set F , i.e. the **ordered set** obtained if we **exclude** from A those elements that do not belong to F ,

Example

if $A = \langle d1, d2, d3 \rangle$ then

- $A(2) = \langle d1, d2 \rangle$
- $A\{2\} = \{d1, d2\}$
- $\{A\} = A\{|A|\} = \{d1, d2, d3\}$
- if $F = \{d1, d3\}$, then $A|_F = \langle d1, d3 \rangle$



Defining Formally Relaxed/Upper/Lower/Exact Names

A query q is a **relaxed name** of an answer A iff :

Case: A is a **set** : $|S(q)\{m\} \cap A| = j$ where $m \geq j > 0$.

Case: A is an **ordered set**: $S(q)(m)_{\{A(j)\}} = A(j)$ where $m \geq j > 0$.

- If $m=j=|A|$ then q is an **exact name**
- If $m \geq j=|A|$ then q is an **upper name** (the **best upper name** if m is the least possible)
- If $m=j < |A|$ then q is a **lower name** (the **best lower name** if m is the greatest possible)

So each query can be characterized by a pair (m,j) . We can now define an ordering over these pairs. The ordering should reflect the quality of the queries (as solutions for the naming problem).



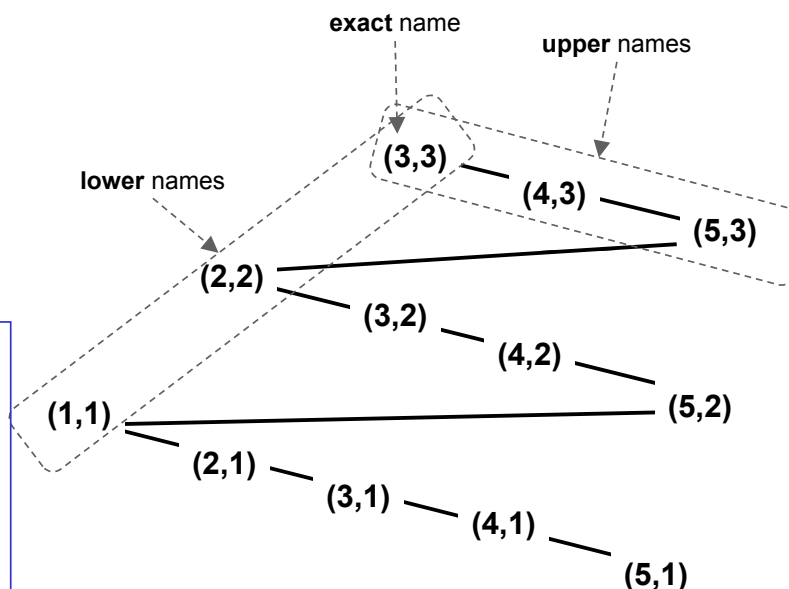
Defining a total ordering of queries

$(m,j) > (m',j')$ iff: $(j > j')$ or $(j = j'$ and $m < m')$

Let $|Obj|=5$ and $|A|=3$

Let $|Obj|=5$ and $|A|=3$.

The ordering of the corresponding pairs are:



So far, we have defined what exact, upper, lower, relaxed names are, and we defined an ordering over them. The next question is whether and how we can compute these names for a given answer A.



Investigating possible approaches for solving the Naming Problem for **Unordered Sets**

Let $A = \{d_1, d_2, d_3, \dots, d_n\}$

Possible name queries

- $q_a = \frac{1}{2}(d_i, d_j)$ where $(d_i, d_j) = \arg \max \{ \text{dist}(d, d') \mid d, d' \in A \}$
- $q_b = 1/|A| \sum \{ d \mid d \in A \}$
- $q_c = 1/|A| \sum \{ d \mid d \in A \} - 1/|\text{Obj}-A| \sum \{ d \mid d \notin A \}$

Notes

- q_a : minimizes the maximum dist from the elements of A
- q_b : minimizes the average dist from the elements of A
- q_c : is the Rocchio method (avg A , - avg $\text{Obj}-A$)

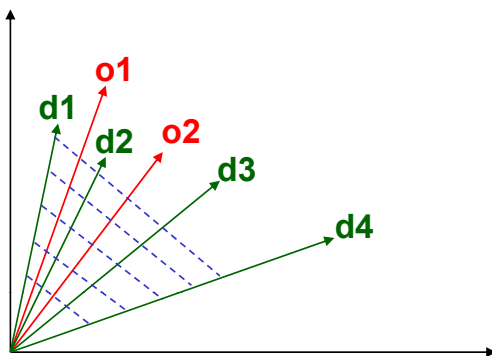
Important remarks:

- **None** is guaranteed to be the exact name (they are all **relaxed** names)
- However, if an exact name does not exist, then q_a is the **best upper** name
- Moreover, the evaluation of q_a (i.e. the computation of its answer) is **faster** than the evaluation of q_b, q_c .



