



HY463 - Συστήματα Ανάκτησης Πληροφοριών
Information Retrieval (IR) Systems

Προχωρημένες Λειτουργίες Επερώτησης
Advanced Query Operations

Κεφάλαιο 5



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



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)

Βελτίωση της ερώτησης με:

- α) αναπροσαρμογή βαρών
- β) επέκταση της ερώτησης



Τεχνικές Βελτίωσης της Αρχικής Επερώτησης

Κατηγορίες:

(α) τεχνικές που **απαιτούν είσοδο από τον χρήστη**

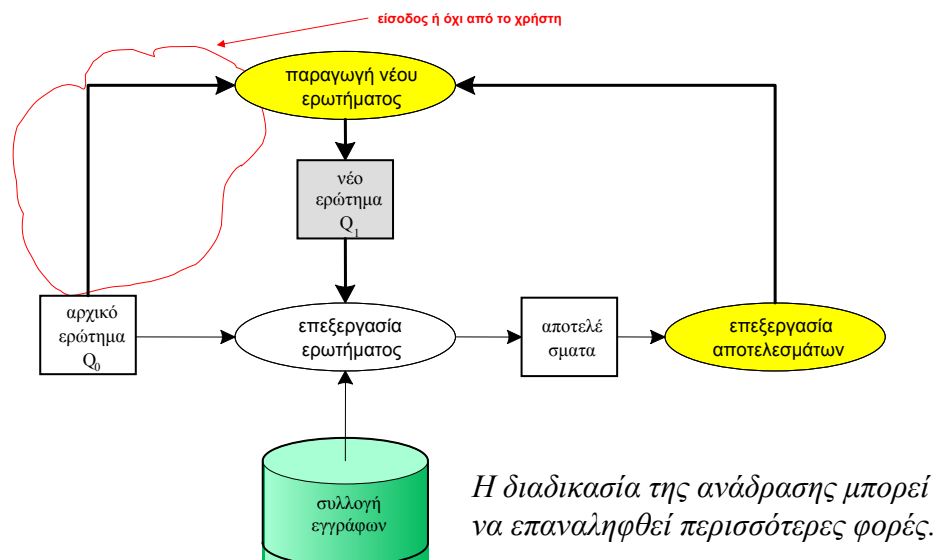
(β) τεχνικές που **δεν απαιτούν** είσοδο

(β1) που βασίζονται στα **κορυφαία έγγραφα** που ανακτήθηκαν

(β2) που βασίζονται σε **όλα τα έγγραφα** της συλλογής



Η Διαδικασία της Ανάδρασης





Ανάδραση Συνάφειας (Relevance Feedback): Η βασική ιδέα

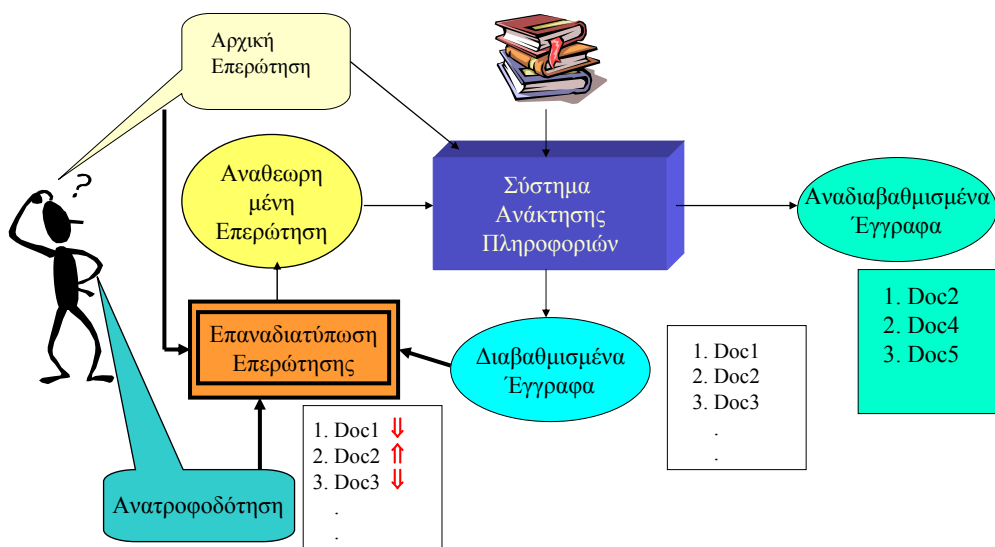
Με είσοδο από το χρήστη

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/

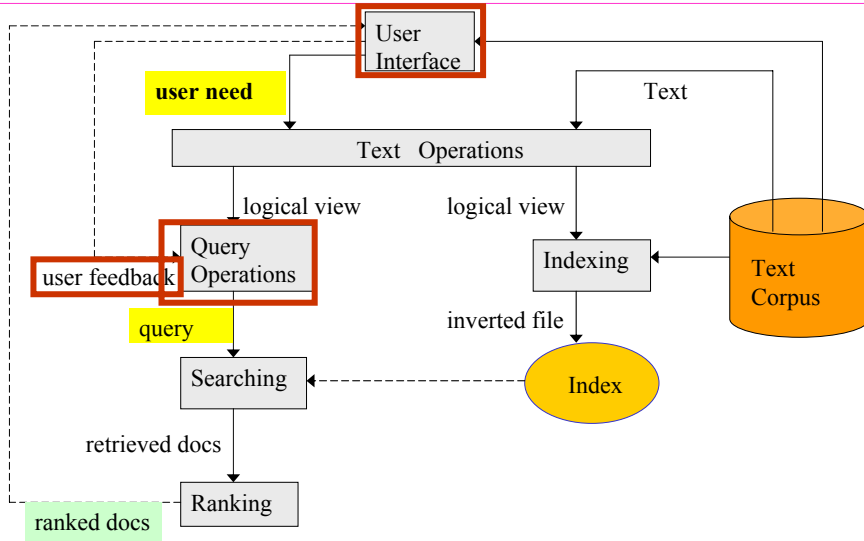


Ανάδραση Συνάφειας (Relevance Feedback): Η βασική ιδέα





Τμήματα της Αρχιτεκτονικής που Εμπλέκονται



Παράδειγμα ανατροφοδότησης συνάφειας σε σύστημα ανάκτησης εικόνων

q = bike








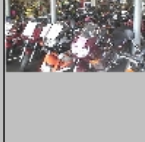






(<http://nayana.ece.ucsb.edu/imsearch/imsearch.html>)



