



Θα μιλήσουμε για  
**ΑΞΙΟΛΟΓΗΣΗ ΤΗΣ ΑΝΑΚΤΗΣΗΣ ΠΛΗΡΟΦΟΡΙΑΣ**

Βασισμένες στις διαφάνειες του καθ. **Γιάννη Τζίτζικα** (Παν. Κρήτης)

<http://www.ics.forth.gr/~tzitzik/>

Κεφάλαιο 3 του βιβλίου



## Διάρθρωση Διάλεξης

- *Τι εξυπηρετεί η αξιολόγηση;*
  - αξιολόγηση αποτελεσματικότητας
- *Δυσκολίες της αξιολόγησης*
- *Αξιολόγηση μέσω Χειρονακτικά Μαρκιαρισμένων Συλλογών*
- *Μέτρα αξιολόγησης αποτελεσματικότητας*
  - *Ανάκληση & Ακρίβεια & (Recall & Precision)*
  - *Καμπύλες Ακρίβειας/Ανάκλησης*
    - *Σύγκριση Συστημάτων*
  - *Εναλλακτικά μέτρα*
    - *R-Precision (Precision Histograms)*
    - *F-Measure*
    - *E-Measure*
    - *Fallout, Expected Search Length*
  - *User-Oriented Measures*
- *Δοκιμασίες Αποτελεσματικότητας-Συλλογές Αναφοράς (TREC)*



## Είδη Αξιολόγησης

### 1. Functional Analysis (ανάλυση των λειτουργικών απαιτήσεων)

Κάθε σύστημα λογισμικού πρέπει να παρέχει τη λειτουργικότητα για την οποία σχεδιάστηκε

Specified system functionalities are tested one by one +  
Error analysis phase

Για τα συστήματα που ικανοποιούν τις λειτουργικές απαιτήσεις, αξιολόγηση της απόδοσης τους



## Είδη Αξιολόγησης

### 2. Performance Evaluation (Αξιολόγηση της Απόδοσης του Συστήματος)

Συνήθως:

Χώρος

Χρόνος (απόκρισης)

Quantifiable

Θέματα στην περίπτωση της ΑΠ: ευρετήρια, ΛΣ, επικοινωνία, κλπ

*Αξιολογούμε την αποδοτικότητα (efficiency) - επιδόσεις του συστήματος*



## Είδη Αξιολόγησης

### 3. Retrieval Performance Evaluation (Αξιολόγηση της Απόδοσης της Ανάκτησης)

Πόσο ακριβές/σωστό είναι το σύνολο απάντησης

*Αξιολογούμε την αποτελεσματικότητα (effectiveness)*



## Why Evaluate?

### 1. Ποια επιλογή είναι καλύτερη:

Υπάρχουν πολλά μοντέλα υπολογισμού του βαθμού συνάφειας, πολλοί αλγόριθμοι και ακόμα περισσότερα συστήματα. Ποιο είναι το καλύτερο;

Ποιος είναι ο καλύτερος τρόπος για:

Επιλογή των όρων του ευρετηρίου (stopword removal, stemming...)

Προσδιορισμό των βαρών των όρων (Term weighting) (TF, TF-IDF,...)

Υπολογισμό του βαθμού συνάφειας (dot-product, cosine, ...) βάσει του οποίου θα γίνει η κατάταξη των εγγράφων;

### 2. Καθορισμός/εκτίμηση κάποιων παραμέτρων

Καθορισμός καλών τιμών για παραμέτρους

Πόσα έγγραφα της απόκρισης ενός συστήματος πρέπει να εξετάσει ο χρήστης προκειμένου να βρει μερικά/όλα τα συναφή έγγραφα;



## Why Evaluate?

3. Σύγκριση συστημάτων

4. Αξιολόγηση μιας νέας αλλαγής



## Παραδείγματα

Πότε είναι καλή (αποτελεσματική) για παράδειγμα:

- Μηχανή αναζήτησης
- Αναζήτηση αρχείων
- Αναζήτηση σε μια βιβλιοθήκη
- Σύστημα προτιμήσεων

Πως θα την αξιολογήσουμε;

- User studies
- Test collection
- Traces (infer from datasets -> click throughs, etc)



## Αξιολόγηση βάσει Χειρονακτικά Μαρκκαρισμένων Συλλογών (Human Labeled Corpora)

Τρόπος:

- 1) Επέλεξε ένα συγκεκριμένο σύνολο εγγράφων  $C$  (κατά προτίμηση του ίδιου γνωστικού πεδίου).
- 2) Διατύπωσε ένα σύνολο επερωτήσεων για αυτά  $Q$
- 3) Βρες έναν ή περισσότερους ειδικούς (experts) του γνωστικού πεδίου, και βάλε τους να μαρκάρουν τα συναφή έγγραφα για κάθε ερώτηση  
Συνήθως, οι κρίσεις τους είναι (Συναφές, Μη-Συναφές)  
Αρα **το αποτέλεσμα της διαδικασίας αυτής είναι ένα σύνολο από πλειάδες της μορφής:  $(c,q,Relevant)$  ή  $(c,q,Irrelevant)$ , όπου  $c \in C, q \in Q$ .**
- 4) Χρησιμοποίησε αυτή τη συλλογή για την αξιολόγηση της αποτελεσματικότητας ενός ΣΑΠ.  
Βάζουμε το ΣΑΠ να ευρετηριάσει τη συλλογή  $C$ , κατόπιν του στέλνουμε επερωτήσεις από το  $Q$  και αξιολογούμε τις αποκρίσεις του βάσει των κρίσεων που έχουν κάνει ήδη οι ειδικοί.

Δυσκολίες:

Η παραπάνω μέθοδος απαιτεί μεγάλη ανθρώπινη προσπάθεια για μεγάλες συλλογές εγγράφων/επερωτήσεων.



## Αξιολόγηση βάσει Χειρονακτικά Μαρκκαρισμένων Συλλογών (Human Labeled Corpora)

- 4) Χρησιμοποίησε αυτή τη συλλογή για την αξιολόγηση της αποτελεσματικότητας ενός ΣΑΠ.

Βάζουμε το ΣΑΠ να ευρετηριάσει τη συλλογή  $C$ , κατόπιν του στέλνουμε επερωτήσεις από το  $Q$  και αξιολογούμε τις αποκρίσεις του βάσει των κρίσεων που έχουν κάνει ήδη οι ειδικοί.

Πως;



## Test Collection

We need a test collection consisting of three things:

1. A **document collection**
2. A **test suite of information needs**, expressible as queries
3. A set of **relevance judgments**, standardly a binary assessment of either *relevant* or *nonrelevant* for each query-document pair.

### Relevance Judgment:

With respect to a user information need:

a document in the test collection is given a binary classification as either *relevant* or *nonrelevant*.

- o Referred to as the *gold standard* or *ground truth judgment* of relevance.
- o Non binary relevance (how much relevant)



## Test Collection

- collection and suite of information needs have to be of a **reasonable size**:

Why? need to average performance over fairly large test sets, as results are highly variable over different documents and information needs).

As a rule of thumb: 50 information needs a sufficient minimum.

- Relevance is assessed relative to an *information need*, not a query.

example, an information need might be:

*Information on whether drinking red wine is more effective at reducing your risk of heart attacks than white wine.*

might be translated into a query such as:

*wine AND red AND white AND heart AND attack AND effective*



## Test Collection

Many systems contain various *parameters* (weights) that can be adjusted to tune system performance.

- Wrong: report results on a test collection which were obtained by tuning these parameters to maximize performance on that collection (such tuning overstates the expected performance of the system)
- Correct: have one or more *development test collections*, and to tune the parameters on the development test collection. The tester then runs the system with those weights on the test collection and reports the results on that collection as an unbiased estimate of performance.



## Οι Δυσκολίες της Αξιολόγησης Αποτελεσματικότητας

- Η αποτελεσματικότητα εξαρτάται από τη *συνάφεια* των ανακτημένων εγγράφων
  - Δεν υπάρχει τυπικός ορισμός της συνάφειας
- Στην ουσία η συνάφεια δεν είναι δυαδική αλλά συνεχής
- Ακόμα και αν ήταν δυαδική, η κρίση της μπορεί να μην είναι εύκολη
- Από την πλευρά του χρήστη η συνάφεια είναι:
  - **υποκειμενική**: διαφορετική από χρήστη σε χρήστη
  - **περιστασιακή** (situational): σχετίζεται με τις τρέχουσες ανάγκες του χρήστη
  - **γνωστική** (cognitive): εξαρτάται από την αντίληψη/συμπεριφορά του χρήστη
  - **δυναμική**: μεταβάλλεται με το χρόνο (δεν είναι αναλλοίωτη)



## Τύποι Αξιολόγησης

### Retrieval Task

- Batch mode
- Interactive session (additional issues: interface, guidance, duration of the session, etc)

### Setting

- Real-life
- Laboratory (repeatability, scalability)

### Type of interface



## Measures of Effectiveness

### Retrieval Performance Evaluation (Αξιολόγηση της Απόδοσης της Ανάκτησης)

Given a retrieval strategy  $S$ , the evaluation measure quantifies the **similarity** between the set of documents retrieved by  $S$  and the **set of relevant documents** provided by the experts.