The Naming Problem for **Unordered Sets** (III)

Method (a): Let $A = \{d_1, d_2, d_3, d_4\}$



1st step:

Find the pair of most distant documents

Here: (d_1, d_4)

2nd step:

Find whether other documents lie in the area specified by (d_1, d_4)

Here: $\{o_1, o_2\} \neq \emptyset$

Thus, $q_a = \frac{1}{2}(d_1, d_4)$ is an **upper name** of A
If there were no other documents then q_a would be an exact name

- **Cost**
 - 1st step: $O(|A|^2)$ computations of similarity
 - 2nd step: $O(|\text{Obj}|)$ or compute the answer of q_a and check if $S(q_a)\{A\}=A$

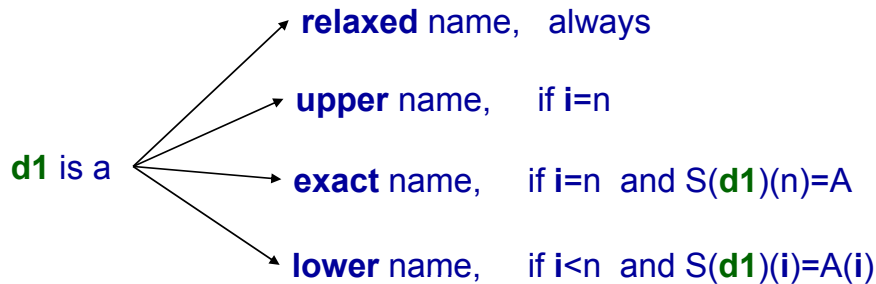


The Naming Problem for Ordered Sets

Let $A = \langle d_1, d_2, d_3, \dots, d_n \rangle$

Let i the max integer for which it holds:

$$\text{sim}(d_1, d_2) > \text{sim}(d_1, d_3) > \dots > \text{sim}(d_1, d_i)$$



Computational cost:

- $O(n)$ computations of similarity to find i .
- One query evaluation in order to decide whether d_1 is lower/exact name.



Experimental Evaluation

- Experiments conducted over the experimental web search engine GRoogle
 - <http://google.csd.uoc.gr:8080/google> (2006-2007)
- The set A was selected randomly, the average times are listed in the Table below
 - $|O|$: number of documents, $|T|$: number of terms

Collection			Naming Functions (in sec)					
			Unordered			Ordered		
$ O $	$ T $	$ \{A\} $	t_a	$t_b(\text{query terms})$	t_c	t_a	$t_b(\text{query terms})$	t_c
1K	40K	10	0.015	1.566 (5) - 3.174 (10)	0.001	0	1.878 (5) - 3.575 (10)	0.008
1K	40K	100	0.328	1.56 (5) - 3.176 (10)	0.005	0	1.624 (5) - 3.192 (10)	0.004
5K	255K	10	0.015	112.1 (5) - 262.5 (10)	0.001	0	131.2 (5) - 264.7 (10)	0.048
5K	255K	100	0.328	116.0 (5) - 251.1 (10)	0.006	0	153.4 (5) - 271.1 (10)	0.152

t_a : time to find the query (fast, depends on $|A|$, not the size of the db)

t_b : time to evaluate the query (this is the only expensive task)

t_c : time to decide what kind of name it is (fast)



Concluding Remarks and Further Research

- The naming problem has several applications (e.g. for relevance feedback or for providing a flexible interaction scheme between the system and the users)
- We expressed formally several variations of this problem and related optimality criteria
- We provided optimal solutions
 - E.g. we provided a method that certainly returns the best upper name for unordered sets (this is not true for the Rocchio method)
- We described (polynomial) algorithms for solving these problems
- Future research
 - Extend the problem statement with an additional parameter: the maximum number of words that a name could have, e.g. “find the best upper name with no more than 3 words”.
 - Indexing structures for efficient computation of names
- For more see
http://www.ics.forth.gr/~tzitzik/publications/2007_TzitzikasTheoharisNaming.pdf



Πανεπιστήμιο Κρήτης, Τμήμα Επιστήμης Υπολογιστών
Άνοιξη 2009

HY463 - Συστήματα Ανάκτησης Πληροφοριών
Information Retrieval (IR) Systems

Προχωρημένες Λειτουργίες Επερώτησης (II) Advanced Query Operations

Γιάννης Τζιτζίκας

Διάλεξη : 12

Ημερομηνία :



Επέκταση Επερώτησης (Query Expansion)

In *relevance feedback*, users give additional input (relevant/non-relevant) on documents.
In *query expansion*, users give additional input (good/bad search term) on words or phrases



Επέκταση Επερώτησης (Query Expansion)

- Τοπική Ανάλυση
 - Αναλύουμε τα (κορυφαία) έγγραφα της απάντησης
- Καθολική Ανάλυση
 - Αναλύουμε όλα τα έγγραφα της συλλογής



Επέκταση Επερώτησης (Query Expansion) Τοπική Ανάλυση (Local Analysis)



Αυτόματη Τοπική (Επιτόπια) Ανάλυση Automatic Local Analysis

- Μετά την διατύπωση της επερώτησης, ανάλυσε (στατιστικά) τις λέξεις που εμφανίζονται **μόνο** στα κορυφαία ανακτημένα έγγραφα
 - π.χ. επιλέγουμε τις 10 πιο συχνά εμφανιζόμενες λέξεις των κορυφαίων 5 εγγράφων
- Το σύστημα παρουσιάζει στο χρήστη τις πιο συχνά εμφανιζόμενες λέξεις και ο αυτός επιλέγει εκείνες που θέλει να προστεθούν στην επερώτηση
 - εναλλακτικά η επιλογή μπορεί να γίνει αυτόματα (χωρίς την παρέμβαση ή συγκατάθεση του χρήστη)
- Επίδραση στην αποτελεσματικότητα της ανάκτησης
 - Οι ασαφείς (ή αμφίσημες) λέξεις δημιουργούν λιγότερα προβλήματα (απ' ότι στην καθολική ανάλυση – την οποία θα αναλύσουμε παρακάτω)
 - Παράδειγμα: με τοπική ανάλυση η επερώτηση “Apple computer” μπορεί να επεκταθεί στην “Apple computer Powerbook laptop”



Παράδειγμα εφαρμογής

YOU ARE HERE > [Home](#) > [My InfoSpace](#) > [Meta-Search](#) > Web Search Results

Web Search Results

Your Search

Select:

[Yellow Pages](#) [White Pages](#) [Classifieds](#)

Are you looking for?

