

# Evaluation of Retrieval Systems

1

## Performance Criteria

1. Expressiveness of query language
  - Can query language capture **information needs**?
2. Quality of search results
  - **Relevance** to users' **information needs**
3. Usability
  - Search Interface
  - Results page format
  - Other?
4. Efficiency
  - Speed affects usability
  - Overall efficiency affects cost of operation
5. Other?

2

## **Quantitative** evaluation

- **Concentrate** on **quality** of search **results**
- Goals for measure
  - Capture **relevance** to user **information need**
  - Allow **comparison** between results of **different systems**
- Measures define for sets of documents returned
- More generally “document” could be any information object

3

## Core measures: **Precision** and **Recall**

- Need binary evaluation by **human judge** of each retrieved document as **relevant/irrelevant**
- Need **know complete set of relevant documents** within collection being searched
- **Recall** =  
$$\frac{\# \text{ relevant documents retrieved}}{\# \text{ relevant documents}}$$
- **Precision** =  
$$\frac{\# \text{ relevant documents retrieved}}{\# \text{ retrieved documents}}$$

4

## Combine recall and precision

**F-score** (aka F-measure) **defined** to be:  
harmonic mean<sup>‡</sup> of precision and recall

$$= \frac{2 * \text{recall} * \text{precision}}{\text{precision} + \text{recall}}$$

<sup>‡</sup> The harmonic mean  $h$  of two numbers  $m$  and  $n$  satisfies  $(n-h)/n = (h-m)/m$ . Also  $= (1/m) - (1/h) = (1/h) - (1/n)$

5

## Use in “modern times”

- Defined in 1950s
- For small collections, these make sense
- For large collections,
  - Rarely know complete set relevant documents
  - Rarely could return complete set relevant documents
- For large collections
  - Rank returned documents
  - **Use ranking!**

6

## Ranked result list

- At any point along ranked list
  - Can look at precision so far
  - Can look at recall so far
    - if know total # relevant docs
    - Google's "about N results" inadequate estimate
- Can focus on points that relevant docs appears
  - If  $m^{\text{th}}$  doc in ranking is  $k^{\text{th}}$  relevant doc so far, precision is  $k/m$ 
    - No a priori ranking on relevant docs

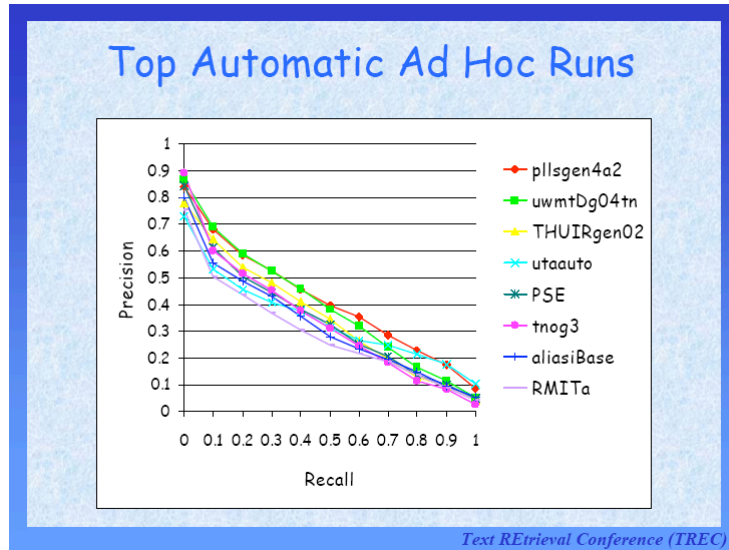
7

## Plot: precision versus recall

- Choose standard recall levels:  $r_1, r_2 \dots$ 
  - Eg 10%, 20% ...
  - Define "precision at recall level  $r_j$ "
$$p(r_j) = \max \text{ over all } r \text{ with } r_j \leq r < r_j + 1 \text{ of}$$
precision when recall  $r$  achieved
    - Similar to *Intro IR* "interpolated precision"

8

Reproduced from presentation "Overview of TREC 2004" by Ellen Voorhees, available from TREC presentations Web site: [trec.nist.gov/presentations/TREC2004/04overview.pdf](http://trec.nist.gov/presentations/TREC2004/04overview.pdf)



9

## Single number characterizations

- Can look at precision at one fixed critical position: "Precision at k"
  - If know are R relevant documents can choose k=R
    - May not want to look that far even if know R
  - Can choose set of S relevant docs, and calc. precision at k=S only with respect to these docs
    - "R-precision" of *Intro IR*
  - For Web search
    - Choose k to be number pages people look at
    - k=? What expecting?