An estimation of the **goodness** of the retrieval strategy  $S$





## Μέτρα αξιολόγησης αποτελεσματικότητας

Θεωρούν μη διαβαθμισμένη απάντηση



## Μέτρα αξιολόγησης αποτελεσματικότητας: Ακρίβεια (Precision) και Ανάκληση(Recall)

- **Ακρίβεια (Precision):**

Διαισθητικά: Η ικανότητα ανάκτησης μόνο συναφών εγγράφων (Πόσα από τα έγγραφα στην απάντηση είναι συναφή - πόσα «σκουπίδια» παίρνω)

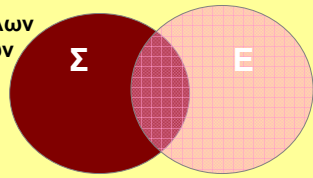
- **Ανάκληση (Recall):**

Διαισθητικά: Η ικανότητα εύρεσης όλων των συναφών εγγράφων της συλλογής (Πόσα από τα συναφή έγγραφα βρίσκονται στην απάντηση - πόσα συναφή έγγραφα χάνω)



## Ακρίβεια (Precision) και Ανάκληση(Recall)

Η συλλογή όλων των εγγράφων



**Σ:** Το «τέλειο» σύνολο, όλα τα συναφή (με το ερώτημα q) έγγραφα (π.χ. μας τα έχουν δώσει οι ειδικοί)

**Ε:** Το σύνολο των εγγράφων στην απάντηση (από το ΣΑΠ)

$$P(\text{recision}) = P(\text{relevant} | \text{retrieved})$$

$$\text{Ακρίβεια} = \frac{|E \cap \Sigma|}{|E|}$$

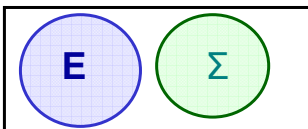
Not relevant	False negatives	True positives
	True negatives	False positive
	Not retrieved	Retrieved

$$R(\text{ecall}) = P(\text{retrieved} | \text{relevant})$$

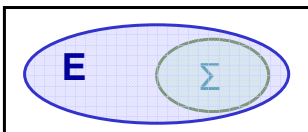
$$\text{Ανάκληση} = \frac{|E \cap \Sigma|}{|\Sigma|}$$



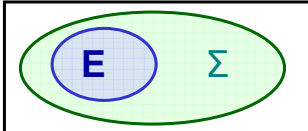
## Περιπτώσεις



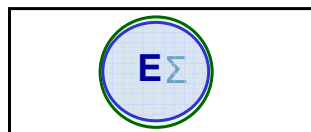
P=0%, R=0%  
(χειρότερη περίπτωση)



P=low, R=100% (η επίτευξη R=1 είναι ευκολότερη)



P=100%, R:low



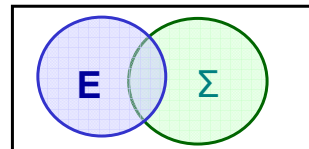
P=100%, R=100% (ιδανική περίπτωση)

Επειδή τα πιο πολλά είναι μη σχετικά - > όλα μη σχετικά (καλή accuracy)

Επιστροφή όλων (καλή ανάκληση)

Όσο αυξάνει το μέγεθος της απάντησης

Recall? Precision?





## What about accuracy?

Classification Problem

**Accuracy:** the fraction of its classification that are correct

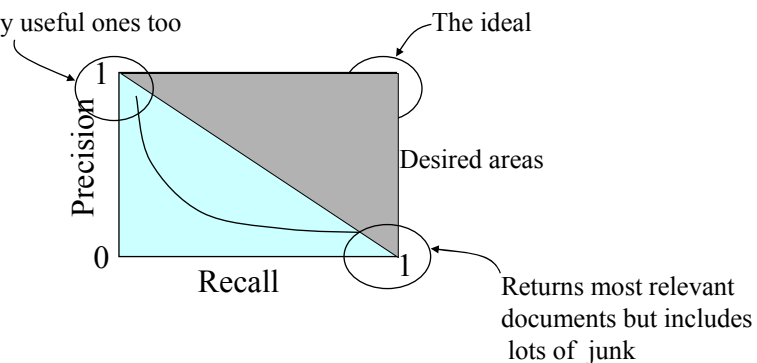
$$\text{Accuracy} = (TP + TN) / (TP + TN + FP + FN)$$

- the data is extremely skewed: normally over 99.9% of the documents are in the nonrelevant category => ?
- precision and recall concentrate the evaluation on the return of true positives (asking what percentage of the relevant documents have been found and how many false positives among them)



## Αντιπαραβολή (trade-off) μεταξύ βαθμού ανάκλησης και βαθμού ακρίβειας

Returns relevant documents but misses many useful ones too





Θα μπορούσαμε με έναν μόνο αριθμό να χαρακτηρίσουμε την αποτελεσματικότητα ενός συστήματος;



## F-Measure

*Συνδυασμός των δύο μέτρων*



## Μέτρα αξιολόγησης αποτελεσματικότητας: F-Measure

Μέτρο που λαμβάνει υπόψη την Ακρίβεια και την Ανάκληση  
Είναι το αρμονικό μέσο (harmonic mean) της ανάκλησης και ακρίβειας για το ερώτημα  $i$ :

$$F(i) = \frac{2P(i)R(i)}{P(i) + R(i)} = \frac{2}{\frac{1}{R(i)} + \frac{1}{P(i)}}$$

Παίρνει τιμές στο  $[0, 1]$

0 όταν έστω ένα από τα 2 είναι μηδέν

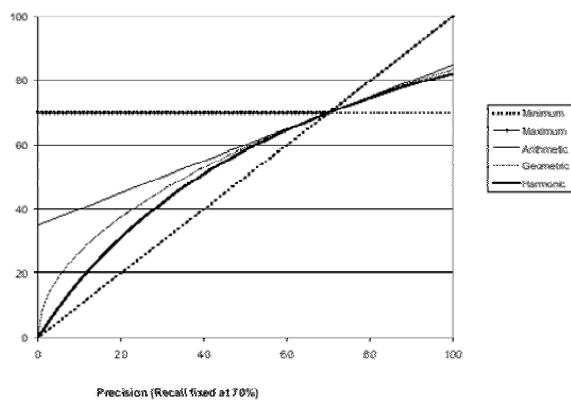
1 όταν και τα 2 είναι 1

**Γιατί αρμονικό μέσο και όχι αριθμητικό;**

100% recall (return all: 10.000 document – 1 relevant) => 50% Arithmetic mean -- Harmonic mean: 0.02%

Για να πάρουμε υψηλή τιμή αρμονικού μέσου χρειαζόμαστε υψηλό P και υψηλό R (το αριθμητικό μέσο είναι γενικά πιο κοντά στο μικρότερο).

Πάντα μικρότερο ή ίσο του αριθμητικού μέσου





## E-Measure

*Πως μπορούμε να δώσουμε βάρος σε ένα από τα δύο μέτρα;*

### Motivation

Web surfers - high precision for results in first page

Professionals (e.g. paralegals) as high recall as possible and tolerate fairly low precision

Individuals searching their hard disk also high recall



Μέτρα αξιολόγησης αποτελεσματικότητας:

## E Measure (παραμετρικό F Measure)

Παραλλαγή του F measure που μας επιτρέπει να δώσουμε περισσότερη έμφαση (βάρος) π.χ. στην ακρίβεια:

$$F(i) = \frac{1}{a \frac{1}{R(i)} + (1-a) \frac{1}{P(i)}}$$

Weight  $\alpha \in [0, 1]$



## Μέτρα αξιολόγησης αποτελεσματικότητας: E Measure (παραμετρικό F Measure)

$$E(i) = \frac{(1 + \beta^2)P(i)R(i)}{\beta^2 P(i) + R(i)} = \frac{(1 + \beta^2)}{\frac{\beta^2}{R(i)} + \frac{1}{P(i)}}$$

$$\beta^2 = (1-\alpha)/\alpha$$

- Η τιμή του  $\beta$  ρυθμίζει το trade-off:
  - $\beta = 1$ : Equally weight precision and recall (E-measure = F-measure).
  - $\beta > 1$ : Weights precision more.
  - $\beta < 1$ : Weights recall more.



