

Modern Information Retrieval

Evaluation in information retrieval¹

Hamid Beigy

Sharif university of technology

November 6, 2022

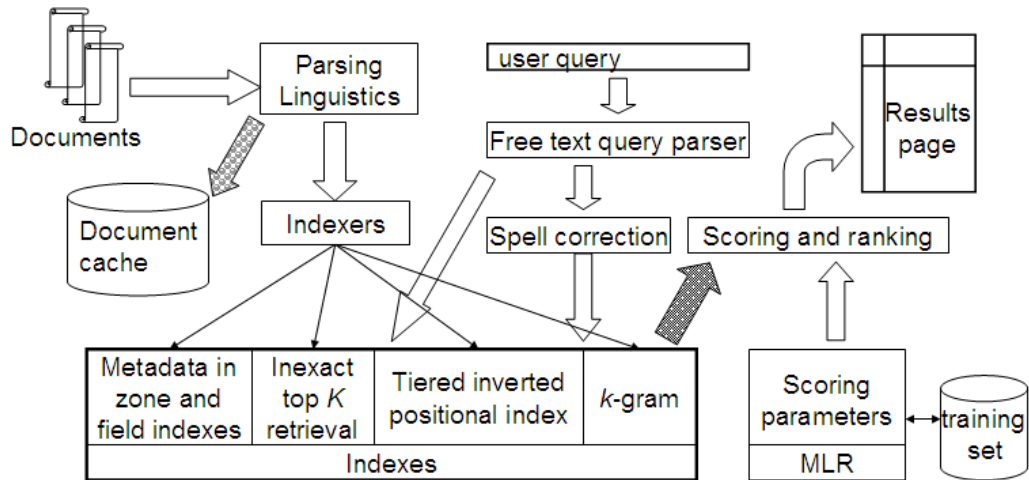


¹Some slides have been adapted from slides of Manning, Yannakoudakis, and Schütze.



1. Introduction
2. Standard test collections
3. Evaluation for unranked retrieval
4. Evaluation for ranked retrieval
5. Assessing relevance
6. System quality and user utility
7. References

Introduction





Framework for the evaluation of an IR system:

1. *Test collection* consisting of
 - ▶ a document *collection*,
 - ▶ a test suite of *information needs*,
 - ▶ a set of *relevance judgments* for each *doc-query* pair
2. *Gold-standard* judgment of relevance.
The classification of a document either as relevant or as irrelevant wrt an information need
3. The test collection must cover at least *50 information needs*
4. The *Development collection* for parameter tuning, if you need.

Standard test collections



1. **Cranfield collection**: 1398 abstracts of journal articles about aerodynamics, gathered in UK in the 1950s, plus 255 queries and exhaustive relevance judgments
2. **TREC** (Text REtrieval Conference): collection maintained by the US National Institute of Standards and Technology since 1992
 - ▶ **TREC Ad Hoc Track**: test collection used for 8 evaluation campaigns led from 1992 to 1999, contains 1.89 million documents and relevance judgments for 450 topics
 - ▶ **TREC 6-8**: test collection providing 150 information needs over 528000 newswires
 - ▶ current state-of-the-art test collection
 - ▶ note that the relevance judgments are not exhaustive



1. **GOV2**: collection also maintained by the NIST, containing 25 millions of web-pages (larger than other test collections, but smaller than current collection supported by WWW search engines)
2. **NTCIR** (Nii Test Collection for IR systems): various test collections focusing on East Asian languages, mainly used for cross-language IR
3. **CLEF** (Cross Language Evaluation Forum): collection focussing on European languages
<http://www.clef-campaign.org>
4. **REUTERS**: Reuters 21578 and REUTERS RCV1 containing respectively 21 578 newswire articles and 806 791 documents, mainly used for text classification

Evaluation for unranked retrieval



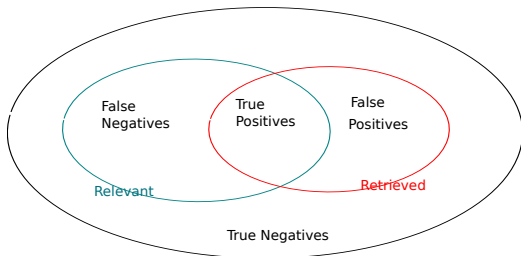
1. Two basic effectiveness measures: *precision* and *recall*

$$Pr = \frac{\#relevant\ retrieved}{\#retrieved}$$

$$Re = \frac{\#relevant\ retrieved}{\#relevant}$$

2. In other terms:

	Relevant	Not relevant
Retrieved	true positive (tp)	false positive (fp)
Not retrieved	false negative (fn)	true negative (tn)



$$Pr = \frac{tp}{tp + fp}$$

$$Re = \frac{tp}{tp + fn}$$



1. **Accuracy**: proportion of the classification relevant/not relevant that is correct

$$accuracy = \frac{tp + tn}{tp + fp + tn + fn}$$

Problem: 99.9% of the collection is usually not relevant to a given query (potential high rate of false positives)

