

# Informativeness for Adhoc IR Evaluation:

## A measure that prevents assessing individual documents

Romain Deveaud<sup>1</sup>, Véronique Moriceau<sup>2</sup>, Josiane Mothe<sup>3</sup>, and Eric SanJuan<sup>1</sup>

<sup>1</sup> LIA, Univ. Avignon, France,

`romain.deveaud@gmail.com` `eric.sanjuan@univ-avignon.fr`

<sup>2</sup> LIMSI-CNRS, Univ. Paris-Sud, Univ. Paris-Saclay, France, `moriceau@limsi.fr`

<sup>3</sup> IRIT UMR5505, Univ. Toulouse, France, `josiane.mothe@irit.fr`

**Abstract.** Informativeness measures have been used in interactive information retrieval and automatic summarization evaluation. Indeed, as opposed to *ad hoc* retrieval, these two tasks cannot rely on the Cranfield evaluation paradigm in which retrieved documents are compared to static query relevance document lists. In this paper, we explore the use of informativeness measures to evaluate *ad hoc* task. The advantage of the proposed evaluation framework is that it does not rely on an exhaustive reference and can be used in a changing environment in which new documents occur, and for which relevance has not been assessed. We show that the correlation between the official system ranking and the informativeness measure is specifically high for most of the TREC *ad hoc* tracks.

**Keywords:** Information retrieval, Evaluation, Informativeness, Adhoc retrieval

## 1 Introduction

Information Retrieval (IR) aims at retrieving the relevant information from a large volume of available documents. Evaluating IR implies to define evaluation frameworks. In *ad hoc* retrieval, Cranfield framework is the prevailing framework [1]; it is composed of documents, queries, relevance assessments and measures. Moreover, document relevance is considered as independent from the document rank and generally as a Boolean function (a document is relevant or not to a given query) even though levels of relevance can be used [7]. Effectiveness measurement is based on comparing the retrieved documents with the reference list of relevant documents. Moreover, it is based on the assessment assumption, that is the relevance of documents is known in advance for each query. It implies that the collection is static since it is assessed by humans. Cranfield paradigm facilitates reproductibility of experiments: at any time it is possible to evaluate a new IR method and to compare it against previous results; this is one of its main strengths. However, such a framework is not usable in changing environments when new documents are continuously added.

As opposed to Cranfield document relevance independency assumption, informativeness expresses the dependency of document relevance and takes into

account the interactive nature of IR [8]. Indeed, one limitation of Cranfield-based evaluation is that relevance is encoded by documents [4]. Moreover, document relevance assessment is a clear limitation in dynamic context, when new documents are continuously added.

Nugget-based evaluation has been introduced to tackle this problem: rather than considering document relevance, it considers information relevance [4]. This method makes it possible to consider documents that have not been evaluated to be labeled as relevant or not, simply because they contain relevant information or not. Similar assumption is considered in automatic translation and automatic summarization evaluation. However this type of measure has not been intensively used in *ad hoc* retrieval evaluation.

Our goal is to develop a method to evaluate *ad hoc* IR using an informativeness measure to ensure reproductibility in dynamic document collections. To evaluate our method, we compare the system rankings we obtained using the informativeness measure proposed in [6] with the official system rankings based on document relevance, considering various TREC collections on *ad hoc* tasks.

We show that the correlation between the official system rankings and the informativeness measure is specifically high for most of the TREC *ad hoc* tracks.

The rest of the paper is organized as follows. Section 2 presents our evaluation framework which makes use of n-grams for informativeness-based evaluation applied to *ad hoc* retrieval. We also present the *ad hoc* retrieval collections we will be using. Section 3 presents and discusses the comparison of system rankings when using informativeness-based measures with the official ranking for the various *ad hoc* retrieval collections/sub-tasks. Finally, Section 4 concludes this paper.

## 2 N-gram based measure for Adhoc Retrieval

The evaluation method we developed makes an informativeness measure being usable in the case of *ad hoc* retrieval. We use it on various TREC *ad hoc* tracks.

We use the generic Log Similarity (LogSim) informativeness measure initially introduced to evaluate tweet contextualization in 2011 CLEF-INEX QA task [5]. LogSim is based on pools of relevant passages extracted from the document collection (Wikipedia in the CLEF/INEX lab case) called t-rels. t-rels are chunks of texts that are marked-up as relevant by human assessors. By considering each n-gram word in these t-rels as a relevant item, the LogSim normalized measure is based on n-gram precision and graded using log frequencies.