[Jacksonville Jaguars](#) [Jaguar Car](#) [Black Jaguar](#) [Jaguar Xk8](#)
[Wild Jaguars](#) [Jaguare](#) [Jaguar Accessories](#) [Jaguar Automobile](#)

Also: see altavista, teoma



YAHOO! SEARCH [Advanced Search](#)

Search Results

1 - 10 of about 45,700,000 for **Jaguar** - 0.21 sec. ([About this page](#))

Also try: [jaguar cars](#), [jaguar animal pictures](#), [jaguar parts](#), [jaguar picture](#)
[More...](#)

- SPONSOR RESULTS
- [Jaguar](#)
[www.Shopping.com](#) - Millions of Products from Thousands of Stores All in One Place.
 - [Jaguar Xk](#)
[Cars.InfoSpot1000.com](#) - Seeking **Jaguar** xk Info? See The Results You Want Now.
 - [Jaguar Cars](#)
[cars.nextag.com](#) - Compare multiple free quotes on a new car from local dealers.

1. [Jaguar](#)
Official site of the Ford Motor Company division featuring new **Jaguar** models and local dealer information.
[www.jaguar.com](#) - [More from this site](#)

SPONSOR RESULTS

[Jaguar](#)
Shop for Car Parts. Compare products, stores & prices.
[www.Dealtime.com](#)

[Jaguar Merchandise Book](#)
Buy **Jaguar** merchandee Book at SHOP.COM.
[www.SHOP.com](#)

[Jaguar Natural Spray on Cataloglink](#)
Find **Jaguar** natural spray on Cataloglink



Web | Images | Video | More >

jaguar

Advanced Search

Narrow

- [Jaguar Cars](#)
- [Black Jaguar](#)
- [Cat Jaguar](#)
- [Jaguar Big Cats](#)
- [Jaguars Habitat](#)
- [What Do Jaguars Eat](#)
- [Panthera Onca](#)
- [Where Do Jaguars Live](#) More >

Expand

- [Cheetah](#)
- [Ferrari](#) More >

Related Names

- [Ford](#)
- [Wolf](#) More >

jaguar

Showing results 1-10 of 10,190,000



[Source](#)

Jaguar | Save

Kingdom: Animalia **Phylum:** Chordata **Class:** Mammalia **Order:** Carnivora **Family:** Felidae

Genus: Panthera **Species:** Panthera onca
 The biggest and most powerful North American cat, the Jaguar is the only one that roars. It moves over a large home range with a diameter of 3 to 15 miles (5-25 km) where prey is abundant, larger where prey is scarce. This cat hunts... [More »](#)

[Key Facts](#) | [Images](#) | [Encyclopedia](#)

Jaguar

Gama actual, concesionarios, historia, noticias, anuncios y servicios financieros.

[www.jaguar.com/](#)

Jaguar (Panthera onca)

Jaguar (Panthera onca) facts, photos and videos. ... The **Jaguar** is the largest cat in the Western Hemisphere and the third largest cat in ...

[www.thebigzoo.com/Animals/Jaguar.asp](#) · [Cached](#)

Jaguar

The **jaguar** measures five to six feet from its nose to the tip of its tail and weighs 140 to 220 pounds (females are slightly smaller).

[www.kidsplanet.org/factsheets/jaguar.html](#) · [Cached](#)

Jaguar

Images



[More >](#)

Dictionary

Definitions of 'jaguar'

(jăgwăr, jăgyŭ-ărTM) - 1 definition
[The American Heritage® Dictiona](#)

jaguar (n.) A large feline mammal (Panthera onca) of Central and South America, closely related to the leopard & having a tawny coat spotted with black rosettes.

All Music Guide



Jaguar

By: **Fred Small**

Whether an artist is conservative, centrist, liberal or downright radical, there's nothing wrong with getting on a



ΣΤΟ google/mitos

A very simple technique is currently supported:

- For each term t_i that appears in the top L (by default L=5) documents returned by the Query Evaluator, we sum its term frequencies (i.e. all tf_{ij} where j in top-L documents) and we recommend to the user the S terms (by default S=5) with the highest accumulative frequency.

[Advanced Search](#)

Results per page

List of documents matching the search

You can expand your query with: ανακτηση τηρη πληροφορια assign project

[HY-463 Συστήματα Ανάκτησης Πληροφορίας - 0.0441005](#)

IR Link Analysis 5 5 5 6 Solutions Project **GRoogle** 24 ...