## Ο Προσδιορισμός της Ανάκλησης είναι καμιά φορά δύσκολος (δυσκολότερος της Ακρίβειας)

**ΠΡΟΒΛΗΜΑ:** Ο συνολικός αριθμός των εγγράφων που είναι συναφή με μια επερώτηση μπορεί να είναι άγνωστος

- Π.χ. Αυτό συμβαίνει με τον Ιστό, και σε recommendation systems

**Τρόποι Αντιμετώπισης αυτού του Προβλήματος**

- Δειγματοληψία (sampling)
  - Sample across the database and perform relevance judgment only on these items.
- Pooling
  - Apply different retrieval algorithms to the same database for the same query. Then the aggregate of relevant items is taken as the total relevant set.

*[Τρόπους συνάθροισης διατάξεων (rank aggregation) θα δούμε στο μάθημα περί μετα-μηχανών αναζήτησης]*



## Fallout



### Μέτρα αξιολόγησης αποτελεσματικότητας: Fallout Rate

#### Προβλήματα της Ακρίβειας και Ανάκλησης:

- Ο αριθμός των μη-συναφών εγγράφων δεν λαμβάνεται υπόψη (μη συμμετρικό μέτρο)
- Η Ανάκληση δεν ορίζεται αν η συλλογή δεν έχει κανένα συναφές έγγραφο.
- Η Ακρίβεια δεν ορίζεται αν δεν ανακληθεί κανένα έγγραφο

$$Fallout = \frac{\text{no. of nonrelevant items retrieved}}{\text{total no. of nonrelevant items in the collection}}$$

Πάλι αφορά ένα ερώτημα





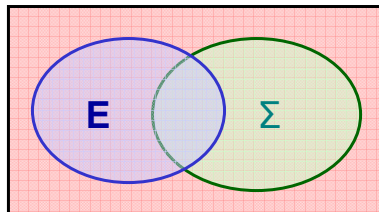
## Fallout

Έστω ένα ερώτημα  $q$

$\Sigma$ : Συναφή (με το ερώτημα  $q$ )

$\Sigma^c$ : Μη-Συναφή

Συλλογή εγγράφων



$E$ : Ευρεθέντα (από το ΣΑΠ)

$$\text{Fallout} = \frac{|E \cap \Sigma^c|}{|\Sigma^c|}$$

Γενικά, θέλουμε το fallout να είναι κοντά στο 0.

Θέλουμε να μεγιστοποιήσουμε την ανάκληση και να ελαχιστοποιήσουμε το fallout.

Εξέταση του γραφήματος fallout-recall graph.



## Καμπύλες Ακρίβειας/Ανάκλησης (Precision/Recall Curves)

Μας ενδιαφέρει και η θέση ενός εγγράφου στο αποτέλεσμα



## Μέτρα αξιολόγησης αποτελεσματικότητας: Σημεία και Καμπύλες Ανάκλησης/Ακρίβειας

### Κίνητρο:

Ο χρήστης δεν «καταναλώνει» όλη την απάντηση μονομιάς.  
Αντίθετα αρχίζει από την κορυφή της λίστας των αποτελεσμάτων

Αυτό δεν λαμβάνεται υπόψη από τα μέτρα Recall και Precision

Θεωρείστε την περίπτωση που:

Answer(System1,q) = <N N N N N N N R R R>

Answer(System2,q) = <R R R N N N N N N N>

N: συμβολίζει ένα non-relevant έγγραφο

R: συμβολίζει ένα relevant έγγραφο

Η Ακρίβεια και η Ανάκληση των δυο συστημάτων είναι η ίδια! :(



## Μέτρα αξιολόγησης αποτελεσματικότητας: Σημεία και Καμπύλες Ανάκλησης/Ακρίβειας

### Παράδειγμα 1

Έστω το σύνολο των σχετικών εγγράφων είναι

$R_q = \{d_3, d_5, d_9, d_{25}, d_{39}, d_{44}, d_{56}, d_{71}, d_{89}, d_{123}\}$

Η απάντηση για την q:

- |              |              |               |
|--------------|--------------|---------------|
| 1. $d_{123}$ | 6. $d_9$     | 11. $d_{38}$  |
| 2. $d_{84}$  | 7. $d_{511}$ | 12. $d_{48}$  |
| 3. $d_{56}$  | 8. $d_{129}$ | 13. $d_{250}$ |
| 4. $d_6$     | 9. $d_{187}$ | 14. $d_{113}$ |
| 5. $d_8$     | 10. $d_{25}$ | 15. $d_3$     |



## Σημεία και Καμπύλες Ανάκλησης/Ακρίβειας (II) (Recall/Precision Points and Curves)

### Αντιμετώπιση Προβλήματος: Χρήση Recall/Precision Curves

#### Τρόπος υπολογισμού:

- 1) Για δοθείσα επερώτηση, παίρνουμε τη διατεταγμένη λίστα από το ΣΑΠ  
Σημείωση: αν δεν πάρουμε όλη την απάντηση αλλά ένα τμήμα της, τότε το σύνολο των Ευρεθέντων αλλάζει, και άρα θα πάρουμε διαφορετικές recall/precision μετρήσεις
- 2) Σημειώνουμε κάθε έγγραφο της λίστας που είναι συναφές (βάσει της χειρονακτικά μαρκαρισμένης συλλογής)
- 3) Υπολογίζουμε ένα ζεύγος τιμών Ανάκλησης/Ακρίβειας για κάθε θέση της διατεταγμένης λίστας που περιέχει ένα συναφές έγγραφο.



## Recall/Precision Points and Curves: Παράδειγμα

### Παράδειγμα 2

Έστω  $|\Sigma \text{συναφή}|=6$

n	doc #	relevant		
1	588	x		
2	589	x		
3	576			
4	590	x		
5	986			
6	592	x		
7	984			
8	988			
9	578			
10	985			
11	103			
12	591			
13	772	x		
14	990			

	Recall	Precision
→	R=1/6=0.167;	P=1/1=1
→	R=2/6=0.333;	P=2/2=1
→	R=3/6=0.5;	P=3/4=0.75
→	R=4/6=0.667;	P=4/6=0.667
→	R=5/6=0.833;	P=5/13=0.38

Missing one  
relevant document.  
Never reach  
100% recall



## Plotting the Recall/Precision Points

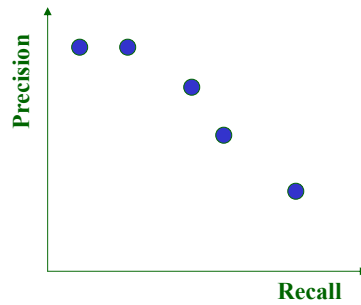
$$R=1/6=0.167; \quad P=1/1=1$$

$$R=2/6=0.333; \quad P=2/2=1$$

$$R=3/6=0.5; \quad P=3/4=0.75$$

$$R=4/6=0.667; \quad P=4/6=0.667$$

$$R=5/6=0.833; \quad P=5/13=0.38$$



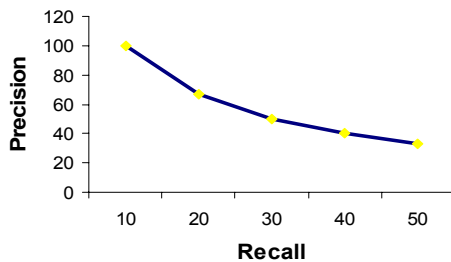
Σχεδιάζουμε την καμπύλη της ακρίβειας (άξονας y) με την ανάκληση (άξονας x)

+ Παράδειγμα 1



## Παράδειγμα

$$R_q = \left\{ \begin{array}{l} d_3, d_5, d_9, d_{25}, d_{39}, \\ d_{44}, d_{56}, d_{71}, d_{89}, d_{123} \end{array} \right\}$$

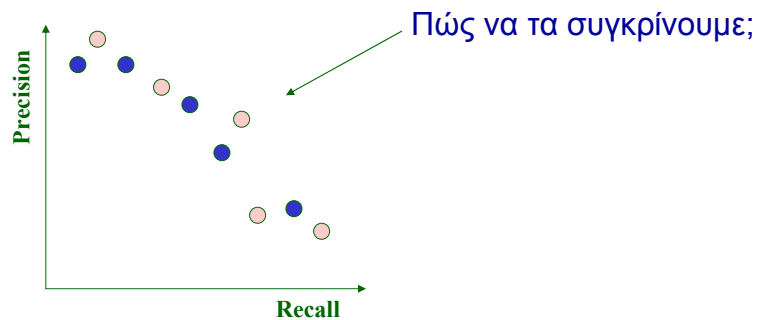


Rank	Doc	Rel	$R_{recall}$	$P_{precision}$
0			0%	0%
1	$d_{123}$	✓	10%	100%
2	$d_{84}$		10%	50%
3	$d_{56}$	✓	20%	67%
4	$d_6$		20%	50%
5	$d_{84}$		20%	40%
6	$d_9$	✓	30%	50%
7	$d_{511}$		30%	43%
8	$d_{129}$		30%	38%
9	$d_{187}$		30%	33%
10	$d_{25}$	✓	40%	40%
11	$d_{38}$		40%	36%
12	$d_{48}$		40%	33%
13	$d_{250}$		40%	31%
14	$d_{113}$		40%	29%
15	$d_3$	✓	50%	33%