Given a reference  $R$  and a summary  $S$ , the Log Similarity on n-grams ( $LogSim$ ) measure stands as:

$$LogSim(S|R) = \sum_{w \in \mathcal{F}_{St}^n(R) \cap \mathcal{F}_{St}^n(S)} \frac{\log(\min(P(w|S), P(w|R)) \cdot |R| + 1)}{\log(\max(P(w|S), P(w|R)) \cdot |R| + 1)} \cdot P(w|R)$$

where  $P(w|X) = \frac{f_X(w)}{|X|}$  corresponds to the frequency  $f_X(w)$  of n-gram  $w$  in  $X$  over the length  $|X|$  of text  $X$ ,  $\cap \mathcal{F}_{St}^n(S(X))$  is the set of n-grams of stem words from  $X$ , and  $X$  is either  $R$  or  $S$ .

To build such textual references over document ad-hoc q-rels in order to easily apply informativeness to *ad hoc* IR tasks, one approach consists in extracting from documents information nugget candidates; it has been shown that this is possible over non-spammed document collections like TREC robust track or Gov collections [4]. This paper aims at showing that similar results can be obtained without requiring a prior extraction of relevant nuggets. Indeed we propose a direct conversion of relevant documents into a textual reference and experiment plain informativeness measures over it.

For that, we introduce the concept of content interpolated precision at length  $\lambda$  ( $cP_\lambda$ ). Assuming that a user reads the retrieved documents following the ranking given by an IR system,  $cP_\lambda$  evaluates the informativeness of the reading after  $\lambda$  words.

## 2.1 Evaluating Precision based on Document Content

Consider an *ad hoc* task and its document q-rels. Assume that runs are ranked according to Mean Average Precision or Interpolated Precision at several recall levels. Runs can be converted into textual outcomes by concatenating ranked documents and q-rels can be converted into a textual reference by merging together all relevant documents per topic. Runs can be then evaluated by applying informativeness metrics to measure the overlap between submission and reference at various recall levels.

Let  $D = (D_i)_{1 \leq i \leq d}$  be a ranked list of  $d$  documents. We consider as text  $T = (t_1, \dots, t_n)$  the concatenation of these documents (where  $n = \sum_{i=1}^{i=d} |D_i|$ ). For each integer  $\lambda$ , we denote by  $T_\lambda$  the truncated text  $T_\lambda = (t_1, \dots, t_\lambda)$  and by  $D_\lambda^{n,k}$  the set of n-grams with gap  $k$ .  $D_\lambda^{n,0}$  or  $D_\lambda^n$  being the set of n-grams.

Similarly, given a set  $R = \{R_i : 1 \leq i \leq r\}$  of  $r$  relevant documents we shall consider:  $R^{n,k}$  the set of n-grams with gap  $k$  occurring in at least one reference.

In the case of TREC tracks,  $D$  is a run, each  $D_\lambda^n$  is the set of n-grams occurring in one of the  $m$  top ranked documents such that  $\sum_{i=1}^{i=m} |D_i| \leq \lambda$  meanwhile  $R^n$  is the set of n-grams appearing at least once in the relevant documents from the corresponding q-rels.

We apply the English Porter stemming algorithm<sup>4</sup> to all documents after removing all stop words and all document identifiers like TREC doc-ID, etc. This is not only to reduce data, but to convert q-rels into reusable textual relevance judgments (t-rels) than can be applied to non official runs including documents not in the initial collection.

Given a run  $D$  and a reference  $R$ , we define for  $n \in \{1, 2\}$ ,  $k \in \{0, 2\}$  the content interpolated precision  $cP_\lambda^{n,k}$  as:

$$cP_\lambda^{n,k}(D, R) = LS(D_\lambda^{n,k} | R^{n,k}) \quad (2)$$

Observe that, if  $D_j \in R$  then :

$$cP_\lambda^{1,0} \sum_{i=1}^{i=d} |D_i| (D, R) \geq \frac{|D_j|}{\sum_{i=1}^{i=d} |D_i|}$$

<sup>4</sup> <http://snowball.tartarus.org/algorithms/english/stemmer.html>

but conversely,  $D \cap R = \emptyset$  does not imply  $cP_\lambda(D) = 0$  since there can be some overlap between n-grams in documents and in the reference.