<http://www.csd.uoc.gr:80/-hy463/2006/en/assignments.html> - 1162813764000 - 6KB [Cached](#) [\[mark as spam\]](#)

Table 9. Query Expansion Examples

Initial Query	Expanded Terms				
1 retrieval	imag	medic	index	storag	system
2 web	system	servic	page	process	cours
3 user	interfac	layer	system	develop	softwar

Table 10. Query Expansion Average Times

L	Time (sec)
5	0.002
10	0.003
15	0.004
20	0.004



Τεχνικές αυτόματης τοπικής ανάλυσης

- Association Matrix
 - based on the co-occurrence of terms in documents
- Metric Correlation Matrix
 - based on the co-occurrence and proximity of terms in documents
- //Scalar Clusters
- //Local context analysis



(a) Association Matrix and Normalized Association Matrix

D: τα έγγραφα της απάντησης και $w_1 \dots w_n$ οι όροι που εμφανίζονται σε αυτά.

	w_1	w_2	w_3	w_n
w_1	c_{11}	c_{12}	c_{13}	c_{1n}
w_2	c_{21}				
w_3	c_{31}				
.	.				
.	.				
w_n	c_{n1}				

c_{ij} : Correlation factor between term i and term j :

$$c_{ij} = \sum_{d_k \in D} f_{ik} \times f_{jk}$$

f_{ik} : frequency of term i in document k

Normalized Association Matrix

- Frequency based correlation factor favors more frequent terms.
- Normalize association scores:

$$s_{ij} = \frac{c_{ij}}{c_{ii} + c_{jj} - c_{ij}}$$

- Normalized score is 1 if two terms have the same frequency in all documents in D.

Από αυτόν τον πίνακα μπορούμε να βρούμε τους όρους που είναι **πιο κοντά σε αυτούς της επερώτησης** (θυμηθείτε και τον πίνακα συσχέτισης στο fuzzy model του μαθήματος 4)



(b) Metric Correlation Matrix

- Association correlation does not account for the **proximity** of terms in documents, just co-occurrence frequencies within documents.
- **Metric correlations** account for term proximity.

V_i : Set of all occurrences of term i in any document in D .

$r(k_u, k_v)$: Distance in words between word occurrences k_u and k_v
 ($=\infty$ if k_u and k_v are occurrences in different documents).

$$c_{ij} = \sum_{k_u \in V_i} \sum_{k_v \in V_j} \frac{1}{r(k_u, k_v)}$$

- **Normalized Metric Correlation Matrix**
 – to account for term frequencies:

$$s_{ij} = \frac{c_{ij}}{|V_i| \times |V_j|}$$



Query Expansion with (Association or Metric) Correlation Matrix

	w_1	w_2	w_3	w_n
w_1	c_{11}	c_{12}	c_{13}	c_{1n}
w_2	c_{21}				
w_3	c_{31}				
.	.				
.	.				
w_n	c_{n1}				

- For each term i in the query q , expand query with n terms, those with the highest value of c_{ij} .
- This adds semantically related terms in the “neighborhood” of the query terms.



Αυτόματη Καθολική Ανάλυση (Automatic Global Analysis)



Αυτόματη Καθολική Ανάλυση Automatic Global Analysis

- Προσδιορισμός βαθμού ομοιότητας μεταξύ των όρων βάσει στατιστικής ανάλυσης ολόκληρης της συλλογής
 - Υπολογισμός πινάκων συσχέτισης (association matrices) που ποσοτικοποιούν την ομοιότητα μεταξύ των όρων ανάλογα με το πόσο συχνά συνεμφανίζονται
- Επέκταση επερώτησης με τους πιο όμοιους όρους.
- Επίδραση στην αποτελεσματικότητα της ανάκτησης
 - Οι ασαφείς (ή αμφίσημες) λέξεις δημιουργούν **περισσότερα προβλήματα** (απ' ότι στην τοπική ανάλυση)
 - Παράδειγμα: με καθολική ανάλυση η επερώτηση “Apple computer” μπορεί να επεκταθεί στην “Apple red fruit orange computer”
- Μια λύση:
 - Query Expansion Based on a Similarity Thesaurus



Query Expansion Based on a Similarity Thesaurus

- Βασική ιδέα
 - Οι όροι που προστίθενται στην επερώτηση καθορίζονται με βάση την απόσταση τους **από ολόκληρη την επερώτηση** (και όχι βάσει της απόστασής τους από κάθε όρο της επερώτησης ξεχωριστά)
- Στην αντίθετη περίπτωση θα είχαμε:
 - “Apple computer” → “Apple red fruit computer”
- Ενώ τώρα
 - “fruit” not added to “Apple computer” since it is far from “computer.”
 - “fruit” added to “apple pie” since “fruit” close to both “apple” and “pie.”



Query Expansion Based on a Similarity Thesaurus

- Τρόπος
 - Έστω N έγγραφα, t όροι $K=\{k_1, \dots, k_t\}$
 - Παριστάνουμε **κάθε όρο** με ένα διάνυσμα στον χώρο των N διαστάσεων
 - Είναι σαν να έχουμε αντιστρέψει το ρόλο των όρων και των εγγράφων: έχουμε λοιπόν μια διανυσματική παράσταση των όρων (κάθε έγγραφο αποτελεί μια διάσταση στο χώρο των διανυσμάτων). Προσαρμόζουμε το σχήμα βάρυνσης TF-IDF βάσει αυτής της θεώρησης.

$$\vec{k}_i = (w_{i1}, \dots, w_{iN})$$

itf: Inverse term frequency (το ανάλογο του idf):

$$w_{ij} = \frac{(0.5 + 0.5 \frac{f_{ij}}{\max_j(f_{ij})}) itf_j}{\sqrt{\sum_{l=1}^N (0.5 + 0.5 \frac{f_{il}}{\max_l(f_{il})})^2 itf_j^2}}$$

Num of terms in the collection
 $itf_j = \frac{\dots}{\dots}$
 Num of distinct terms in d_j

← ανάλογο της βάρυνσης TF*IDF μόνο που εδώ χρησιμοποιούμε το το inverse term frequency.