## Σύγκριση δύο συστημάτων

- Σύστημα 1
- Σύστημα 2



Information Retrieval 2009-2010

41



## Interpolating a Recall/Precision Curve

Σκοπός: Δυνατότητα σύγκρισης διαφορετικών συστημάτων

Τρόπος:

Χρήση κανονικοποιημένων επιπέδων ανάκλησης (standard recall levels)

Παράδειγμα καθιερωμένων επιπέδων ανάκλησης (πλήθος επιπέδων: 11):

Standard Recall levels at 0%, 10%, 20%, ..., 100%

$$r_j \in \{0.0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0\}$$

$$r_0 = 0.0, r_1 = 0.1, \dots, r_{10} = 1.0$$

Information Retrieval 2009-2010

42



## Interpolating a Recall/Precision Curve

Υπολογίζουμε μια τιμή ακρίβειας για κάθε *standard recall level*:

Πως;

ως ακρίβεια στο  $j$  επίπεδο ανάκλησης ορίζουμε τη **μέγιστη ακρίβεια** που εμφανίζεται μεταξύ των βαθμών ανάκλησης  $j$  και  $j+1$

$$P(r_j) = \max_{r_j \leq r \leq r_{j+1}} P(r)$$

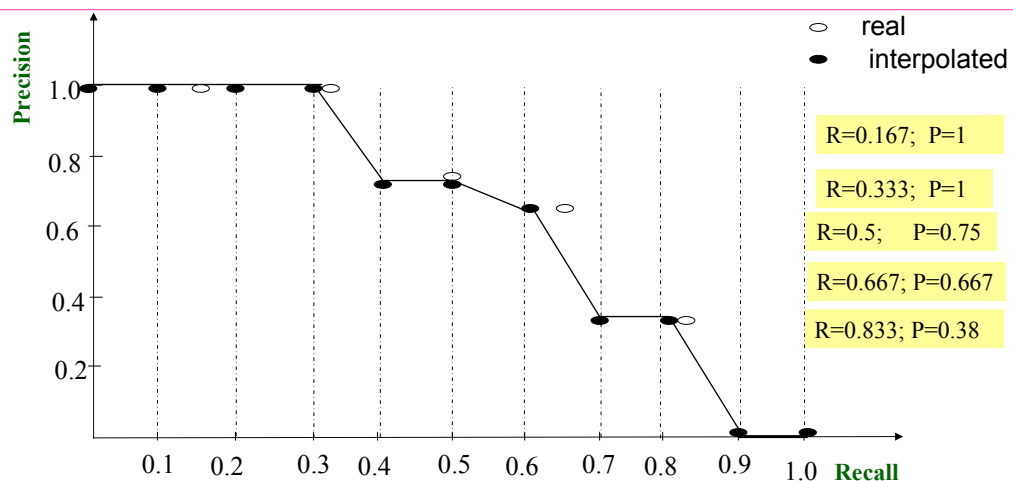
Γιατί;

Αυτό στηρίζεται στην παρατήρηση ότι όσο η ανάκληση μεγαλώνει τόσο η ακρίβεια μειώνεται

- Άρα, είναι λογικό να στοχεύουμε προς μια καμπύλη παρεμβολής (interpolation) που δίδει μια μονότονα φθίνουσα συνάρτηση – *όσο μεγαλώνει η ανάκληση τόσο πέφτει η ακρίβεια*



## Interpolating a Recall/Precision Curve: Παράδειγμα

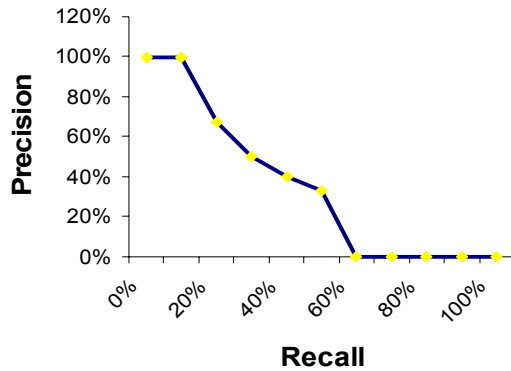


Σημείωση: Από τα 5 ζεύγη (P,R) που είχαμε πήγαμε στα 11



## Παρεμβολή

Εφαρμογή παρεμβολής για τα 11 επίπεδα ανάκλησης



Level	Iterpolated Precision
0%	100%
10%	100%
20%	67%
30%	50%
40%	40%
50%	33%
60%	0%
70%	0%
80%	0%
90%	0%
100%	0%

Information Retrieval 2009-2010

45



## Μέτρα αξιολόγησης αποτελεσματικότητας: Σημεία και Καμπύλες Ανάκλησης/Ακρίβειας

### Παράδειγμα 3

Έστω το σύνολο των σχετικών ερωτημάτων είναι

$$R_q = \{d_3, d_{56}, d_{129}\}$$

Η απάντηση για την  $q$ :

- |              |              |               |
|--------------|--------------|---------------|
| 1. $d_{123}$ | 6. $d_9$     | 11. $d_{38}$  |
| 2. $d_{84}$  | 7. $d_{511}$ | 12. $d_{48}$  |
| 3. $d_{56}$  | 8. $d_{129}$ | 13. $d_{250}$ |
| 4. $d_6$     | 9. $d_{187}$ | 14. $d_{113}$ |
| 5. $d_8$     | 10. $d_{25}$ | 15. $d_3$     |

Information Retrieval 2009-2010

Yannis Tzitzikas, U. of Crete

46



Τι κάνουμε αν έχουμε πολλά ερωτήματα στη συλλογή αξιολόγησης;



## Μέση Καμπύλη Ανάκλησης/Ακρίβειας

- Προκύπτει αξιολογώντας την αποτελεσματικότητα του συστήματος με ένα μεγάλο πλήθος επερωτήσεων
- Υπολογίζουμε τη **μέση ακρίβεια** σε κάθε **standard recall level** για όλες τις επερωτήσεις

$$\bar{P}(r) = \sum_{i=1}^{N_q} \frac{P_i(r)}{N_q}$$

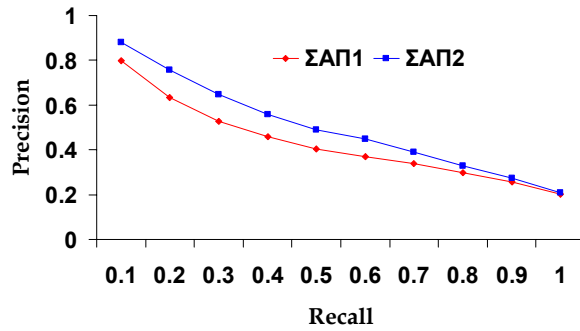
- $N_q$  – number of queries
- $P_i(r)$  - precision at recall level  $r$  for  $i^{\text{th}}$  query

- Σχεδιάζουμε τη μέση precision/recall καμπύλη η οποία εκφράζει την επίδοση του συστήματος στη συλλογή





## Σύγκριση Συστημάτων

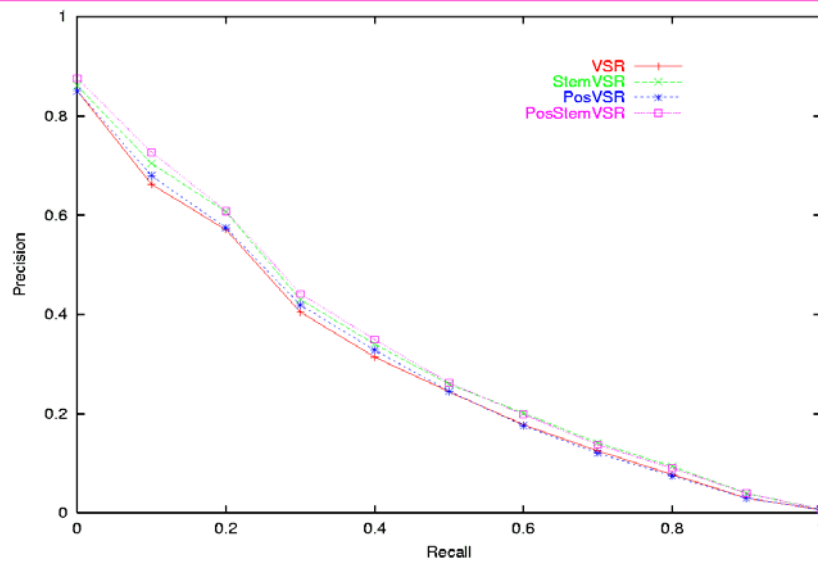


- Η καμπύλη που είναι πιο κοντά στην πάνω-δεξιά γωνία του γραφήματος υποδηλώνει καλύτερη επίδοση