So our approach based on document contents instead of document IDs does not require exhaustive references and therefore, can be applied to incomplete references based on pools of relevant documents. However, meanwhile *ad hoc* IR returns a ranked list of documents independently of their respective lengths, relevance judgments can be used to automatically generate text references (*t-rels*) by concatenating the textual content of relevant documents.

## 2.2 Data Sets and Ground Truth

Among international evaluation collections, we chose TREC collections composed of news articles (Robust2004) and Web (Web and Terabyte). Doing so, we also focused on the quality of the collections with large amounts of runs and a comprehensive set of relevance judgments. The number of retrieval systems to rank ranges from 56 to 129, while the number of topics is typically 50 and increases to 150 for Terabyte2006 and 250 for Robust2004 (Table 1). All runs can be downloaded from the TREC web site, and document collections can be obtained on the web site for active participants or through track organizers.

Name	# runs	# topics	Corpus	$ D_n^{1,1} \cup R^{1,1} $
<b>TREC-5</b>	106	50	TREC Vol. 4+5	$15 \times 10^6$
<b>TREC-6</b>	107	50	TREC Vol. 4+5	$12 \times 10^6$
<b>TREC-7</b>	103	50	TREC Vol. 4+5	$28 \times 10^6$
<b>TREC-8</b>	129	50	TREC Vol. 4+5	$35 \times 10^6$
<b>Web2000</b>	104	50	WT10g	$137 \times 10^6$
<b>Web2001</b>	97	50	WT10g	$195 \times 10^6$
<b>Robust2004</b>	110	250	TREC Vol. 4+5	$150 \times 10^6$
<b>Terabyte2004</b>	70	50	GOV2	$46 \times 10^6$
<b>Terabyte2005</b>	58	50	GOV2	$46 \times 10^6$
<b>Terabyte2006</b>	80	150	GOV2	$46 \times 10^6$
<b>Web2010</b>	56	50	ClueWeb09-B	$66 \times 10^6$
<b>Web2011</b>	62	50	ClueWeb09-B	$64 \times 10^6$

Table 1: **Summary of TREC test collections and size in tokens of generated t-rels used for evaluation.**

We take the same experimental approach as in [3] and [2], and reproduced the official rankings of all these retrieval systems for these various collections using the official measure. For all collections, the official measure is the Mean Average Precision (MAP), except for Web2010 and Web2011 where Expected Reciprocal Rank (ERR@20) was preferred. These official rankings constitute the ground truth ranking, against which we will compare the rankings produced by:

- $cP_{10^3}$  based on the 1000 tokens of each run and on their log frequencies.

–  $cP_n$  based on all tokens of each run.

By comparing averaged measures, we evaluate if the average informativeness of all documents retrieved by a system is correlated to the official ranking. Let us emphasize that this does not necessarily imply that document informativeness is correlated to individual document relevance for a given query. We use the Kendall’s  $\tau$  rank correlation coefficient to identify correlations between the ground truth ranking and the informativeness ranking.

### 3 Results and discussion

In this section we report the correlation results of the ground truth ranking (TREC official measure depending on the track) and the content-based ranking produced by the  $cP_\lambda$  informativeness measure. All correlations reported are significantly different from zero with a p-value  $< 0.001$ . While we chose Kendall’s  $\tau$  as the correlation measure, we also report the Pearson’s linear correlation coefficient for convenience. A  $\tau > 0.5$  typically indicates a strong correlation since it implies an agreement between the two measures over more than half of all ordered pairs.