Query Expansion Based on a Similarity Thesaurus (II)

- Υπολογισμός ομοιότητας δυο όρων

– (π.χ. με εσωτερικό γινόμενο)

$$c_{u,v} = \vec{k}_u \cdot \vec{k}_v$$

- Τα βήματα για την επέκταση της επερώτησης

– (1) Represent query in the concept space that we used to represent terms

$$\vec{q} = \sum_{k_i \in q} w_{iq} \vec{k}_i$$

– (2) Compute $\text{sim}(q, k_u)$ for each k_u

$$\text{sim}(q, k_u) = \vec{q} \cdot \vec{k}_u$$

– (3) Expand q with the top r ranked terms. The weight of each added term k_u is set

$$w_{uq'} = \frac{\text{sim}(q, k_u)}{\sum_{k_i \in q} w_{iq}}$$

- Results

– 20% improved retrieval performance



Καθολική vs. Επιτόπια Ανάλυση

- Η καθολική ανάλυση έχει μεγάλο υπολογιστικό κόστος αλλά μόνο στην αρχή
 - υποθέτοντας ότι τα έγγραφα της συλλογής είναι σταθερά
- Η τοπική ανάλυση έχει αρκετό υπολογιστικό κόστος για κάθε επερώτηση
 - (παρόλο που το πλήθος των όρων και των εγγράφων είναι μικρότερο αυτού της καθολικής)
- Η τοπική ανάλυση δίδει καλύτερα αποτελέσματα



Επέκταση επερωτήσεων: Συμπεράσματα

- Η επέκταση των επερωτήσεων με σχετιζόμενους όρους μπορεί να βελτιώσει την αποτελεσματικότητα της ανάκτησης, ιδιαίτερα την ανάκληση (recall).
- Η αλόγιστη επιλογή σχετιζόμενων όρων μπορεί να μειώσει την ακρίβεια (precision).



Θησαυροί Όρων και Καθολική Ανάλυση



Θησαυροί Όρων

- Ένας θησαυρός παρέχει πληροφορίες για συνώνυμα και σημασιολογικά κοντινές λέξεις και φράσεις [see also Sec 7.2.5]

- Παράδειγμα:

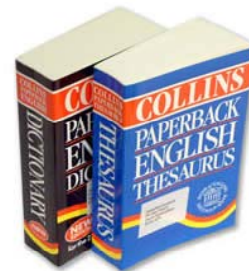
physician

syn: ||croaker, doc, doctor, MD, medical, mediciner, medico, ||sawbones

rel: medic, general practitioner, surgeon,

- Online-θησαυροί:

- **Roget's** thesaurus
- **INSPEC** thesaurus
- **WordNet** (<http://wordnet.princeton.edu/>)
- The free dictionary <http://www.thefreedictionary.com/>



Χρήσεις Θησαυρού

- Ευρετηρίαση κειμένων/βιβλίων με επιλογή όρου από θησαυρό
- Αναζήτηση χρησιμοποιώντας όρους του θησαυρού
 - (αυτόματη ή ύστερα από επιλογή του χρήστη)
- Για βελτίωση της ανάκτησης
 - Αν η απάντηση μιας επερώτησης είναι **μικρή**, μπορούμε να προσθέσουμε όρους βάσει των σχέσεων του θησαυρού (συνώνυμα, ..)
 - Αν απάντηση είναι **πολύ μεγάλη**, μπορούμε να συμβουλευτούμε το θησαυρό και να αντικαταστήσουμε κάποιους όρους της επερώτησης με πιο ειδικούς.

MyStuff | Settings

Ask.com

Web | Images | Video | More ▶

jaguar

Advanced Search

Narrow

- Jaguar Cars
- Black Jaguar
- Cat Jaguar
- Jaguar Big Cats
- Jaguars Habitat
- What Do Jaguars Eat
- Panthera Onca
- Where Do Jaguars Live

More ▶

Expand

- Cheetah
- Ferrari

More ▶

Related Names

- Ford
- Wolf

More ▶



- **Γλωσσικοί Θησαυροί**
 - Designed to assist the writer in creatively selecting vocabulary
 - Παράδειγμα:
 - **Roget's** thesaurus.
- **Θησαυροί κατάλληλοι για Information Retrieval**
 - για το «συντονισμό» των διαδικασιών ευρετηρίασης και αναζήτησης
 - σχεδιάζονται για συγκεκριμένες θεματικές περιοχές (π.χ. βιβλιογραφία πληροφορικής, αρχιτεκτονικής, κλπ)
 - Παράδειγμα:
 - **INSPEC**
 - Στην κατηγορία αυτή μπορούμε να πούμε ότι εντάσσονται και οι Οντολογίες του Σημασιολογικού Ιστού (Semantic Web Ontologies)



Παράδειγμα: **INSPEC** thesaurus (for IR)

- Πεδίο: **physics, electrical engineering, electronics, computers**
- **Τύποι Συσχετίσεων** μεταξύ δύο όρων:
 - **UF**: ακρώνυμο του **Used For (converse: USE)** // π.χ. USE X σημαίνει ότι ο X είναι ο δόκιμος όρος
 - **BT**: ακρώνυμο του **Broader Term (converse NT)**
 - **TT**: ακρώνυμο του **Top Term** (δείχνει στον κορυφαίο όρο της ιεραρχίας)
 - **RT**: ακρώνυμο του **Related Term**
- **Παράδειγμα:**
 - computer-ai ded i nstructi on
 - see also educati on
 - UF teachi ng machi nes
 - BT educati onal computi ng
 - TT computer appl i cati ons
 - RT educati on , teachi ng



Παράδειγμα: **WordNet** (<http://wordnet.princeton.edu/>)

- A detailed database of semantic relationships between English words that was developed by the famous cognitive psychologist George Miller and a team at Princeton University.
- About 144,000 English words. Nouns, adjectives, verbs, and adverbs grouped into about 109,000 synonym sets called **synsets**.