Παράδειγμα ανατροφοδότησης συνάφειας
σε σύστημα ανάκτησης εικόνων

Answer("bike")=

Browse Search Prev Next Random

					
(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

CS463 - Information Retrieval Systems Yannis Tzitzikas, U. of Crete






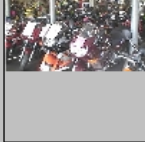
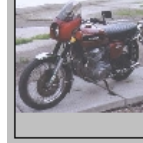
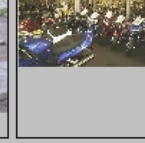




11



Παράδειγμα ανατροφοδότησης συνάφειας
σε σύστημα ανάκτησης εικόνων

Μαρκάρισμα των Συναφών (η Επιθυμητών) από τον Χρήστη

Browse Search Prev Next Random

					
(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













CS463 - Information Retrieval Systems Yannis Tzitzikas, U. of Crete

12



Παράδειγμα ανατροφοδότησης συνάφειας σε σύστημα ανάκτησης εικόνων

Απάντηση της αναδιατυπωμένης απάντησης =

Browse Search Prev Next Random					
					
(144538, 523493) 0.54182 0.231944 0.309876	(144538, 523835) 0.56219296 0.267304 0.295889	(144538, 523529) 0.584279 0.280881 0.303398	(144456, 253569) 0.64501 0.351395 0.293615	(144456, 253568) 0.650275 0.411745 0.23853	(144538, 523799) 0.66709197 0.358033 0.309059
					
(144473, 16249) 0.6721 0.393922 0.278178	(144456, 249634) 0.675018 0.4639 0.211118	(144456, 253693) 0.676901 0.47645 0.200451	(144473, 16328) 0.700339 0.309002 0.391337	(144483, 265264) 0.70170796 0.36176 0.339948	(144478, 512410) 0.70297 0.469111 0.233859

CS463 - Information Retrieval Systems

Yannis Tzitzikas, U. of Crete

13



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.

CS463 - Information Retrieval Systems

Yannis Tzitzikas, U. of Crete

14



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

Επέκταση επερώτησης (Query Expansion):

- Προσθήκη νέων όρων στην επερώτηση (π.χ. από γνωστά συναφή έγγραφα)

Υπάρχουν πολλοί αλγόριθμοι για επαναδιτύπωση επερώτησης



Ανάδραση στο Διανυσματικό Μοντέλο

Βέλτιστη ερώτηση: Η πιο γνωστή μέθοδος ανάδρασης στο Διανυσματικό μοντέλο είναι η μέθοδος του **Rocchio**.

A query vector, q_{opt} that

- maximizes similarity with relevant documents
- minimizes similarity with nonrelevant documents.

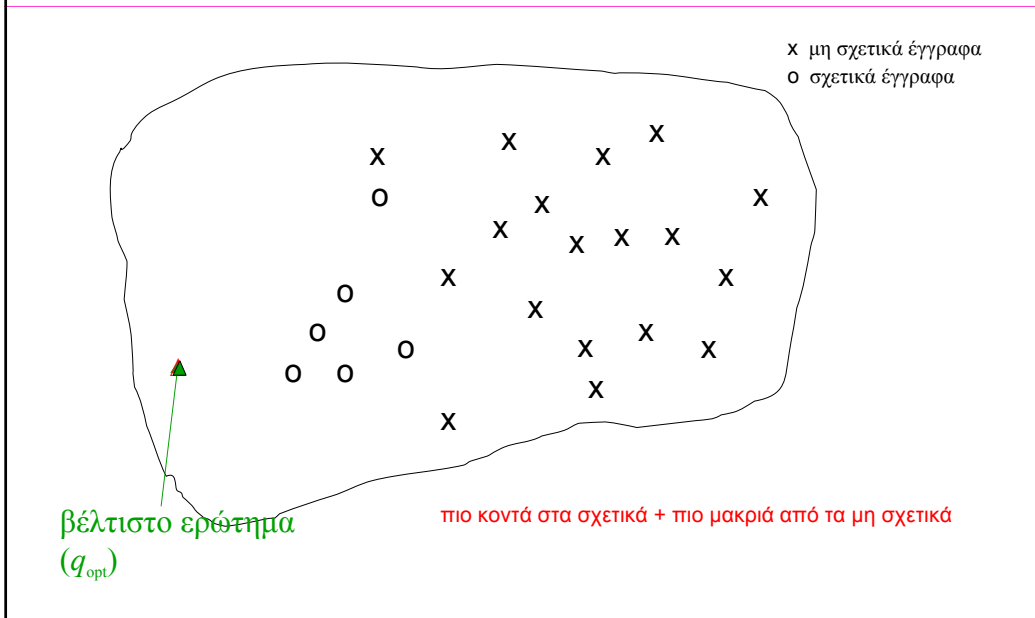
$$\bar{q}_{opt} = \arg \max_{\bar{q}} [\text{sim}(\bar{q}, C_r) - \text{sim}(\bar{q}, C_{nr})],$$

C : συλλογή εγγράφων, C_r : σύνολο σχετικών, C_{nr} : σύνολο μη σχετικών

Υπενθύμιση:
$$\text{sim}(d_1, d_2) = \frac{\vec{V}(d_1) \cdot \vec{V}(d_2)}{|\vec{V}(d_1)| |\vec{V}(d_2)|}$$



