Chapter IR:VIII

VIII. Evaluation

- □ Laboratory Experiments
- Logging
- □ Effectiveness Measures
- □ Efficiency Measures
- □ Training and Testing

Retrieval Tasks

Ad hoc retrieval:

- One question, one result set.
- → Amenable to laboratory environments
- → Canonical measurement of retrieval performance
- Reproducibility and scalability

Interactive retrieval / task-based retrieval:

- The user has a goal or a task that requires many queries and exploration, refining the information need along the way.
- Depends not only on result quality, but also on human factors, context, user interface and experience, and the search engine's supporting facilities.
- Performance is difficult to be measured

Experimental Setup

A laboratory experiment for ad hoc retrieval requires three items:

1. A document collection (corpus)

A "representative" sample of documents from the domain of interest. The sampling method of how documents are drawn from the population determines a corpus's validity. Statistical representativeness may be difficult to achieve, e.g., in case of the web. In that case, the larger a corpus can be built for a given domain, the better.

2. A set of information needs (topics)

Formalized, written descriptions of users' tasks, goals, or gaps of knowledge. Alternatively, descriptions of desired search results. Often accompanied by specific queries the users (would) have used to search for relevant documents.

3. A set of relevance judgments (ground truth)

Pairs of topics and documents, where each documents has been manually assessed with respect to its relevance to its associated topic. Ideally, the judgments are obtained from the same users who supplied the topics, but in practice judgments are collected from third parties. Judgments may be given in binary form, or on a Likert scale.

Every search engine has parameters. Parameter optimization must use an experimental setup (training) different from that used for evaluation (test). Performance deltas between training and test hint at overfitting.

- This setup is sometimes referred to as an experiment under the Cranfield paradigm, in reference to Cleverdon's series of projects at Cranfield University in the 1960s, which first used this evaluation methodology.

 [codalism.com 1] [codalism.com 2]
- In linguistics, a corpus (plural corpora) or text corpus is a large and structured set of texts.

 They are used to do statistical analysis and hypothesis testing, checking occurrences or validating linguistic rules within a specific language territory.

 [Wikipedia]

The term has been adopted in various other branches of human language technologies.

Public Resources

For ad hoc retrieval, the <u>Text Retrieval Conference (TREC)</u> has organized evaluation tracks as of 1992, inviting scientists to compete.

Key document collections used:

Collection	Documents	Size	Words/Doc.	Topics	Words/Query	Jdgmts/Query
CACM	3,204	2.2 MB	64	64	13.0	16
AP	242,918	0.7 GB	474	100	4.3	220
GOV2	25 million	426.0 GB	1073	150	3.1	180
ClueWeb09	1 billion	25.0 TB	304	200	2-6	100-200

- CACM: Communications of the ACM 1958-1979 (only titles and abstracts)
- AP: Associated Press newswire documents 1988–1990
- □ GOV2: Crawl of .gov domains early 2004
- ClueWeb09: Web crawl from 2009

Reusing an existing experimental setup when developing a web search engine allows for direct comparisons with previously evaluated approaches.

- □ TREC is organized by the United States National Institute of Standards and Technology (NIST). It has been key to popularize laboratory evaluation of search engines, organizing evaluation tracks on many different retrieval-related tasks every year: trec.nist.gov.
- Ad hoc retrieval has been studied in the ad hoc tracks, the terabyte tracks and the web tracks.
- Several similar initiatives have formed, namely CLEF, NTCIR, and FIRE.

Topic Descriptions

```
<top>
<num> Number: 794
<title> pet therapy
<desc> Description:
```

How are pets or animals used in therapy for humans and what are the benefits?

```
<narr> Narrative:
```

Relevant documents must include details of how pet or animal-assisted therapy is or has been used. Relevant details include information about pet therapy programs, descriptions of the circumstances in which pet therapy is used, the benefits of this type of therapy, the degree of success of this therapy, and any laws or regulations governing it.

</top>

- ☐ This format has been used for many TREC tracks; it is an invalid form of XML where tags within <top> are not closed, their semantics explicitly repeated in writing.
- <title>'s are typically used as queries.
- More recent tracks resort to well-formed XML:

□ At TREC, every year usually 50 topics are provided. The ones from previous years can be used for training.

Relevance Judgments

```
<top>
<num> Number: 794

<title> pet therapy

<desc> Description:
```

How are pets or animals used in therapy for humans and what are the benefits?

```
<narr> Narrative:
```

Relevant documents must include details of how pet or animal-assisted therapy is or has been used. Relevant details include information about pet therapy programs, descriptions of the circumstances in which pet therapy is used, the benefits of this type of therapy, the degree of success of this therapy, and any laws or regulations governing it.

</top>

Relevance Judgments

<top>

<num> Number: 794

<title> pet therapy

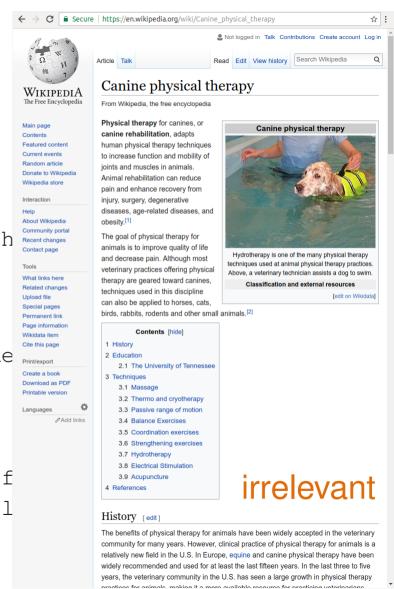
<desc> Description:

How are pets or animals used in th are the benefits?