2. Recall and precision are inter-dependent measures:
 - ▶ precision usually decreases while the number of retrieved documents increases
 - ▶ recall increases while the number of retrieved documents increases



1. Measure relating precision and recall:

$$F = \frac{1}{\alpha \times \frac{1}{Pr} + (1 - \alpha) \times \frac{1}{Re}} = \frac{(\beta^2 + 1)Pr \times Re}{\beta^2 Pr + Re}, \beta = \frac{1 - \alpha}{\alpha}$$

2. Most frequently used: balanced F_1 with $\beta = 1$ (or $\alpha = 0.5$):

$$F_1 = \frac{2 \times Pr \times Re}{Pr + Re}$$

3. Uses a harmonic mean rather than an arithmetic one for dealing with extreme values



	Relevant	Not relevant	
Retrieved	20	40	60
Not retrieved	60	1,000,000	1,000,060
	80	1,000,040	1,000,120

$$Pr = \frac{tp}{tp + fp} = \frac{20}{20 + 40} = \frac{1}{3}$$

$$Re = \frac{tp}{tp + fn} = \frac{20}{20 + 60} = \frac{1}{4}$$

$$F_1 = \frac{2 \times \frac{1}{3} \times \frac{1}{4}}{\frac{1}{3} + \frac{1}{4}} = \frac{2}{7}$$

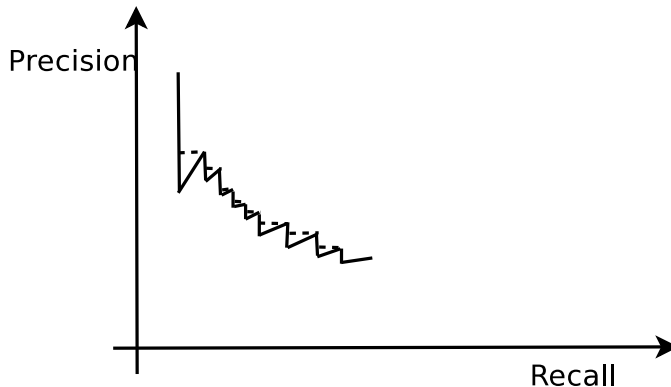
Evaluation for ranked retrieval



1. precision, recall and F-measure are set-based measures (order of documents not taken into account)
2. if we consider the first k retrieved documents, we can compute the precision and recall values
we can plot the relation between precision and recall for each value of k
3. if the $(k + 1)^{\text{st}}$ is not relevant then recall is the same, but precision decreases
4. if the $(k + 1)^{\text{st}}$ is relevant then recall and precision increase



1. Precision-recall curve:



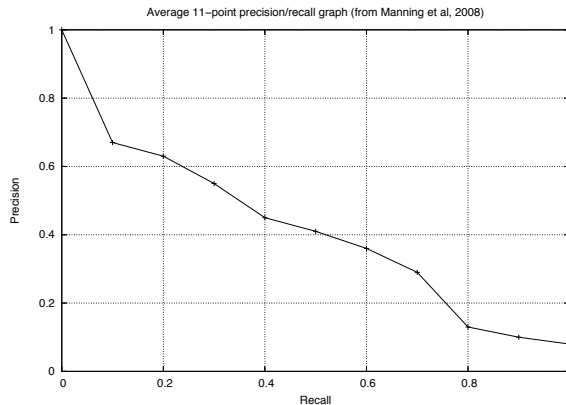
2. For removing jiggles, interpolation of the precision (smoothing):

$$P_{inter}(r) = \max_{r' \geq r} P(r')$$



1. 11-point interpolated average precision:

For each information need, the value P_{inter} is measured for the 11 recall values 0.0, 0.1, 0.2, ... 1.0. The arithmetic mean of P_{inter} for a given recall value over the information needs is then computed.





1. **Precision at k:**

For www search engines, we are interested in the proportion of good results among the k first answers (say the first 3 pages)

This means **precision at a fixed level**

Pros : does not need an estimate of the size of the set of relevant documents

Cons : unstable measure, does not average well because the number of relevant documents for a query has a strong influence on precision at k .