Ανάδραση στο Διανυσματικό Μοντέλο



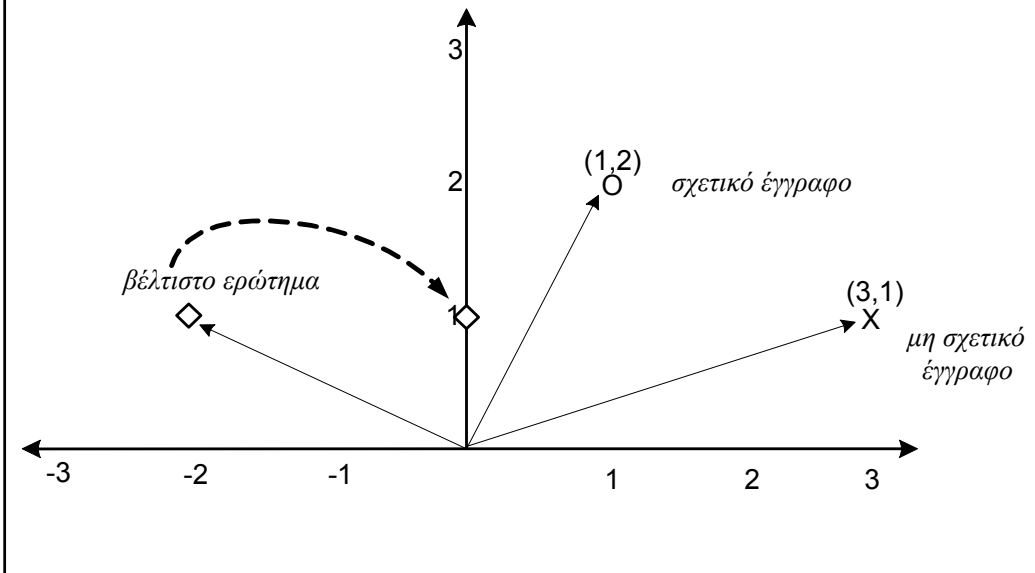
Ανάδραση στο Διανυσματικό Μοντέλο

Έστω ότι έχουμε ένα μόνο σχετικό έγγραφο (έστω r) και ένα μόνο μη σχετικό έγγραφο (έστω nr). Για να μπορέσουμε να διαχωρίσουμε το ένα από το άλλο, το διάνυσμα του βέλτιστου ερωτήματος Q_{opt} θα είναι (για συνημίτονο):

$$\text{vec}(Q_{opt}) = \text{vec}(r) - \text{vec}(nr)$$



Ανάδραση στο Διανυσματικό Μοντέλο



Ανάδραση στο Διανυσματικό Μοντέλο

Το βέλτιστο ερώτημα που πρέπει να διατυπώσουμε ώστε να διαχωριστούν τα σχετικά έγγραφα από τα μη σχετικά είναι (όταν χρησιμοποιείται ομοιότητα συνημιτόνου):

$$\vec{Q}_{opt} = \frac{1}{|R|} \sum_{\vec{d}_j \in R} \vec{d}_j - \frac{1}{N - |R|} \sum_{\vec{d}_j \notin R} \vec{d}_j$$

C : συλλογή εγγράφων, R : σύνολο σχετικών, $N-R$: σύνολο μη σχετικών

The optimal query is the vector difference between the centroids of the relevant and nonrelevant documents



Ανάδραση στο Διανυσματικό Μοντέλο

Ομαδοποίηση (Clustering)

Relevant documents have similarities among themselves ->
Ιδανική ερώτηση στο κέντρο (**centroid**) της συστάδα τους

Irrelevant documents have term-weight vectors which are
dissimilar from the ones for the relevant documents



Ανάδραση στο Διανυσματικό Μοντέλο

Πρόβλημα: το βέλτιστο ερώτημα δεν μπορεί να βρεθεί στην πράξη.

Γιατί?

Διότι δε γνωρίζουμε εκ των προτέρων το σύνολο των
σχετικών εγγράφων. Αν τα γνωρίζαμε ποιος ο λόγος να
εκτελέσουμε το ερώτημα?



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

Rocchio 1971 Algorithm (SMART)



(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

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

α : Tunable weight for initial query.

β : Tunable weight for relevant documents.

γ : Tunable weight for irrelevant documents.

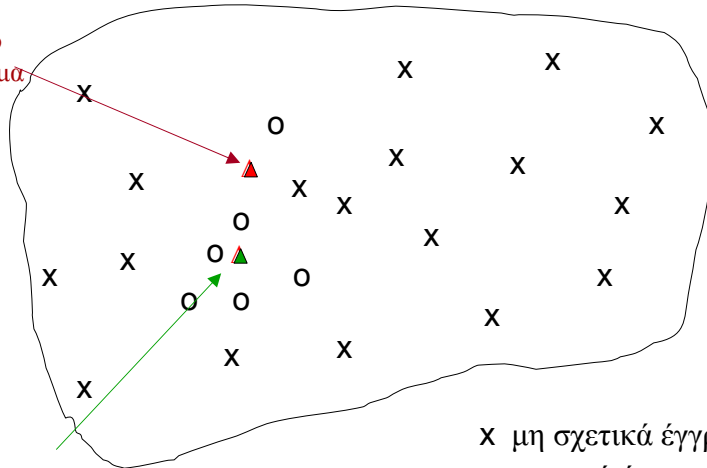
Usually $\gamma < \beta$ (the relevant docs are more important)

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



Ανάδραση στο Διανυσματικό Μοντέλο

αρχικό
ερώτημα



αναθεωρημένο ερώτημα

X μη σχετικά έγγραφα
O σχετικά έγγραφα



Ανάδραση στο Διανυσματικό Μοντέλο

query vect or = $\alpha \cdot$ αρχικό διάνυσμα ερωτήματος

+ $\beta \cdot$ θετική ανάδραση

- $\gamma \cdot$ αρνητική ανάδραση

Συνήθως, $\gamma < \beta$

αρχικό ερώτημα

0	4	0	8	0	0
---	---	---	---	---	---

 $\alpha = 1.0$

0	4	0	8	0	0
---	---	---	---	---	---

θετική ανάδραση

2	4	8	0	0	2
---	---	---	---	---	---

 $\beta = 0.5$

1	2	4	0	0	1
---	---	---	---	---	---

αρνητική ανάδραση

8	0	4	4	0	16
---	---	---	---	---	----

 $\gamma = 0.25$

2	0	1	1	0	4
---	---	---	---	---	---

νέο ερώτημα

-1	6	3	7	0	-3
----	---	---	---	---	----

(+)

(-)

Term weight can go negative

Negative term weights are ignored (set to 0)



Αναδιατύπωση επερώτησης στο Διανυσματικό Χώρο

Τρόποι αξιοποίησης της ανατροφοδότησης του χρήστη

- (I) **Rocchio** Method
- (II) **Ide** Method
- (III) **DeHi** Method



(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.