Synset Relationships (τύπτοι συσχετίσεων μεταξύ synsets)

- **Antonym:** front → back
- **Attribute:** benevolence → good (noun to adjective)
- **Pertainym:** alphabetical → alphabet (adjective to noun)
- **Similar:** unquestioning → absolute
- **Cause:** kill → die
- **Entailment:** breathe → inhale
- **Holonym:** chapter → text (part-of)
- **Meronym:** computer → cpu (whole-of)
- **Hyponym:** tree → plant (specialization)
- **Hypernym:** fruit → apple (generalization)



Παράδειγμα: **AAT** (Art and Architecture Thesaurus)

- Πεδίο: fine art, architecture, decorative art, and material culture.
- Almost 120,000 terms for objects, textual materials, images, architecture and culture from all periods and all cultures.
- Used by archives, museums, and libraries to describe items in their collections.
- Used to search for materials.
- Used by computer programs, for information retrieval, and natural language processing.



Χαρακτηριστικά Θησαυρών

- **Coordination Level (βαθμός συντονισμού)**
 - refers to the construction of phrases from individual terms
 - **Pre-coordination**: the thesaurus **contains phrases**
 - + the vocabulary is very precise
 - - the user has to be aware of the phrase construction rules, large size
 - **Post-coordination**: the thesaurus **does not contain phrases**. They are constructed while indexing/searching
 - + user does not worry about the order of the words
 - - precision may fall
- **Term Relationships**
 - **equivalence** relations (e.g. synonymy)
 - **hierarchical** relations (e.g. dogs BT animals,)
 - **nonhierarchical** relations (e.g. RT)



Χαρακτηριστικά Θησαυρών (2)

- **Number of Entries per Term**
 - preferably: a single entry for each thesaurus term
 - however homonyms does not make this possible
 - solution: usage of **parenthetical qualifiers**:
 - bonds(chemical), bonds(adhesive) // χημικός δεσμός / υλικό συγκόλλησης
- **Specificity of Vocabulary**
 - high specificity -> large vocabulary size
- **Control of Term Frequency of Class Members (for statistical thesauri)**
 - the terms of a thesaurus should have roughly equal frequencies
 - the total frequency in each class (of terms) should be equal
- **Normalization of Vocabulary**
 - terms should be in noun form
 - other rules related to singularity of terms, spelling, capitalization, abbreviations, initials, acronyms, punctuation



Τρόποι Κατασκευής Θησαυρών

[A] Χειροποίητη Δημιουργία

[B] Αυτόματη Κατασκευή

[B.1] από συλλογή κειμένων

Προϋπόθεση: Να υπάρχει μια μεγάλη και αντιπροσωπευτική συλλογή κειμένων

[B.2] από συγχώνευση άλλων θησαυρών

Προϋπόθεση: Να υπάρχουν >2 διαθέσιμοι θησαυροί για την περιοχή που μας ενδιαφέρει



[A] Χειροποίητοι Θησαυροί (διαδικασία κατασκευής)

Διαδικασία κατασκευής

1. Προσδιορισμός εμβέλειας (define subject boundaries)
2. Διαμέριση σε θεματικές ενότητες (partition into divisions and subject areas)
3. Συλλογή Όρων (collection of terms)
 - Πηγές: encyclopedias, handbooks, textbooks, journal titles, catalogues, other thesauri, subject experts, potential users
4. Ανάλυση Όρων
 - συνώνυμα, ιεραρχική δόμηση όρων, ορισμοί όρων, περιγραφή εφαρμοσιμότητας όρων (scope notes)
5. Αξιολόγηση και Αναθεώρηση (reviewing phase)
6. Παράδοση (εκτύπωση) του Θησαυρού σε ιεραρχική και αλφαβητική οργάνωση
7. Συντήρηση Θησαυρού (προσθήκη νέων όρων, αλλαγές στη δομή, κλπ)

→ Πολύ χρονοβόρα, κοπιαστική και ακριβή διαδικασία



Επέκταση επερωτήσεων βάσει Θησαυρού Thesaurus-based Query Expansion

- Τρόπος:
 - Για κάθε όρο t της επερώτησης, πρόσθεσε στην επερώτηση τα συνώνυμα και τις σχετικές λέξεις (related terms) του t
 - Τα βάρη των νέων λέξεων μπορεί να είναι **χαμηλότερα** των βαρών των λέξεων της αρχικής επερώτησης
 - E.g. of a WordNet-based Query Expansion
 - Add synonyms in the same synset.
 - Add hyponyms to add specialized terms.
 - Add hypernyms to generalize a query.
 - Add other related terms to expand query.
- Αποτέλεσμα
 - **Αυξάνει** την ανάκληση (recall.)
 - Μπορεί να **μειώσει** την ακρίβεια (precision), ιδιαίτερα όταν η επερώτηση περιέχει αμφίσημες λέξεις. Παράδειγμα μη επιτυχούς επέκτασης:
 - “interest rate” → “interest rate fascinate evaluate”