Track	$cP_{10^3}^{1,1}$		$cP_{10^3}^{2,0}$		$cP_{10^3}^{2,2}$		$cP_n^{1,1}$	
	$\tau$	$\rho$	$\tau$	$\rho$	$\tau$	$\rho$	$\tau$	$\rho$
<b>TREC-5</b>	57.91%	71.04%	56.22%	56.30%	56.15%	56.25%	<b>69.91%</b>	<b>88.62%</b>
<b>TREC-6</b>	72.49%	84.18%	76.50%	92.90%	<b>76.50%</b>	<b>93.01%</b>	58.78%	68.18%
<b>TREC-7</b>	61.27%	82.17%	70.18%	90.59%	<b>70.27%</b>	<b>90.60%</b>	63.45%	53.49%
<b>TREC-8</b>	54.80%	84.04%	65.79%	92.25%	65.85%	<b>92.30%</b>	<b>67.46%</b>	72.94%
<b>Web2000</b>	46.48%	68.00%	61.42%	83.40%	62.05%	77.82%	<b>70.83%</b>	<b>86.68%</b>
<b>Web2001</b>	31.65%	56.51%	36.94%	57.89%	36.66%	56.66%	<b>77.45%</b>	<b>88.90%</b>
<b>Robust2004</b>	40.97%	64.33%	58.67%	85.43%	59.50%	86.22%	<b>74.71%</b>	<b>90.88%</b>
<b>Terabyte2004</b>	48.56%	60.12%	61.25%	73.65%	61.45%	74.75%	<b>76.37%</b>	<b>86.12%</b>
<b>Terabyte2005</b>	59.45%	85.34%	69.65%	88.98%	69.80%	<b>89.18%</b>	<b>76.01%</b>	84.28%
<b>Terabyte2006</b>	41.14%	50.24%	54.75%	70.70%	55.28%	70.90%	<b>65.04%</b>	<b>89.37%</b>
<b>Web</b>	28.56%	44.00%	<b>44.38%</b>	<b>69.11%</b>	44.17%	69.09%	-	-
<b>Web2011</b>	56.03%	<b>80.98%</b>	55.68%	79.14%	<b>56.35%</b>	79.42%	<b>34.50%</b>	-

Table 2: Retrieval systems ranking correlations between the official ground truth and the  $cP_\lambda$  informativeness measure.  $cP_\lambda^{1,1}$  stands for uniterms while  $cP_\lambda^{2,2}$  corresponds to bigrams with skip. We use either  $10^3$  terms or all the terms from the ordered list of retrieved documents.

When looking at Table 2, we see that  $cP_\lambda$  accurately reproduces official ranking based on MAP for early TREC tracks (TREC6-7-8, Web2000) as well as for Robust2004, Terabyte2004-5 and Web2011. *LogSim*-score applied to all tokens in runs is often the most effective whenever systems are ranked based on MAP.

However on early TREC tracks,  $cP_{10^3}^{2,2}$  can perform better even though only the first 1000 tokens of each run are considered after concatenating ranked retrieved documents. Indeed, the traditional TREC *ad hoc* and Robust tracks used newspaper articles as document collection. Since a single article often deals with a single subject, relevant concepts are likely to occur together, which might be less the case in web pages for example. A relevant news article is very likely to contain only relevant information, whereas a long web document that deals with several subjects might not be relevant as a whole.

## 4 Conclusion

In this paper, we proposed a framework to evaluate ad hoc IR using the LogSim informativeness measure based on token n-grams. To evaluate this measure, we compared the ranks of the systems we obtained with the official rankings based on document relevance, considering various TREC collections on *ad hoc* tasks. We showed that 1) rankings obtained based on n-gram informativeness and with Mean Average Precision are strongly correlated; and 2) *LogSim* informativeness can be estimated on top ranked documents in a robust way. The advantage of this evaluation framework is that it does not rely on an exhaustive reference and can be used in a changing environment in which new documents occur, and for which relevance has not been assessed. In future work, we will evaluate various *LogSim* parameters influence.

## References

1. C. Cleverdon. The Cranfield tests on index languages devices. 1967.
2. C. Hauff. Predicting the effectiveness of queries and retrieval systems. PhD thesis, Enschede, January 2010. SIKS Dissertation Series No. 2010-05.
3. R. Nuray and F. Can. Automatic ranking of information retrieval systems using data fusion. *Inf. Process. Manage.*, 42(3):595–614, May 2006.
4. V. Pavlu, S. Rajput, P. B. Golbus, and J. A. Aslam. IR System Evaluation Using Nugget-based Test Collections. In Proceedings of WSDM, 2012.
5. E. SanJuan, V. Moriceau, X. Tannier, P. Bellot, and J. Mothe. Overview of the INEX 2011 Question Answering track (QA@INEX). In Focused Retrieval of Content and Structure. Springer, 2012.
6. E. SanJuan, V. Moriceau, X. Tannier, P. Bellot, and J. Mothe. Overview of the INEX 2012 Tweet Contextualization Track. In CLEF, 2012.
7. A. Spink, H. Greisdorf, and J. Bateman. From highly relevant to not relevant: examining different regions of relevance. *Information Processing and Management*, 34(5):599 – 621, 1998.
8. J. Tague-Sutcliffe. Measuring the informativeness of a retrieval process. In Proceedings of the International ACM SIGIR conference on research and development in Information Retrieval, pages 23–36, 1992.