Το ΣΑΠ2 είναι καλύτερο από το ΣΑΠ1

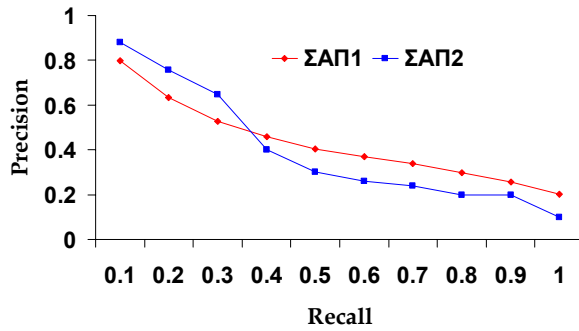


## Σύγκριση Συστημάτων (II)





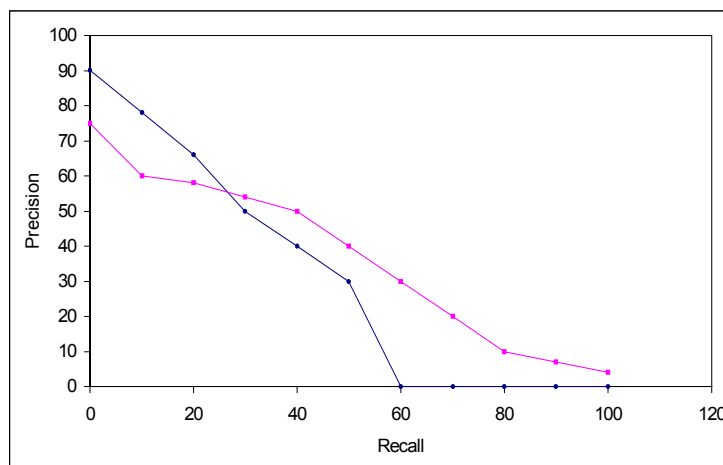
### Σύγκριση Συστημάτων (III)



Το ΣΑΠ2 έχει καλύτερη ακρίβεια στα χαμηλά επίπεδα ανάκλησης  
Το ΣΑΠ1 έχει καλύτερη ακρίβεια στα υψηλά επίπεδα ανάκλησης



### Μέτρα αξιολόγησης αποτελεσματικότητας: Σημεία και Καμπύλες Ανάκλησης/Ακρίβειας





## Cut-off Values – precision at k

### Document Cutoff Values

Μέση ακρίβεια όταν έχουν έχουμε δει top 5, 10, 15, 20, 30, 50 ή 100 έγγραφα

Περισσότερη πληροφορία για τον ranking αλγόριθμο

Precision at k:

**Advantage:** not requiring any estimate of the size of the set of relevant documents

**Disadvantages:**

the least stable of the commonly used evaluation measures

does not average well, since the total number of relevant documents for a query has a strong influence on precision at k.



## Average precision at seen documents

Μέση Ακρίβεια στα σχετικά έγγραφα της απάντησης (Average Precision at Seen Relevant Documents)

Για το Παράδειγμα 1:

$$(1 + 0.66 + 0.5 + 0.4 + 0.3)/5 = 0.57$$

Ευνοεί τα συστήματα που ανακτούν τα σχετικά έγγραφα γρήγορα

Φυσικά, ο αλγόριθμος μπορεί να έχει καλή συμπεριφορά σε αυτά τα έγγραφα αλλά όχι συνολικά



## R-precision

- **R-Precision:** Η ακρίβεια στην R θέση της διάταξης της απάντησης μιας επερώτησης όπου R είναι το πλήθος των συναφών στην ερώτηση εγγράφων

n	doc #	relevant
1	588	x
2	589	x
3	576	
4	590	x
5	986	
6	592	x
7	984	
8	988	
9	578	
10	985	
11	103	
12	591	
13	772	x
14	990	

Δηλαδή, στην «τέλεια» απάντηση  
R-Precision = 1

R = # of relevant docs = 6

R-Precision =  $4/6 = 0.67$



## R-precision

Στο παράδειγμα 1:

Μέγεθος απάντησης 10, άρα κοιτάμε το precision στα πρώτα 10 έγγραφα της απάντησης



## Ένα μέτρο ανά ερώτηση

Υπάρχουν περιπτώσεις που δεν θέλουμε να πάρουμε μέση τιμή σε όλες τις ερωτήσεις γιατί:

- Με αυτό τον τρόπο μπορεί να μη δούμε τις ιδιαιτερότητες (anomalies) του ΣΑΠ – να κρύβονται στο μέσο όρο
- Στην περίπτωση που συγκρίνουμε 2 ΣΑΠ θέλουμε να δούμε τη σχετική τους απόδοση σε διαφορετικά σύνολα ερωτήσεων (ένας καλύτερος του άλλου σε συγκεκριμένου τύπου ερωτήσεις)

**ΣΤΟΧΟΣ:** Μια τιμή ανά ερώτηση που να λαμβάνει υπ' όψει και τη σειρά των σχετικών εγγράφων στην απάντηση



## R-precision

- Ερωτήματα:
  - Αν έχουμε πολλές επερωτήσεις αξιολόγησης, πώς υπολογίζεται αυτό το μέτρο;
  - Πως μπορούμε να συγκρίνουμε 2 συστήματα βάσει του R-Precision ;
- Απάντηση:
  - Χρησιμοποιώντας πολλές επερωτήσεις αξιολόγησης μπορούμε να σχεδιάσουμε το **Ιστόγραμμα Ακρίβειας (Precision Histogram)**.



## R-precision

Έστω 2 συστήματα A και B και k επερωτήσεις αξιολόγησης  $q_1..q_k$

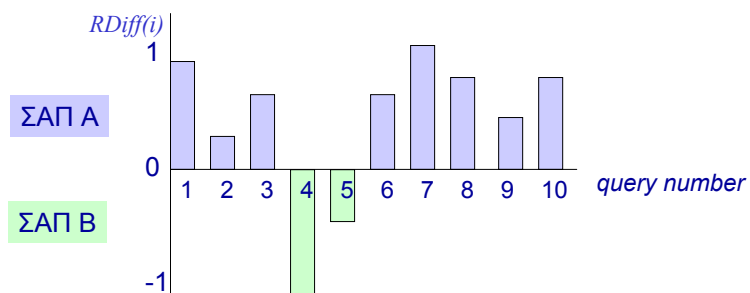
Τρόπος υπολογισμού του ιστογράμματος ακρίβειας:

- Για κάθε  $i = 1$  έως  $k$ 
  - $R(i) :=$  το πλήθος των συναφών εγγράφων της επερώτησης  $q_i$
  - $RPA(i) :=$  Το Ri-precision του συστήματος A για την  $q_i$
  - $RPB(i) :=$  Το Ri-precision του συστήματος B για την  $q_i$
  - Ορίζουμε τη διαφορά ως εξής:  $RDiff(i) := RPA(i) - RPB(i)$
- Κάνουμε την γραφική παράσταση των  $(i, RDiff(i))$  (για  $i = 1..k$ )

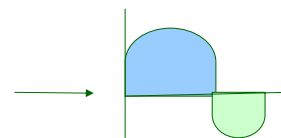


## R-precision + histogram

Παράδειγμα με 10 επερωτήσεις:



Μπορούμε κατόπιν να ταξινομήσουμε ως προς RiDiff ώστε να πάρουμε ένα πιο παραστατικό διάγραμμα





## Expected Search Length



## Expected Search Length Αναμενόμενο μήκος αναζήτησης [Cooper 68]

- **Ορισμός**
  - Το **μέσο** πλήθος εγγράφων που πρέπει να εξεταστούν προκειμένου να ανακτήσουμε ένα **συγκεκριμένο** πλήθος συναφών εγγράφων.
- **Παρατηρήσεις**
  - Δεν είναι ένας αριθμός αλλά συνάρτηση του αριθμού των συναφών εγγράφων που επιθυμούμε
  - Μπορεί να παρασταθεί με πίνακα ή με γράφημα



## Expected Search Length

– Μπορεί να παρασταθεί με πίνακα ή με γράφημα

• Πίνακας

Rel Docs	Search Length
1	2.0
2	4.2
3	5.4
4	6.6
5	7.8

(στα πρώτα δυο στοιχεία της απάντησης υπάρχει ένα συναφές)

Στα 4.2 πρώτα έγγραφα υπάρχουν δύο συναφή.

Το 4.2 είναι είτε μέσος όρος (που προέκυψε κάνοντας πολλές μετρήσεις) ή/και λόγω εγγράφων που έλαβαν ίδιο βαθμό συνάφειας.