Αρνητική Ανάδραση

- Η αρνητική ανάδραση είναι μια μορφή ανάδρασης που χρησιμοποιεί πληροφορία από έγγραφα που έχουν χαρακτηριστεί ως **μη-σχετικά** από το χρήστη.
- Η αρνητική ανάδραση θεωρείται προβληματική:
 - Πότε ένας χρήστης θα πρέπει να χαρακτηρίζει ένα έγγραφο ως μη-σχετικό; Είναι πιο δύσκολο από το χαρακτηρισμό ως σχετικό.
 - Το να περιμένει κανείς από τους χρήστες να επιλέξουν έγγραφα ως μη-σχετικά στην αναζήτηση μπορεί να είναι δύσκολο στην υλοποίηση πρακτικά.



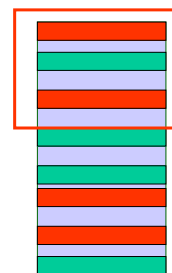
(III) IDE “Dec Hi” Method

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)

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

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

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

- Γενικά, τα πειραματικά δεδομένα δεν δίνουν καθαρό προβάδισμα σε κάποια τεχνική.
- Όλες οι τεχνικές βελτιώνουν την απόδοση (recall & precision)
- Συνήθως $\alpha = \beta = \gamma = 1$. Αν $\gamma = 0$, μόνο θετική ανάδραση



Relevance Feedback Example: Initial Query and Top 8 Results

Query: New space satellite applications

Note: want high recall

- + 1. 0.539, 08/13/91, NASA Hasn't Scrapped Imaging Spectrometer
- + 2. 0.533, 07/09/91, NASA Scratches Environment Gear From Satellite Plan
3. 0.528, 04/04/90, Science Panel Backs NASA Satellite Plan, But Urges Launches of Smaller Probes
4. 0.526, 09/09/91, A NASA Satellite Project Accomplishes Incredible Feat: Staying Within Budget
5. 0.525, 07/24/90, Scientist Who Exposed Global Warming Proposes Satellites for Climate Research
6. 0.524, 08/22/90, Report Provides Support for the Critics Of Using Big Satellites to Study Climate
7. 0.516, 04/13/87, Arianespace Receives Satellite Launch Pact From Telesat Canada
- + 8. 0.509, 12/02/87, Telecommunications Tale of Two Companies



Relevance Feedback Example: Expanded Query

- 2.074 new 15.106 space
 - 30.816 satellite 5.660 application
 - 5.991 nasa 5.196 eos
 - 4.196 launch 3.972 aster
 - 3.516 instrument 3.446 arianespace
 - 3.004 bundespost 2.806 ss
 - 2.790 rocket 2.053 scientist
 - 2.003 broadcast 1.172 earth
 - 0.836 oil 0.646 measure
- Νέα ερώτηση με 18 όρους και νέα βάρη



Top 8 Results After Relevance Feedback

- + 1. 0.513, 07/09/91, NASA Scratches Environment Gear From Satellite Plan
- + 2. 0.500, 08/13/91, NASA Hasn't Scrapped Imaging Spectrometer
3. 0.493, 08/07/89, When the Pentagon Launches a Secret Satellite, Space Sleuths Do Some Spy Work of Their Own
4. 0.493, 07/31/89, NASA Uses 'Warm' Superconductors For Fast Circuit
- + 5. 0.492, 12/02/87, Telecommunications Tale of Two Companies
6. 0.491, 07/09/91, Soviets May Adapt Parts of SS-20 Missile For Commercial Use
7. 0.490, 07/12/88, Gaping Gap: Pentagon Lags in Race To Match the Soviets In Rocket Launchers
8. 0.490, 06/14/90, Rescue of Satellite By Space Agency To Cost \$90 Million



Αξιολόγηση Αποτελεσματικότητας Τεχνικών Ανάδρασης Συνάφειας

Remarks

- Relevance feedback is most useful for increasing *recall* in situations where recall is important
 - Users can be expected to review results and to take time to iterate

Incremental Refinement

- Empirically, one round of relevance feedback is often very useful. Two rounds is sometimes marginally useful.



Αξιολόγηση Αποτελεσματικότητας Τεχνικών Ανάδρασης Συνάφειας

Αξιολόγηση:

Χρήση αρχικής q_0 και υπολογισμός precision and recall graph

Χρήση βελτιωμένης q_m και υπολογισμός precision and recall graph

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.

Συμπέρασμα: θα αγνοήσουμε τα έγγραφα για τα οποία έχουμε ήδη τη γνώμη του χρήστη;



Fair Process for Evaluating the Effectiveness of Relevance Feedback

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



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	
Ide (dec hi)							
expand by	Rank	3					
	Improvement	+49%	+44%	+92%	+32%	+79%	+59%
Rocchio (standard $\beta = .75, \alpha = .25$)							
expand by	Rank	2	39	8	14	17	16
	Precision	.2552	.1404	.2955	.1821	.5630	
expand by	Improvement	+75%	+19%	+156%	+33%	+68%	
most common terms	Rank	3	12	12	10	24	
	Precision	.2491	.1623	.2534	.1861	.5279	
	Improvement	+71%	+37%	+119%	+36%	+55%	
Probabilistic (adjusted revised derivatives)							