[B1] Αυτόματη Κατασκευή Θησαυρών από Κείμενα

- Η κατασκευή (από ανθρώπους) ενός θησαυρού είναι πολύ **χρονοβόρα** και δεν υπάρχουν θησαυροί για όλες τις γλώσσες
- Οι πληροφορίες που μπορούμε να χρησιμοποιήσουμε από έναν θησαυρό περιορίζονται στις σχέσεις που είναι ρητώς καταγεγραμμένες στον θησαυρό
- **Ιδέα: Μπορούμε να ανακαλύψουμε σημασιολογικές σχέσεις μεταξύ λέξεων αναλύοντας στατιστικά μια μεγάλη συλλογή κειμένων**
- Στάδια ενός αυτόματου τρόπου κατασκευής θησαυρών
 - **1/ Κατασκευή λεξιλογίου**
 - **2/ Υπολογισμός ομοιότητας μεταξύ όρων**
 - **3/ Οργάνωση (συνήθως ιεραρχική) του λεξιλογίου**



Αυτόματη Κατασκευή Θησαυρών από Κείμενα (II)

1/ Κατασκευή Λεξιλογίου

- Απόφαση: Ποιος θέλουμε να είναι ο βαθμός εξιδίκευσης (desired specificity)
 - if high then emphasis will be given on identifying precise phrases
- Οι όροι (terms) μπορούν να επιλεγούν από τους *τίτλους*, τις *περιλήψεις* (abstracts), ή ακόμα και από το *πλήρες κείμενο* (full text)
- Normalization: stemming, stoplists
- Criteria for selecting a term:
 - frequency of occurrence (**divide words to 3 categories: low, medium, high, select terms with medium frequency**)
 - discrimination value ~ idf
- Κατασκευή φράσεων (phrase construction) αν κάτι τέτοιο είναι επιθυμητό (θυμηθείτε coordination level)

2/ Υπολογισμός Ομοιότητας μεταξύ όρων

- Παραδείγματα μετρικών: Cosine, Dice



Αυτόματη Κατασκευή Θησαυρών από Κείμενα (III)

3/ Οργάνωση (συνήθως ιεραρχική) του λεξιλογίου

- Οποιοσδήποτε αλγόριθμος clustering μπορεί να χρησιμοποιηθεί για αυτό το βήμα

Ένας αλγόριθμος για ιεραρχική οργάνωση ενός λεξιλογίου:

1/ Identify a set of frequency ranges

2/ Group the vocabulary terms into different classes based on their frequencies and the ranges selected in Step 1. There will be one term class for each frequency range

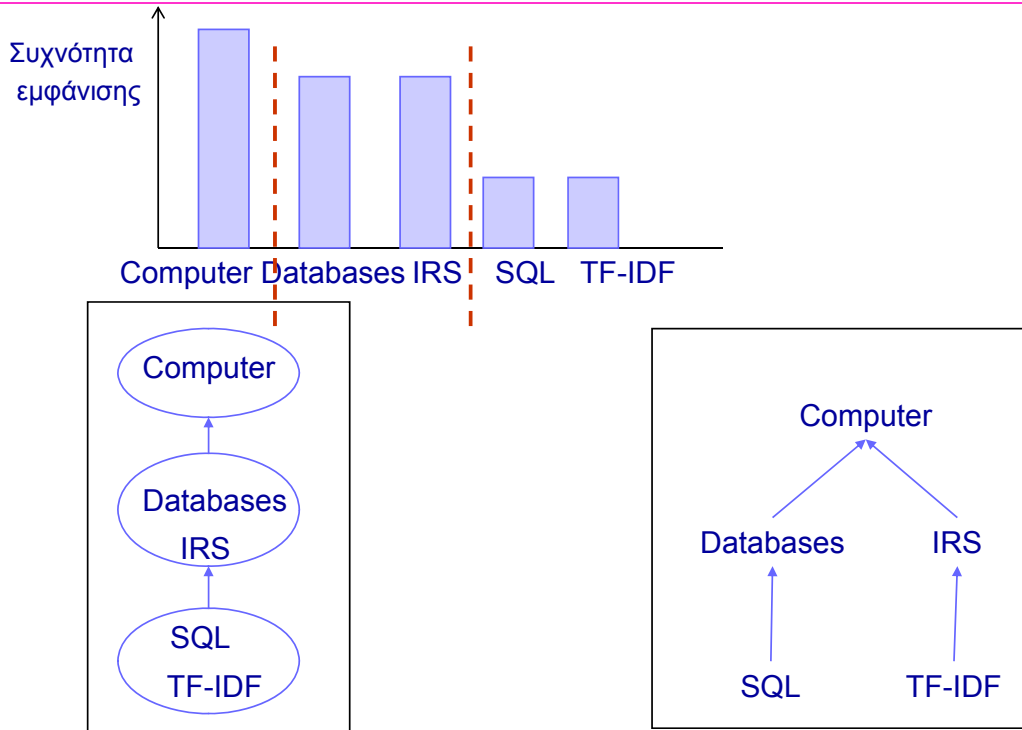
3/ The highest frequency class is assigned level 0, the next level 1, and so on

4/ Parent-child links: The parent(s) of a term at level i is the most similar term in level $i-1$ (a term is allowed to have multiple parents)

5/ Continue until reaching level 1



Παράδειγμα με 3 κλάσεις συχνοτήτων



CS463 - Information Retrieval Systems

Yannis Tzitzikas, U. of Crete

80



Case: grOOGLE'2007

- (1) Compute the minimum and maximum frequency of the words in the lexicon (denoted by df_{mn} and df_{mx} respectively).
- (2) Partition the interval $[df_{mn}, df_{mx}]$ into L successive intervals (where L is administrator-provided), i.e. $[df_{mn}, df_1], \dots, [df_{L-1}, df_{mx}]$. We will refer to them with lev_1, \dots, lev_L respectively.
- (3) Ignore the intervals corresponding to low frequencies, specifically keep only the M intervals with the highest frequencies (M is administrator-provided and it should be $M < L$), i.e. keep only $lev_{L-M+1}, \dots, lev_L$.
- (4) Assign to each of these M intervals those words whose frequency falls to that interval.
- (5) For each word w_i of level z (where $z \leq L - 1$) connect it with the most "correlated" word of the level $z + 1$ (that word will be the "parent" of w_i).