• Μπορούμε όμως να υπολογίσουμε και έναν «μέσο όρο»:

–  $(2/1 + 4.2/2 + 5.4/3 + 6.6/4 + 7.8/5) / 5 =$

$(2 + 2.1 + 1.8 + 1.65 + 1.56)/5 = 9.11/5=1.82$

– Χονδρικά, ο χρήστης χρειάζεται να εξετάζει 82% παραπάνω έγγραφα από τα επιδιωκόμενα συναφή (π.χ. αν θέλει 20 θα χρειαστεί να εξετάσει τα πρώτα  $1.82*20=36$  έγγραφα)



## Expected Search Length

Η έννοια του expected search length μας είναι επίσης χρήσιμη προκειμένου να κάνουμε ακριβείς μετρήσεις στην περίπτωση που οι απαντήσεις του συστήματος δεν είναι μια γραμμική ακολουθία εγγράφων, αλλά μια γραμμική ακολουθία συνόλων εγγράφων.

### Παράδειγμα

•  $\text{Answer}(\text{System1}, q) = \langle d8, d2, \{d3, d4\}, d1 \rangle$

– Αυτό σημαίνει ότι τα d3 και d4 έλαβαν τον ίδιο βαθμό συνάφειας (άρα βρίσκονται και τα δύο στην 3<sup>η</sup> θέση της κατάταξης)

•  $\text{Answer}(\text{System2}, q) = \langle d1, \{d2, d3\}, d8 \rangle$

• Ερώτηση: Αν ξέρουμε ότι η q έχει δύο συναφή έγγραφα, συγκεκριμένα τα d1 και d3, ποιά είναι η R-Precision του System1 και ποια του System2 ?





## Ο Προσδιορισμός της Ανάκλησης είναι καμιά φορά δύσκολος (δυσκολότερος της Ακρίβειας)

**ΠΡΟΒΛΗΜΑ:** Ο συνολικός αριθμός των εγγράφων που είναι συναφή με μια επερώτηση μπορεί να είναι άγνωστος

- Π.χ. Αυτό συμβαίνει με τον Ιστό, και σε recommendation systems

**Τρόποι Αντιμετώπισης αυτού του Προβλήματος**

- **Δειγματοληψία (sampling)**
  - Sample across the database and perform relevance judgment only on these items.
- **Pooling**
  - Apply different retrieval algorithms to the same database for the same query. Then the aggregate of relevant items is taken as the total relevant set.

*[Τρόπους συνάθροισης διατάξεων (rank aggregation) θα δούμε στο μάθημα περί μετα-μηχανών αναζήτησης]*



## *User-Oriented Measures*



## Πιο υποκειμενικά μέτρα Συνάφειας

- **Novelty Ratio (ποσοστό “καινοτομίας”):**  
Το ποσοστό των ανακτημένων και συναφών εγγράφων ( $E \cap S$ ) των οποίων την ύπαρξη ο χρήστης αγνοούσε (πριν την αναζήτηση).
  - Μετράει την ικανότητα εύρεσης νέας πληροφορίας σε ένα θέμα.
- **Coverage Ratio (ποσοστό κάλυψης):**  
Το ποσοστό των ανακτημένων και συναφών εγγράφων ( $E \cap S$ ) σε σχέση με το σύνολο των συναφών εγγράφων τα οποία είναι γνωστά στο χρήστη πριν την αναζήτηση.
  - Relevant when the user wants to locate documents which they have seen before (e.g., the budget report for Year 2007).

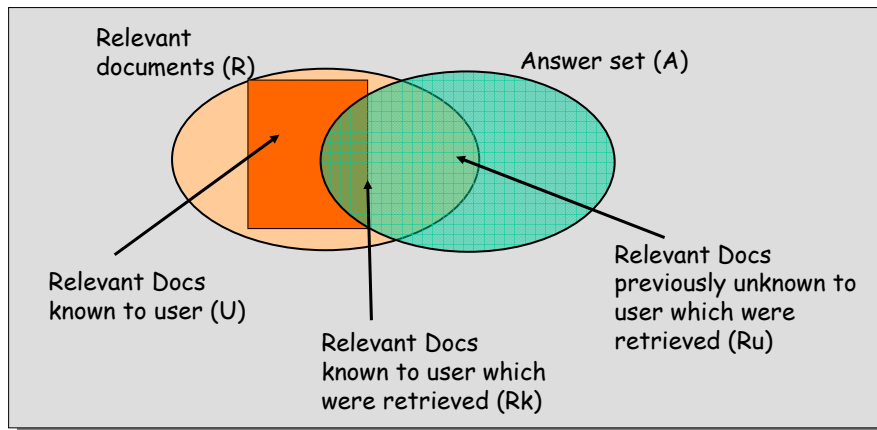


## Άλλοι παράγοντες αξιολόγησης

- **Ανθρώπινη προσπάθεια (User effort):**  
Work required from the user in formulating queries, conducting the search, and screening the output.
- **Χρόνος απόκρισης (Response time):**  
Time interval between receipt of a user query and the presentation of system responses.
- **Μορφή παρουσίασης (Form of presentation):**  
Influence of search output format on the user’s ability to utilize the retrieved materials.
- **Κάλυψη συλλογής (Collection coverage):**  
Extent to which any/all relevant items are included in the document corpus.



## Μέτρα Προσανατολισμένα στο Χρήστη



Information Retrieval 2009-2010

69



## Μέτρα Προσανατολισμένα στο Χρήστη

- Different users might have a different interpretation of which document is relevant or not
- User-oriented measures:
  - Coverage ratio:
  - Novelty ratio:
  - Relative ratio – the ratio between the # of relevant docs found and the # of relevant docs the user expected to find
  - Recall effort – the ratio between the # of relevant docs the user expected to find and the # of docs examined in an attempt to find the expected relevant docs.

$$\text{novelty} = \frac{|R_u|}{|R_u| + |R_k|}$$

$$\text{coverage} = \frac{|R_k|}{|U|}$$

Information Retrieval 2009-2010

70



## Result snippets

Present a **results list** that will be informative to the user

In many cases the user will not want to examine *all* the returned documents and so we want to make the results list informative => Provide a

**snippet**: a short summary of the document designed to allow the user to decide its relevance

Title + Short Summary

- Static (same for all queries)
- Dynamic



## Static summary

comprised of either or both:

- a subset of the document

simplest form: the first two sentences or 50 words of a document, or particular zones of a document (such as the title and author)

- metadata associated with the document

an alternative way to provide an author or date, or may include elements which are designed to give a summary, such as the description metadata which can appear in the meta element of a web HTML page.

typically **extracted and cached at indexing time**, in such a way that it can be retrieved and presented **quickly** when displaying search results, whereas having to access the actual document content



## Static summary

### text summarization

#### NLP

choose to present sentences from the original document

how to select good sentences, combine:

- **positional factors,**

favoring the first and last paragraphs of documents and the first and last sentences of paragraphs, with

- **content factors,**

emphasizing sentences with key terms, which have low document frequency in the collection as a whole, but high frequency and good distribution across the particular document being returned.

In sophisticated NLP approaches, the system **synthesizes sentences** for a summary, either by doing full text generation or by editing and perhaps combining sentences used in the document.



## Dynamic summary

display one or more **“windows”** on the document: pieces that have the most utility to the user in evaluating the document with respect to their information need.

Usually these windows contain one or several of the query terms: *keyword-in-context (KWIC) snippets*

If the query is found as a **phrase**, occurrences of the phrase in the document will be shown as the summary.

**If not**, windows within the document that contain multiple query terms will be selected. Commonly just stretch some number of words to the left and right of the query terms.

NLP techniques: users prefer snippets that read well because they contain complete phrases.



## Dynamic summary

greatly improving the usability of IR systems, but they present a complication for IR system design.

cannot be precomputed,

1. locally cache all the documents at index time (notwithstanding that this approach raises various legal, information security and control issues that are far from resolved)
2. a system can simply scan a document which is about to appear in a displayed results list to find snippets containing the query words.

Given a variety of keyword occurrences in a document, the goal is to choose fragments which are: (i) maximally informative about the discussion of those terms in the document, (ii) self-contained enough to be easy to read, and (iii) short enough to fit within the normally strict constraints on the space available for summaries.



## Dynamic summary

It is common to cache only a generous but fixed size prefix of the document, such as perhaps 10,000 characters. Summaries of documents whose length exceeds the prefix size will be based on material in the prefix only

If a document has been updated since it was last processed by a crawler and indexer, these changes will be neither in the cache nor in the index, but it is the differences between the summary and the actual document content that will be more glaringly obvious to the end user.



## A/B testing

For such a test, precisely one thing is changed between the current system and a proposed system, and a small proportion of traffic (say, 1–10% of users) is randomly directed to the variant system, while most users use the current system.