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



The screenshot shows the Rutgers INQUERY interface. At the top, there are buttons for 'Reset All', 'UNDO LAST RUN QUERY', 'Show Search Topic Text', 'Show Tutorial', and 'EXIT RU INQUERY'. Below these is a search input field with the text 'Enter (next) query term below and hit <RETURN>'. The current query is displayed as 'automobil* manufactur* car* defect* recal*'. A 'Run Query' button is located below the query. The results section shows a list of five items, each with a checkbox and a title: 1. GM Plans to Recall 62,000 1988-89 Cars With Quad 4 Engines; 2. GM, Ford Recall Vehicles to Repair Defective Parts; 3. Isuzu Motors, Honda Commence Car Recalls; 4. Ford and GM Recall Series Of Pickup Trucks, Coupes; 5. General Motors Corp. Recalls 196,000 Cars For Defective Brakes. Below the list, it states 'Total of 6747 documents retrieved' and 'Document # 1 of 6747'. The selected document is displayed in a large text area, showing a news article about GM's recall of 62,000 1988-89 cars with Quad 4 engines.



Πειραματικά Αποτελέσματα

- **Opaque (black box)**
 - Ο χρήστης δεν μπορεί να δει τη διαδικασία ανάδρασης.
- **Transparent**
 - Ο χρήστης μπορεί να δει τους όρους που δημιουργήθηκαν από την ανάδραση αλλά δεν μπορεί να μεταβάλει το ερώτημα.
- **Penetrable**
 - Ο χρήστης μπορεί να μεταβάλει το ερώτημα.



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

Criterion: $p@30$ (precision at 30 documents)

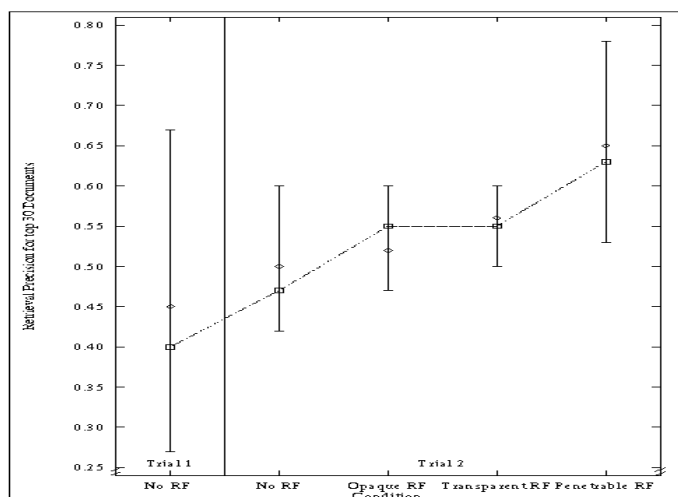
Compare:

- $p@30$ 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





Σιωπηρές Υποθέσεις της Ανάδρασης Συνάφειας

A1: User has sufficient knowledge for initial query.

A2: Relevance prototypes are “well-behaved”.

- Term distribution in relevant documents will be similar
- Term distribution in non-relevant documents will be different from those in relevant documents
 - Either: All relevant documents are tightly clustered around a single prototype.
 - Or: There are different prototypes, but they have significant vocabulary overlap.
 - Similarities between relevant and irrelevant documents are small



Violation of A1

- User does not have sufficient initial knowledge.
- Examples:
 - Misspellings (Brittany Speers).
 - Cross-language information retrieval (higado).
 - Mismatch of searcher’s vocabulary vs. collection vocabulary
 - Cosmonaut/astronaut – laptop/notebook



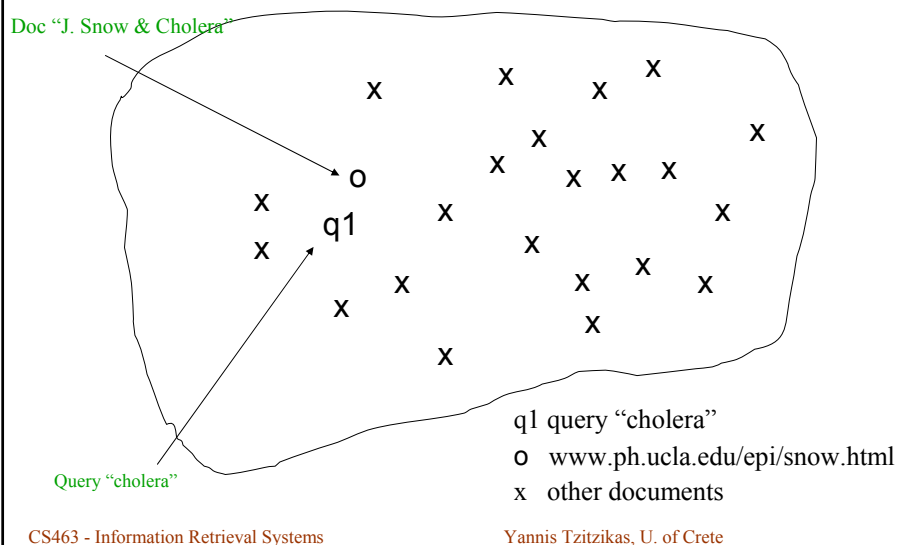
Violation of A2

However, there are several relevance prototypes.

- Examples:
 - Burma/Myanmar
 - Contradictory government policies
 - Pop stars that worked at Burger King (query whose answer set is inherently disjunctive)
 - Often: instances of a general concept
- Good editorial content can address problem
 - Report on contradictory government policies



Aside: Vector Space can be Counterintuitive.





High-dimensional Vector Spaces

- The queries “cholera” and “john snow” are far from each other in vector space.
- How can the document “John Snow and Cholera” be close to both of them?
- Our intuitions for 2- and 3-dimensional space don't work in >10,000 dimensions.
- 3 dimensions: If a document is close to many queries, then some of these queries must be close to each other.
- Doesn't hold for a high-dimensional space.