Rank n	Doc
1	d_{12}
2	d_{123}
3	d_4
4	d_{57}
5	d_{157}
6	d_{222}
7	d_{24}
8	d_{26}
9	d_{77}
10	d_{90}

- Blue documents are relevant.
- $P@n$: $P@3=0.33$, $P@5=0.2$, $P@8=0.25$
- $R@n$: $R@3=0.33$, $R@5=0.33$, $R@8=0.66$



1. **Mean Average Precision** (MAP): For an information need, the average precision is the arithmetic mean of the precisions for the set of top k documents retrieved after each relevant document is retrieved

$q_j \in Q$: information need

$\{d_1 \dots d_{m_j}\}$: relevant documents for q_j

R_{jk} : set of ranked retrieved documents from top to d_k

$$MAP(Q) = \frac{1}{|Q|} \sum_{j=1}^{|Q|} \frac{1}{m_j} \sum_{k=1}^{m_j} Pr(R_{jk})$$

when d_l ($1 \leq l \leq j$) is not retrieved, $Pr(R_{jl}) = 0$

Mean Average Precision (MAP) (example)



Query 1		
Rank		$P(doc_j)$
1	X	1.00
2		
3	X	0.67
4		
5		
6	X	0.50
7		
8		
9		
10	X	0.40
11		
12		
13		
14		
15		
16		
17		
18		
19		
20	X	0.25
AVG:		0.564

Query 2		
Rank		$P(doc_j)$
1	X	1.00
2		
3	X	0.67
4		
5		
6		
7		
8		
9		
10		
11		
12		
13		
14		
15	X	0.2
AVG:		0.623

$$MAP = \frac{0.564 + 0.623}{2} = 0.594$$



1. Normalized Discounted Cumulative Gain (NDCG):

Evaluation made for the top k results

$$NDCG(Q, k) = \frac{1}{|Q|} \sum_{j=1}^{|Q|} Z_k \sum_{m=1}^k \frac{2^{R(j,m)} - 1}{\log(1 + m)}$$

where $R(j, d)$ is the score given by assessors to document d for query j Z_k is a normalization factor (perfect ranking at $k = 1$)

Assessing relevance



1. How good is an IR system at satisfying an information need ?
2. Needs an agreement between judges → computable via the **kappa** statistic:

$$kappa = \frac{P(A) - P(E)}{1 - P(E)}$$

where:

$P(A)$ is the proportion of agreements within the judgments

$P(E)$ is the probability that two judges agreed by chance



Consider the following judgments (from Manning et al., 2008):

		Judge 2		
		Yes	No	Total
Judge 1	Yes	300	20	320
	No	10	70	80
	Total	310	90	400

$$P(A) = \frac{300 + 70}{400} = \frac{370}{400} = 0.925$$

$$P(\text{rel}) = \frac{320 + 310}{400 + 400} = 0.7878$$

$$P(\text{notrel}) = \frac{80 + 90}{400 + 400} = 0.2125$$

$$P(E) = P(\text{rel})^2 + P(\text{notrel})^2 = (0.2125)^2 + (0.7878)^2 = 0.665$$

$$\text{kappa} = \frac{P(A) - P(E)}{1 - P(E)} = \frac{0.925 - 0.665}{1 - 0.665} = 0.776$$



1. Interpretation of the kappa statistic k :
 - ▶ $k \geq 0.8$: good agreement
 - ▶ $0.67 \leq k < 0.8$: fair agreement
 - ▶ $k < 0.67$: bad agreement
2. Note that the kappa statistic can be negative if the agreements between judgments are worse than random
3. In case of large variations between judgments, one can choose an assessor as a gold-standard

System quality and user utility



1. Ultimate interest: how satisfied is the user with the results the system gives for each of its information needs ?
2. Evaluation criteria for an IR system:
 - ▶ fast indexing
 - ▶ fast searching
 - ▶ expressivity of the query language
 - ▶ size of the collection supported
 - ▶ user interface (clearness of the input form and of the output list, e.g. snippets, etc)



1. Quantifying user happiness ?

- ▶ For www search engines: **do the users find the information they are looking for?** can be quantified by evaluating the proportion of users getting back to the engine.
- ▶ For intranet search engines: this efficiency can be evaluated by the time spent searching for a given piece of information.
- ▶ General case: user studies evaluating the adequacy of the search engine with the expected usage (eCommerce, etc).

References



1. Chapters 8 of [Information Retrieval Book](#)²

²Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze (2008). *Introduction to Information Retrieval*. New York, NY, USA: Cambridge University Press.



Manning, Christopher D., Prabhakar Raghavan, and Hinrich Schütze (2008).
Introduction to Information Retrieval. New York, NY, USA: Cambridge
University Press.