10

## Single number characterizations, cont.

- 1) Record precision at each point a relevant document encountered through ranked list
  - Don't need know *all* relevant docs
  - Can cut off ranked list at predetermined rank
- 2) Average the recorded precisions in (1)  
= average precision for a query result

### Mean Average Precision (MAP):

For a [set of test queries](#), take the mean (i.e. average)

Of the [average precision for each query](#)

- Compare retrieval systems with MAP

11

## Using Measures

- [Statistical significance](#) versus [meaningfulness](#)
- Use more than one measure
- Need some set of relevant docs even if don't have complete set  
[How?](#)
  - Look at TREC studies

12

## Relevance by TREC method

Text Retrieval Conference 1992 to present

- Fixed collection per “track”
  - E.g. “.gov”, CACM articles
- Each competing search engine for a track asked to retrieve documents on several “topics”
  - Search engine turns topic into query
  - Topic description has clear statement of what is to be considered *relevant* by *human judge*

13

### Sample TREC 3 topic:

<num> Number: 168

<title> Topic: Financing AMTRAK

<desc> Description:

A document will address the role of the Federal Government in financing the operation of the National Railroad Transportation Corporation (AMTRAK).

<narr> Narrative: A relevant document must provide information on the government’s responsibility to make AMTRAK an economically viable entity. It could also discuss the privatization of AMTRAK as an alternative to continuing government subsidies. Documents comparing government subsidies given to air and bus transportation with those provided to AMTRAK would also be relevant.

</top>

As appeared in “Overview of the Sixth Text REtrieval Conference (TREC-6),” E. M. Voorhees and D. Harman, in NIST Special Publication 500-240: The Sixth Text REtrieval Conference , 1997.

14

**Sample TREC 7 topic:**

<num>Number: 396

<title> sick building syndrome

<desc>Description:

Identify documents that discuss sick building syndrome or building-related illnesses.

<narr> Narrative:

A relevant document would contain any data that refers to the sick building or building-related illnesses, including illnesses cause by asbestos, air conditioning, pollution controls. Work-related illnesses not caused by the building, such as carpal tunnel syndrome, are not relevant.

From "Overview of the Seventh Text REtrieval Conference (TREC-7)," E. M. Voorhees and D. Harman, in *NIST Special Publication 500-242: The Seventh Text REtrieval Conference*, 1998.

15

## Pooling

- Human judges **can't look at all docs** in collection: thousands to millions
- Pooling **chooses subset of docs** of collection for human judges to rate relevance of
- Assume **docs not in pool not relevant**

16

How construct pool for a topic?  
Let competing search engines decide:

- Choose a parameter  $k$  (typically 100)
- Choose the **top  $k$  docs** as ranked by **each search engine**
- Pool = **union** of these sets of docs  
Between  $k$  and  $(\# \text{ search engines}) * k$  docs in pool
- Give pool to judges for relevance scoring

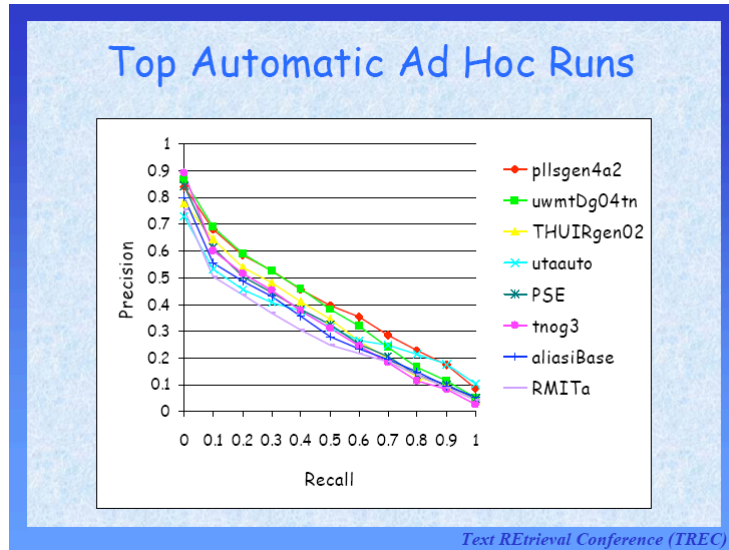
17

## Pooling cont.

- $(k+1)^{\text{st}}$  doc returned by one search engine either irrelevant or ranked higher by another search engine in competition
- In competition, each search engine is judged on **results for top  $r > k$  docs** returned

18

Reproduced from presentation "Overview of TREC 2004" by Ellen Voorhees, available from TREC presentations Web site: [trec.nist.gov/presentations/TREC2004/04overview.pdf](http://trec.nist.gov/presentations/TREC2004/04overview.pdf)



19

## Web search evaluation

- Are different kinds of queries – identified in TREC Web Track – some are:
  - Ad hoc
  - Topic distillation: set of key resources small, 100% recall?
  - Home page: # relevant pages = 1 (except mirrors)
  - Distinguish for competitors or just judges?
- Andrei Broder gave similar categories
  - Information
    - Broad research or single fact?
  - Transaction
  - Navigation

20

## More web/online issues

- Are browser-dependent and presentation dependent issues:
  - On first page of results?
  - See result without scrolling?

21

## Other issues in evaluation

- Degrees of relevance?
  - Discounted Cumulative Gain tries to measure
    - Uses degree of relevance and position in ranking
- Does retrieving highly relevant documents really satisfy users?
  - Subjectivity?
- Are there dependences not accounted for?
- Many searches are interactive

22