Γιατί η ανάδραση συνάφειας δεν χρησιμοποιείται ευρέως;

- Long queries are inefficient for typical IR engine.
 - Long response times for user.
 - High cost for retrieval system.
 - Partial solution:
 - Only reweight certain prominent terms
 - Perhaps top 20 by term frequency
- Users are often reluctant to provide explicit feedback
- It's often harder to understand why a particular document was retrieved after applying relevance feedback





Ανάδραση Συνάφειας στον Παγκόσμιο Ιστό

The screenshot shows a Google search interface. The search term is 'Information Retrieval'. The results show a link to 'INFORMATION RETRIEVAL' with a description: 'An online book by CJ Rijsbergen, University of Glasgow. www.dcs.gla.ac.uk/Keith/Preface.html - 7k'. Below the link, there is a red circle around the text 'Παρόμοιες σελίδες' (Similar/related pages).

- 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.



Relevance Feedback on the Web

[in 2003: now less major search engines, but same general story]

- Some search engines offer a similar/related pages feature (this is a trivial form of relevance feedback)
 - Google (link-based)
 - Altavista
 - Stanford WebBase
- But some don't because it's hard to explain to average user:
 - Alltheweb
 - msn
 - Yahoo
- Excite initially had true relevance feedback, but abandoned it due to lack of use.



Excite Relevance Feedback

Spink et al. 2000

- Only about 4% of query sessions from a user used relevance feedback option
 - Expressed as “More like this” link next to each result
- But about 70% of users only looked at first page of results and didn't pursue things further
 - So 4% is about 1/8 of people extending search
- Relevance feedback improved results about 2/3 of the time



Other Uses of Relevance Feedback

- Following a changing information need (name of car models change over time)
- Maintaining an information filter (e.g., for a news feed)
- Active learning
 - [Deciding which examples it is most useful to know the class of to reduce annotation costs]



Δυναμική Αναζήτηση

- Κατά ένα μεγάλο μέρος στη δουλειά που γίνεται για το RF υπάρχει η παραδοχή ότι η πληροφορία που αναζητεί ο χρήστης δε μεταβάλλεται κατά τη διάρκεια της αναζήτησης.
- Αν αυτό δεν συμβαίνει, τότε τα έγγραφα που χαρακτηρίστηκαν σχετικά στην αρχή της αναζήτησης μπορεί να μην είναι καλά παραδείγματα για το τι θεωρεί ο χρήστης σχετικό μετά.
- Με χρήση ageing component μπορεί να μειώνεται το βάρος ενός όρου όσο περνάει ο χρόνος ώστε να έχει μικρότερη επίδραση στην εύρεση εγγράφων.



Ψευτοανάδραση Συνάφειας

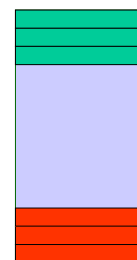
Ανάδραση χωρίς είσοδο από το χρήστη



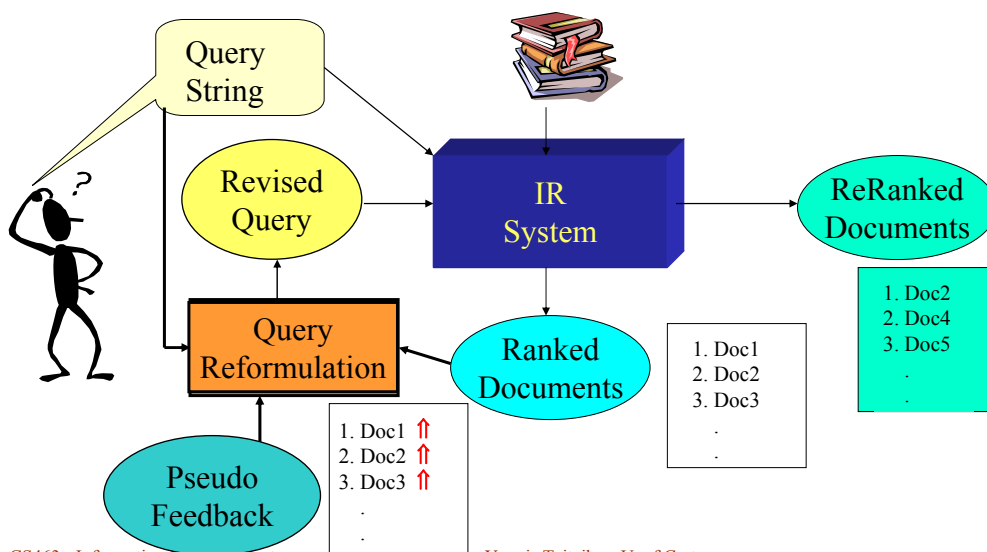
Ψευδοανάδραση Συνάφειας Pseudo Relevance Feedback

- 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





Σχετικά με την Ψευδοανάδραση

- Γνωστή και ως **blind** ή **ad-hoc** RF, χρησιμοποιεί τεχνικές RF για να **βελτιώσει αυτόματα** την κατάταξη πριν παρουσιαστούν τα έγγραφα στο χρήστη.
- Η pseudo RF τεχνική λειτουργεί ικανοποιητικά για «καλά» αρχικά ερωτήματα αλλά είναι αναποτελεσματική για «κακά» αρχικά ερωτήματα.



Αξιολόγηση Ψευδοανάδρασης

- Βρέθηκε να βελτιώνει την απόδοση στο διαγωνισμό του TREC (ad-hoc retrieval task)
- Δίνει τη δυνατότητα να βρεθεί ένα καλύτερο ερώτημα **χωρίς** να απαιτείται από το χρήστη να διατυπώσει ένα νέο ερώτημα από την αρχή.
- Εύκολη υλοποίηση και εφαρμογή στα πιο δημοφιλή μοντέλα ανάκτησης (boolean, vector, probabilistic)
- 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)



Ανάδραση στο Πιθανοκρατικό Μοντέλο

Η συνάρτηση ομοιότητας του πιθανοκρατικού μοντέλου είναι:

$$S_{prob}(q, d) = \sum_i \log \frac{p_i \cdot (1 - r_i)}{r_i \cdot (1 - p_i)}$$

Όπου η άθροιση αφορά στους όρους που βρίσκονται **και στο ερώτημα και στο έγγραφο**.