Regarding step (5), the correlation c_{ij} between two words w_i and w_j is computed using the formula:

$$c_{ij} = \sum_{d_k \in D} tf_{ik} \times tf_{jk} \quad (1)$$

where tf_{ik} is the frequency of term i in document k .

CS463 - Information Retrieval Systems

Yannis Tzitzikas, U. of Crete

81



(cont)

As an example, Table 11 describes the partitioning obtained assuming $L = 20$ (for each level the table shows the number of words that belong to that level). To construct the taxonomy we have considered only the last 5 groups (empty groups, like level 19, are considered as non existant). So the taxonomy includes 35 words in total. After creating the connections between words we realized that each word has an average of 1.4 child nodes.

	Low frequency										High frequency										
Level	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	
Num. Of Words	217	142	1103	523	292	199	128	83	83	53	52	25	18	18	14	14	7	8	2	0	4

The reason for partitioning words into groups (according to their frequency) is for avoiding computing the correlation matrix between all pairs of words (which would be formidably expensive⁷). In addition, ignoring those words that occur rarely further improves efficiency (as more than 95% of the vocabulary has a very small document frequency) and does not harm the quality of the result as these words do not describe the main concepts of the document corpus, and we have not anyway adequate statistical information to connect them right in a hierarchy.



(cont)

Resulting taxonomy:

- <1> http
- <2> system
- <3> us
- <3.1> new
- <3.1.1> url
- <3.1.1.1> map
- <3.1.1.2> network
- <3.1.1.2.1> commun
- <3.1.1.2.2> gener
- <3.1.1.2.3> site
- <3.1.1.2.4> scienc
- <3.1.1.2.5> web
- <3.1.1.3> page
- <3.1.1.3.1> link
- <3.1.1.3.2> applic
- <3.1.1.4> document
- <3.1.1.5> access
- <3.1.1.6> univers
- <3.1.1.7> data
- <3.1.1.8> process
- <3.1.2> time
- <3.1.3> base
- <4> inform

As you can see this taxonomy is **not very good/useful**

Possible improvements:

- **Better vocabulary construction**
 - The terms with high frequency are not very informative as you can see (e.g.. http, system, url, ...). Therefore we should try the **middle** levels.
 - Furthermore if we had selected words that appear only in titles/abstracts then we would avoid words like: http, url, ..
 - The user at run-time could even specify how specific/general the taxonomy should be (his/her choice would determine the visible part of the taxonomy)
- **Other improvements**
 - It's better to show the original words (rather than stems)
 - Use phrases instead of single words as terms

	Low frequency										High frequency										
Level	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	
Of Words	217	142	1103	523	292	199	128	83	83	53	52	25	18	18	14	14	7	8	2	0	4



Αυτόματη Κατασκευή Ιεραρχιών

[M. Sanderson and W. B. Croft. Deriving concept hierarchies from text. In SIGIR'1999]

- Given two terms x and y from a document collection, we say that x *subsumes* y , and we write $x \rightarrow y$ if: $P(x|y) > 0.8$ and $P(y|x) < 1$

$P(x|y)$ is the probability that term x occurs in a document, given that term y does

- This technique leads to creation of a hierarchy of terms, where
 - General terms appear as top-level categories
 - More specific terms appear as lower-level categories
- Pros
 - Simple and effective
- Cons
 - Requires n^2 computations of conditional probabilities, where n is the number of terms in the collection
 - Requires the terms to have a unique meaning
 - However If we use this technique only on query results and by using only terms that appear more frequently in the query results than in the whole collection. then this lessens the problem of ambiguity and reduces the number of terms that form the subsumption hierarchy.



Αυτόματη Κατασκευή Πολυεδρικών Ιεραρχιών

[Automatic Construction of Multifaceted Browsing Interfaces, W. Dakka, P. Ipeirotis, K. Wood, CICK'05]

- It describes an approach for constructing multifaceted hierarchies.
- Includes methods for selecting the best parts of the generated hierarchies when it is not possible to fit all the categories on screen
- Experiments with real-life data sets indicate that automatic construction of multifaceted interfaces is feasible, and generates high-quality hierarchies



A Data Mining approach (to organize the set of terms hierarchically)

- Let $I = \{i_1, \dots, i_m\}$ be a set of items
- Let D be a set of transactions where each transaction is a subset of I
- An association rule is an implication of the form $X \rightarrow Y$ where X, Y are subsets of I and $X \cap Y = \emptyset$
- A rule $X \rightarrow Y$ holds in the transaction set D with
 - confidence c if $c\%$ of the transactions in D that contain X also contain Y
 - support s if $s\%$ of the transactions in D contain $X \cup Y$

Consider the case of an IR system. In that case

- The set I could be the set of all terms (the vocabulary)
- The set D could be the set of binary vectors of the documents
- A rule $X \rightarrow Y$ would be an implication between set of terms
 - If $|X|=|Y|=1$ then the implications are between single terms
 - If $|X|=|Y|=2$ then the implications are between pairs of terms
- So we could exploit **data mining** algorithms to get a taxonomy from an IR system