## Δοκιμασίες Αποτελεσματικότητας-Συλλογές Αναφοράς (TREC)



## Συλλογές Αναφοράς

- TREC collection (Text REtrieval Conference) – 5.8 GB, >1.5 Million Docs.
- CACM - computer science – 3024 articles.
- ISI (CISI) – library science – 1460 docs.
- Cystic Fibrosis (CF) - medicine – 1239 docs.
- CRAN – aeronautics – 1400 docs.
- Time – general articles – 423 docs.
- NPL – electrical engineering – 11429 docs.
- κλπ



## Δοκιμασίες επιδόσεων (Benchmarking)

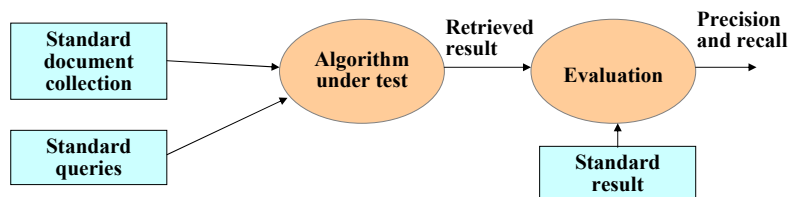
- Η *αναλυτική αξιολόγηση* επίδοσης είναι δύσκολη στα ΣΑΠ διότι πολλά χαρακτηριστικά (συνάφεια, κατανομή λέξεων, κλπ) δύσκολα προσδιορίζονται με μαθηματική ακρίβεια
- Η επίδοσεις συνήθως μετρώνται με *Δοκιμασίες Επιδόσεων (benchmarking)*. Η αξιολόγηση της αποτελεσματικότητας αξιολογείται σε συγκεκριμένες συλλογές εγγράφων, επερωτήσεων και κρίσεις συνάφειας
- Τα αποτελέσματα είναι έγκυρα μόνο στο περιβάλλον που έγινε η αξιολόγηση.





## Δοκιμασίες Επιδόσεων

- A benchmark collection contains:
  - A set of standard documents and queries/topics.
  - A list of relevant documents for each query.
- Standard collections for traditional IR:
  - **Smart collection**: <ftp://ftp.cs.cornell.edu/pub/smart>
  - **TREC**: <http://trec.nist.gov/>



## Τα προβλήματα του Benchmarking

- Τα αποτελέσματα της αξιολόγησης είναι έγκυρα μόνο για τη συγκεκριμένη δοκιμασία αξιολόγησης
- Ο κατασκευή ενός benchmark είναι δύσκολη και χρονοβόρα
- Αφορούν κυρίως κείμενα στα ΑΓΓΛΙΚΑ



## Early Test Collections

- Previous experiments were based on the SMART collection which is fairly small. (<ftp://ftp.cs.cornell.edu/pub/smart>)

Collection Name	Number Of Documents	Number Of Queries	Raw Size (Mbytes)
CACM	3,204	64	1.5
CISI	1,460	112	1.3
CRAN	1,400	225	1.6
MED	1,033	30	1.1
TIME	425	83	1.5

- Different researchers used different test collections and evaluation techniques.



**TREC**  
<http://trec.nist.gov/>



## The TREC Benchmark

- TREC: **T**ext **R**etrieval **C**onference (<http://trec.nist.gov/>)  
Originated from the TIPSTER program sponsored by Defense Advanced Research Projects Agency (DARPA).
- Became an annual conference in 1992, co-sponsored by the National Institute of Standards and Technology (NIST) and DARPA.
- Participants are given parts of a standard set of documents and TOPICS (from which queries have to be derived) in different stages for training and testing.
- Participants submit the P/R values for the final document and query corpus and present their results at the conference.



## Οι στόχοι του TREC

- Provide a common ground for comparing different IR techniques.
  - Same set of documents and queries, and same evaluation method.
- Sharing of resources and experiences in developing the benchmark.
  - With major sponsorship from government to develop large benchmark collections.
- Encourage participation from industry and academia.
- Development of new evaluation techniques, particularly for new applications.
  - Retrieval, routing/filtering, non-English collection, web-based collection, question answering.



## Τα πλεονεκτήματα του TREC

- Large scale (compared to a few MB in the SMART Collection).
- Relevance judgments provided.
- Under continuous development with support from the U.S. Government.
- Wide participation:
  - TREC 1: 28 papers 360 pages.
  - TREC 4: 37 papers 560 pages.
  - TREC 7: 61 papers 600 pages.
  - TREC 8: 74 papers.



## TREC Tasks

- **Ad hoc**: New questions are being asked on a static set of data.
- **Routing**: Same questions are being asked, but new information is being searched. (news clipping, library profiling).
- New tasks added after TREC 5 - Interactive, multilingual, natural language, multiple database merging, filtering, very large corpus (20 GB, 7.5 million documents), question answering.



## TREC Tracks

- **Cross-Language Track**
  - the ability of retrieval systems to find documents that pertain to a topic **regardless of the language** in which the document is written.
  - Also studied in CLEF (Cross-Language Evaluation Forum), and the NTCIR workshops.
- **Filtering Track**
  - user's information need is **stable** (and some relevant documents are known) but there is a **stream of new documents**. For each document, the system must make a binary decision as to whether the document should be retrieved (as opposed to forming a ranked list).
- **Genomics Track**
  - study retrieval tasks in a specific domain, where the domain of interest is **genomics data** (broadly construed to include not just gene sequences but also supporting documentation such as research papers, lab reports, etc.)

Information Retrieval 2009-2010

89



## TREC Tracks (II)

- **HARD Track**
  - achieve **High Accuracy** Retrieval from Documents by leveraging additional information about the searcher and/or the search **context**, through techniques such as **passage retrieval** and using very targeted interaction with the searcher.
- **Interactive Track**
  - A track studying user **interaction** with text retrieval systems. Participating groups develop a consensus experimental protocol and carry out studies with real users using a common collection and set of user queries.
- **Novelty Track**
  - ability to locate **new** (i.e., non-redundant) information.
- **Question Answering Track**
  - a step closer to information retrieval rather than document retrieval. Focus on definition, list, and factoid questions.

Information Retrieval 2009-2010

90



## TREC Tracks (III)

- **Terabyte Track**
  - investigate whether/how the IR community can scale traditional IR test-collection-based evaluation to significantly **larger document collections** than those currently used in TREC. The retrieval task will be an ad hoc task using a static collection of approximately **1 terabyte of spidered web pages** (probably from the .GOV domain).
- **Video Track**
  - research in automatic segmentation, indexing, and content-based retrieval of digital video. Beginning in 2003, the track became an independent evaluation (TRECVID).
- **Web Track**
  - A track featuring search tasks on a document set that is a snapshot of the World Wide Web.



## Χαρακτηριστικά της συλλογής TREC

- Both long and short documents (from a few hundred to over one thousand unique terms in a document).
- Test documents consist of:

WSJ Wall Street Journal articles (1986-1992)	550 M
AP Associate Press Newswire (1989)	514 M
ZIFF Computer Select Disks (Ziff-Davis Publishing)	493 M
FR Federal Register	469 M
DOE Abstracts from Department of Energy reports	190 M



## More Details on Document Collections

- Volume 1 (Mar 1994) - Wall Street Journal (1987, 1988, 1989), Federal Register (1989), Associated Press (1989), Department of Energy abstracts, and Information from the Computer Select disks (1989, 1990)
- Volume 2 (Mar 1994) - Wall Street Journal (1990, 1991, 1992), the Federal Register (1988), Associated Press (1988) and Information from the Computer Select disks (1989, 1990)
- Volume 3 (Mar 1994) - San Jose Mercury News (1991), the Associated Press (1990), U.S. Patents (1983-1991), and Information from the Computer Select disks (1991, 1992)
- Volume 4 (May 1996) - Financial Times Limited (1991, 1992, 1993, 1994), the Congressional Record of the 103rd Congress (1993), and the Federal Register (1994).
- Volume 5 (Apr 1997) - Foreign Broadcast Information Service (1996) and the Los Angeles Times (1989, 1990).



## TREC Disk 4,5

TREC Disk 4	Congressional Record of the 103rd Congress approx. 30,000 documents approx. 235 MB
	Federal Register (1994) approx. 55,000 documents approx. 395 MB
	Financial Times (1992-1994) approx. 210,000 documents approx. 565 MB
TREC Disk 5	Data provided from the Foreign Broadcast Information Service approx. 130,000 documents approx. 470 MB
	Los Angeles Times (randomly selected articles from 1989 & 1990) approx. 130,000 document approx. 475 MB



## Δείγμα Εγγράφου (σε SGML)

```
<DOC>
<DOCNO> WSJ870324-0001 </DOCNO>
<HL> John Blair Is Near Accord To Sell Unit, Sources Say </HL>
<DD> 03/24/87</DD>
<SO> WALL STREET JOURNAL (J) </SO>
<IN> REL TENDER OFFERS, MERGERS, ACQUISITIONS (TNM) MARKETING, ADVERTISING (MKT)
TELECOMMUNICATIONS, BROADCASTING, TELEPHONE, TELEGRAPH (TEL) </IN>
<DATELINE> NEW YORK </DATELINE>
<TEXT>
  John Blair & Co. is close to an agreement to sell its TV station advertising representation operation
  and program production unit to an investor group led by James H. Rosenfield, a former CBS Inc.
  executive, industry sources said. Industry sources put the value of the proposed acquisition at more
  than $100 million. ...
</TEXT>
</DOC>
```



## Δείγμα επερώτησης (with SGML)

```
<top>
<head> Tipster Topic Description
<num> Number: 066
<dom> Domain: Science and Technology
<title> Topic: Natural Language Processing
<desc> Description: Document will identify a type of natural language processing technology which is
being developed or marketed in the U.S.
<narr> Narrative: A relevant document will identify a company or institution developing or marketing a
natural language processing technology, identify the technology, and identify one of more features
of the company's product.
<con> Concept(s): 1. natural language processing ;2. translation, language, dictionary
<fac> Factor(s):
<nat> Nationality: U.S.</nat>
</fac>
<def> Definitions(s):
</top>
```





## TREC Properties

- Both documents and queries contain many different kinds of information (fields).
- Generation of the formal queries (Boolean, Vector Space, etc.) is the responsibility of the system.
  - A system may be very good at querying and ranking, but if it generates poor queries from the topic, its final P/R would be poor.