p_i πιθανότητα ο όρος t_i να εμφανίζεται σε σχετικό έγγραφο

r_i πιθανότητα ο όρος t_i να εμφανίζεται σε μη σχετικό έγγραφο



Ανάδραση στο Πιθανοκρατικό Μοντέλο

Αρχικά θέτουμε τιμές στις πιθανότητες :

$$p_i = P(x_i | R) = c$$

$$r_i = P(x_i | \bar{R}) = n_i / N$$

όπου:

c είναι μία τυχαία σταθερά (π.χ., 0.5)

n_i είναι το πλήθος των εγγράφων που περιέχουν τον i -οστό όρο

N πλήθος εγγράφων συλλογής



Ανάδραση στο Πιθανοκρατικό Μοντέλο

- Για τη βελτίωση της ποιότητας των αποτελεσμάτων οι πρώτες εφαρμογές του Πιθανοκρατικού μοντέλου χρειάζονται την παρέμβαση του χρήστη για την αναπροσαρμογή των τιμών (πχ Rocchio).

Εναλλακτικά μπορεί να χρησιμοποιηθεί και αυτοματοποιημένος τρόπος (ψευτο-ανάδραση).

Αρχικά εκτελείται το ερώτημα με τις αρχικές εκτιμήσεις. Επιλέγονται τα k καλύτερα έγγραφα.

Έστω k_i ο αριθμός των εγγράφων που περιέχουν τον i -οστό όρο. Θέτουμε:

$$p_i = P(x_i | R) = k_i / k \quad + 0.5 \text{ adjustment}$$

$$r_i = P(x_i | \bar{R}) = (n_i - k_i) / (N - k)$$



Ανάδραση στο Πιθανοκρατικό Μοντέλο

Για μικρές τιμές του k και k_i

$$p_i = P(x_i | R) = (k_i + 0.5) / (k + 1) \quad + 0.5 \text{ adjustment}$$

$$r_i = P(x_i | \bar{R}) = [(n_i - k_i) + 0.5] / [(N - k) + 1]$$

Άλλο πιθανό + k_i/N



Χρήσιμοι Σύνδεσμοι

Σύστημα ανάκτησης εικόνων με δυνατότητα ανάδρασης

<http://amazon.ece.utexas.edu/~qasim/cires.htm>

Survey of relevance feedback over VSM, Probabilistic and Logic Model

http://www.dcs.qmul.ac.uk/~mounia/CV/Papers/ker_ruthven_lalmas.pdf



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 Θησαυρώ



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



Διαδραστική Επαύξηση Ερωτήματος

- Οι μέθοδοι για τροποποίηση ερωτήματος που έχουν περιγραφεί ως τώρα επιλέγουν όρους από έγγραφα και προσθέτουν κάποιους στο ερώτημα.
- Ένας εναλλακτικός τρόπος είναι οι χρήστες να επιλέγουν τους όρους που θα προστεθούν στο ερώτημα (**IQE - interactive query expansion**).
 - Το σύστημα επιλέγει τους πιο σχετικούς όρους
- Αναζητήσεις σε περιπτώσεις που έχει ανακτηθεί λίγη σχετική πληροφορία ωφελούνται από IQE.
- Είναι πιο πιθανό οι χρήστες να χρησιμοποιήσουν IQE σε μια πολύπλοκη ή δύσκολη αναζήτηση.



Επέκταση Επερώτησης (Query Expansion)

Πως θα βρούμε τους πιο σχετικούς όρους

- 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

1. 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
2. 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' ό,τι 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](#) [Jaquar Car](#) [Black Jaguar](#) [Jaguar Xk8](#)
[Wild Jaguars](#) [Jaquare](#) [Jaguar Accessories](#) [Jaguar Automobile](#)

Also: see altavista, teoma

Ο χρήστης βλέπει τις κορυφαίες λέξεις
(και όχι έγγραφα)


[Web](#) | [Images](#) | [Video](#) | [Local](#) | [Shopping](#) | [more >](#)

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** merchandise Book at SHOP.COM.
[www.SHOP.com](#)
- Jaguar Natural Spray on Cataloglink**
 Find **Jaguar** natural spray on Cataloglink

CS463 - Information Retrieval Systems Yannis Tzitzikas, U. of Crete 71

MyStuff | Settings



[Web](#) | [Images](#) | [Video](#) | [More >](#)

[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



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)
 The **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'
 (jäg-wär, jäg-yü-är)^(*) - 1 definition
 The American Heritage® Diction...
 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...

CS463 - Information Retrieval Systems Yannis Tzitzikas, U. of Crete



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.

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



Αυτόματη Τοπική (Επιτόπια) Ανάλυση

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 $t_1 \dots t_n$ οι όροι που εμφανίζονται σε αυτά.

	t_1	t_2	t_3	t_n
t_1	c_{11}	c_{12}	c_{13}	c_{1n}
t_2	c_{21}				
t_3	c_{31}				
.	.				
.	.				
t_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:
Normalized score is 1 if two terms have the same frequency in all documents in D.

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

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)



(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

	t_1	t_2	t_3	t_n
t_1	c_{11}	c_{12}	c_{13}	c_{1n}
t_2	c_{21}				
t_3	c_{31}				
.	.				
.	.				
t_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.



Query Expansion with Scalar Clusters

	t_1	t_2	t_3	t_n
t_1	c_{11}	c_{12}	c_{13}	c_{1n}
t_2	c_{21}				
t_3	c_{31}				
.	.				
.	.				
t_n	c_{n1}				

For each query term t_q : consider the terms that have similar correlation values with it

Inner product of the related vectors



Αυτόματη Καθολική Ανάλυση (Automatic Global Analysis)



Αυτόματη Καθολική Ανάλυση Automatic Global Analysis

1. 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
2. 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 στο 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):

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

Num of terms in the collection

$$\text{itf}_j = \frac{\text{Num of terms in the collection}}{\text{Num of distinct terms in } d_j}$$

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



Query Expansion Based on a Similarity Thesaurus

Υπενθύμιση

idf_i = inverse document frequency of term i := $\log(N/df_i)$

$$tf_{ij} = \text{freq}_{ij} / \max_k \{ \text{freq}_{kj} \}$$

Όπου

- freq_{ij} = πλήθος εμφανίσεων του όρου i στο έγγραφο j
- $\max_k \{ \text{freq}_{kj} \}$ το μεγαλύτερο πλήθος εμφανίσεων ενός όρου στο έγγραφο j

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



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 δίνει 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.

CS463 - Information Retrieval Systems

Yannis Tzitzikas, U. of Crete

The screenshot shows the Ask.com search interface. At the top, there are links for 'MyStuff' and 'Settings'. The Ask.com logo is prominent. Below it, there are navigation links for 'Web', 'Images', 'Video', and 'More'. A search bar contains the text 'jaguar' and a search button. Below the search bar, there are sections for 'Advanced Search', 'Narrow' (with sub-sections like Jaguar Cars, Black Jaguar, Cat Jaguar, Jaguar Big Cats, Jaguars Habitat, What Do Jaguars Eat, Panthera Onca, Where Do Jaguars Live), 'Expand' (with sub-sections like Cheetah, Ferrari), and 'Related Names' (with sub-sections like Ford, Wolf). Each sub-section has a 'More >' link.



Διάκριση Θησαυρών

- Γλωσσικοί Θησαυροί
 - Πχ Roget's thesaurus. Designed to assist the writer in creatively selecting vocabulary
- Θησαυροί κατάλληλοι για Information Retrieval
 - for coordinating the basic processes of indexing and retrieval
 - designed for specific subject areas and are therefore domain dependent
 - Examples
 - INSPEC

CS463 - Information Retrieval Systems

Yannis Tzitzikas, U. of Crete

90



INSPEC thesaurus (for IR)

- Domain: physics, electrical engineering, electronics, computers
- Types of relationships between two terms
 - **UF: Used For (converse: USE)** // π.χ. USE X σημαίνει ότι ο X είναι ο δόκιμος όρος
 - **BT: Broader Term (converse NT)**
 - **TT: Top Node, i.e. root of the hierarchy)**
 - **RT: Related Term**
- Example:
 - computer-aided instruction
 - see also education
 - UF teaching machines (UF: Used For, converse: USE)
 - BT educational computing (BT: Broader Term)
 - TT computer applications (TT: Top Node, i.e. root of the hierarchy)
 - RT education, teaching (RT: Related Term)



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

- A more detailed database of semantic relationships between English words. Developed by 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

- **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)



AAT (Art and Architecture Thesaurus)

- Controlled vocabulary for describing and retrieving information: 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 contain 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
 - 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

CS463 - Information Retrieval Systems

Yannis Tzitzikas, U. of Crete

95



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
 - “interest rate” → “interest rate fascinate evaluate”

CS463 - Information Retrieval Systems

Yannis Tzitzikas, U. of Crete

96



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



[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 υποστηρίζει ο θησαυρός
- **Ιδέα: Μπορούμε να ανακαλύψουμε σημασιολογικές σχέσεις μεταξύ λέξεων αναλύοντας στατιστικά μια μεγάλη συλλογή κειμένων**
- Στάδια
 - 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 \sim 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

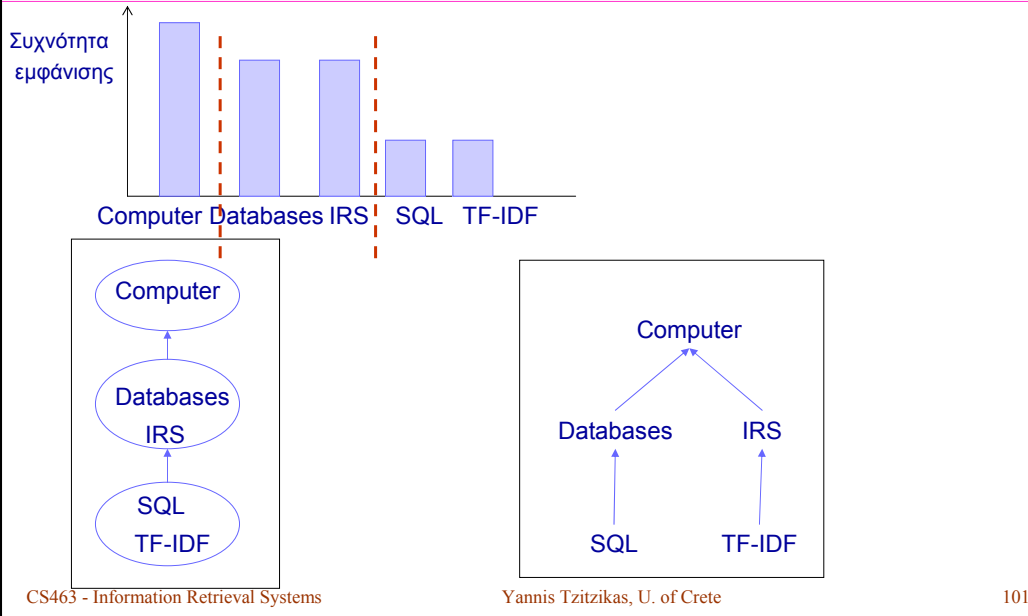
- Οποιοσδήποτε αλγόριθμος clustering μπορεί να χρησιμοποιηθεί

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 συχνότητας



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 .



(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										
Of Words	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	
	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



Περίληψη

- The complete landscape
 - Global methods
 - Query expansion
 - Thesauri
 - Automatic thesaurus generation
 - Local methods
 - Relevance feedback
 - Pseudo relevance feedback



Relevance Feedback -- Summary

- Relevance feedback has been shown to be very effective at improving relevance of results.
 - Requires enough judged documents, otherwise it's unstable (≥ 5 recommended)
 - Requires queries for which the set of relevant documents is medium to large
- Full relevance feedback is painful for the user.
- Full relevance feedback is not very efficient in most IR systems.
- Other types of interactive retrieval may improve relevance by as much with less work.