## Two more TREC Document Examples

ZIFF Communications Company	San Jose Mercury News
<pre> &lt;DOC&gt; &lt;DOCNO&gt; ZF109-706-077 &lt;/DOCNO&gt; &lt;DOCID&gt;09 706 077 &amp;O;&lt;/DOCID&gt; &lt;JOURNAL&gt;Business Week Dec 31 1990 n3194 p93(12) &amp;M; &lt;/JOURNAL&gt; &lt;TITLE&gt;Fujitsu means business for America. (Special Advertising Section by Fujitsu Ltd.) (includes related articles on the company's business relationships with Pepsi-Cola, Convex Computer, Greenville EMS, and Sequent Computer Systems)&amp;M; &lt;/TITLE&gt; &lt;TEXT&gt; &lt;ABSTRACT&gt;In establishing itself as a major manufacturer in the computer hardware market, Fujitsu Ltd boasts a long list of corporate customers.&amp;P; The company's client base includes: MCI Telecommunications Corp., Page Composition, Johns Hopkins Hospital, Tiara Computer Systems Inc., Pepsi-Cola, Convex Computer, Greenville EMS, and Sequent Computer Systems Inc. The company stresses its good customer relations and product development aspects, as well as its telecommunications products.&amp;O; &lt;/ABSTRACT&gt; &lt;/TEXT&gt; &lt;DESCRIPT&gt; Company: Fujitsu Ltd. (Marketing) &amp;O; Topic: Marketing Strategy Customer Relations photogra.ph. &amp;M; &lt;/DESCRIPT&gt; &lt;/DOC&gt; </pre>	<pre> &lt;DOC&gt; &lt;DOCNO&gt; SJMN91-06364024 &lt;/DOCNO&gt; &lt;ACCESS&gt; 06364024 &lt;/ACCESS&gt; &lt;CAPTION&gt; Photo; PHOTO: Associated Press; ANOTHER TURNOVER - Kansas City's Leonard Griffin (96) closes in on Raiders quarterback Todd Marinovich, who fumbled on the play. Marinovich also threw four interceptions. &lt;/CAPTION&gt; &lt;DESCRIPT&gt; PROFESSIONAL; FOOTBALL; PLAYOFF; GAME; RESULT; BRIEF &lt;/DESCRIPT&gt; &lt;LEADPARA&gt; Too much excitement on top of too much cold medication may have caused the rapid heartbeat that forced Kansas City linebacker Derrick Thomas out of the ... reliable place-kicker, kicked an 18-yard field goal at 10:26 of the fourth quarter, but he missed two field goals in the first half, from 33 and 47 yards. ... &lt;/TEXT&gt; &lt;FEATURE&gt; PHOTO &lt;/FEATURE&gt; &lt;STATE&gt; CA &lt;/STATE&gt; &lt;WORD_CT&gt; 539 &lt;/WORD_CT&gt; &lt;DATELINE&gt; Sunday, December 29, 1991 00364024.SJ1 &lt;/DATELINE&gt; &lt;COPYRIGHT&gt; Copyright 1991, San Jose Mercury News &lt;/COPYRIGHT&gt; &lt;LANGUAGE&gt; ENG &lt;/LANGUAGE&gt; &lt;/DOC&gt; </pre>



## Another Example of TREC Topic/Query

```
<top>
<head> Tipster Topic Description
<num> Number: 101
<dom> Domain: Science and Technology
<title> Topic: Design of the "Star Wars" Anti-missile Defense System
<desc> Description:
Document will provide information on the proposed configuration, components, and
technology of the U.S.'s "star wars" anti-missile defense system.
<narr> Narrative:
proposed configuration, components, and technology of the U.S.'s "star wars" anti-missile
defense system. The design and technology to be used in the anti-missile defense system
advocated by the Reagan administration, the Strategic Defense Initiative (SDI), also known
as "star wars." Changes of constituent technologies, are also relevant documents.
<con> Concept(s):
1. Strategic Defense Initiative, SDI, star wars, peace shield
2. kinetic energy weapon, kinetic kill, directed energy weapon, laser, particle beam, ERIS
(exoatmospheric reentry-vehicle interceptor system), phased-array radar, microwave
3. anti-satellite (ASAT) weapon, spaced-based technology, strategic defense technologies
<fac> Factor(s):
<nat> Nationality: U.S.
</nat>
<def> Definition(s):
</top>
```



## Αποτελέσματα Αξιολόγησης ενός ΣΑΠ βάσει του TREC

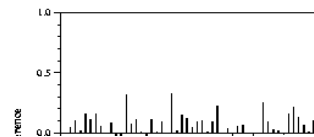
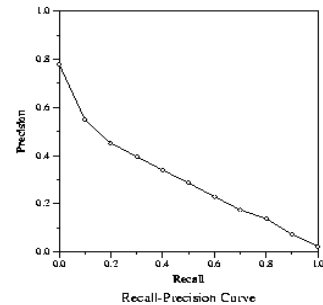
- **Summary table statistics:** Number of topics, number of documents retrieved, number of relevant documents.
- **Recall-precision average:** Average precision at 11 recall levels (0 to 1 at 0.1 increments).
- **Document level average:** Average precision when 5, 10, ..., 100, ... 1000 documents are retrieved.
- **Average precision histogram:** Difference of the R-precision for each topic and the average R-precision of all systems for that topic.



Summary Statistics	
Run Number	Flab8atd2
Run Description	Automatic, title + desc
Number of Topics	50
Total number of documents over all topics	
Retrieved:	50000
Relevant:	4728
Rel ret:	2990

Recall Level Precision Averages	
Recall	Precision
0.00	0.7796
0.10	0.5490
0.20	0.4517
0.30	0.3954
0.40	0.3397
0.50	0.2863
0.60	0.2291
0.70	0.1745
0.80	0.1381
0.90	0.0720
1.00	0.0224
Average precision over all relevant docs	
non interpolated	0.2930

Document Level Averages	
	Precision
At 5 docs	0.5480
At 10 docs	0.4880
At 15 docs	0.4587
At 20 docs	0.4200
At 30 docs	0.3887
At 100 docs	0.2490
At 200 docs	0.1777
At 500 docs	0.1011
At 1000 docs	0.0598
R Precision (precision after R docs retrieved (where R is the number of relevant documents))	
Exact	0.3203



## TREC 2008

<http://trec.nist.gov/call08.html>

- **Blog Track**
  - The purpose of the blog track is to explore information seeking behavior in the blogosphere.
- **Enterprise Track**
  - The purpose of the enterprise track is to study enterprise search: satisfying a user who is searching the data of an organization to complete some task.
- **Legal Track**
  - The goal of the legal track is to develop search technology that meets the needs of lawyers to engage in effective discovery in digital document collections.
- **Million Query Track**
  - The goal of the "million query" track is to test the hypothesis that a test collection built from very many very incompletely judged topics is a better tool than a collection built using traditional TREC pooling.
- **Relevance Feedback Track**
  - The goal of the relevance feedback track is to provide a framework for exploring the effects of different factors on the success of relevance feedback.



## Διάρθρωση Διάλεξης

- Τι εξυπηρετεί η αξιολόγηση;
  - αξιολόγηση αποτελεσματικότητας
- Δυσκολίες της αξιολόγησης
- Αξιολόγηση μέσω Χειρονακτικά Μαρκιαρισμένων Συλλογών
- Μέτρα αξιολόγησης αποτελεσματικότητας
  - Ακρίβεια & Ανάκληση (*Recall & Precision*)
  - Καμπύλες Ακρίβειας/Ανάκλησης
    - Σύγκριση Συστημάτων
  - Εναλλακτικά μέτρα
    - *R-Precision (Precision Histograms)*
    - *F-Measure*
    - *E-Measure*
    - *Fallout, Expected Search Length*
  - *User-Oriented Measures*
- Δοκιμασίες Αποτελεσματικότητας-Συλλογές Αναφοράς (*TREC*)