<narr>> Narrative:

Relevant documents must include de animal-assisted therapy is or has details include information about descriptions of the circumstances used, the benefits of this type of success of this therapy, and any 1 governing it.

</top>



Relevance Judgments

<top>

<num> Number: 794

<title> pet therapy

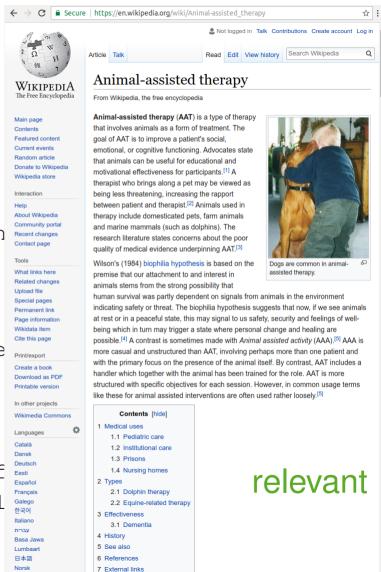
<desc> Description:

How are pets or animals used in th are the benefits?

<narr>> Narrative:

Relevant documents must include de animal-assisted therapy is or has details include information about descriptions of the circumstances used, the benefits of this type of success of this therapy, and any 1 governing it.

</top>



Polski

Relevance Judgments

A relevance judgment requires the manual assessment whether a document returned by a search engine for a given query is relevant under a given topic.

Assessment depth

Assessing every document's relevance for every topic quickly becomes infeasible with growing collection size, and growing numbers of topics and (variants of) search engines to be evaluated. Partial assessments are made based on a sampling strategy called pooling.

Assessment scale

Typically, binary assessments are made, judging documents as relevant or irrelevant. Less often, degrees of relevance are distinguished (e.g., relevant and highly relevant).

Assessor selection and instruction

Ideally, the persons who supplied topics also assess relevance, presuming they genuinely perceived the underlying information need. In practice, this usually cannot be achieved. Hence, a topic's description must be sufficiently exhaustive to serve as instruction.

Assessor agreement

Humans are unreliable judges, their judgments depending on many outside influences. One cannot expect objective truths from just one assessment per document for a given topic. Multi-assessments yield more reliable judgments; assessor agreement can be measured.

□ At TREC, assessors are recruited from retired NIST staff:



IR:VIII-13 Evaluation © HAGEN/POTTHAST/STEIN 2018

Assessment Sampling Strategy: Pooling

Pooling is a sampling strategy for documents retrieved by to-be-evaluated search engines for a given set of topics:

For each topic:

- 1. Collect the top k results returned by each search engine (variant).
- 2. Merge the results, omitting duplicates, obtaining a "pool" of documents.
- 3. Present the pool of documents in random order to assessors.

Observations:

- Presuming a certain correlation between search engines' results on a topic, only documents considered relevant by search engines need be analyzed.
- New retrieval algorithms that are evaluated later may retrieve unjudged documents, requiring new judgments, probably from different assessors.
- All documents ranked below the threshold are deemed irrelevant by default, regardless the truth.

Measuring Annotator Agreement

Annotator agreement is the degree to which annotators agree with each other in a series of decisions. If no agreement can be measured, the scale may be inappropriate, or the annotators need to be retrained.

Annotator agreement is measured in cases where ambiguous or subjective decisions have to be made. Relevance is a subjective notion.

Measuring Annotator Agreement: Kappa Statistics

Given the judgments of two annotators on a given topic, a kappa statistic measures their agreement as follows:

$$\kappa = \frac{p_o - p_e}{1 - p_e},$$

where p_o denotes the proportion of agreement observed, and p_e the expected proportion of agreement by chance.

Properties:

- \Box 1 p_e denotes the agreement attainable above chance

- \Box $\kappa = 0$ indicates random agreement
- $\alpha < 0$ indicates disagreement worse than chance

Measuring Annotator Agreement: Kappa Statistics

Given the judgments of two annotators on a given topic, a kappa statistic measures their agreement as follows:

$$\kappa = \frac{p_o - p_e}{1 - p_e},$$

where p_o denotes the proportion of agreement observed, and p_e the expected proportion of agreement by chance.

Suppose A and B are two annotators asked to make n relevance judgments. Then a simple kappa statistic can be computed as follows:

		В		\sum	$p_o = \frac{a+e}{}$
		yes	no		n
Λ	yes	a	b	c	$p_e = P(extsf{yes})^2 + P(extsf{no})^2$
^	no	d	$rac{b}{e}$	f	
\sum		g	h	n	$P(\mathbf{yes}) = \frac{c+g}{2n}, P(\mathbf{no}) = \frac{f+g}{2n}$

Measuring Annotator Agreement: Kappa Statistics

Given the judgments of two annotators on a given topic, a kappa statistic measures their agreement as follows:

$$\kappa = \frac{p_o - p_e}{1 - p_e},$$

where p_o denotes the proportion of agreement observed, and p_e the expected proportion of agreement by chance.

Suppose A and B are two annotators asked to make n=400 relevance judgments. Then a simple kappa statistic can be computed as follows:

		В	\sum	
		yes	no	
Α	yes	300	20	320
	no	10	70	80
\sum		310	90	400

$$p_o = rac{300 + 70}{400}$$
 $p_e = P(extsf{yes})^2 + P(extsf{no})^2$ $P(extsf{yes}) = rac{320 + 310}{2 \cdot 400}, \quad P(extsf{no}) = rac{80 + 90}{2 \cdot 400}$

$$\kappa = 0.776$$

- \Box Well-known kappa statistics include Cohen's κ , Scott's π , and Fleiss' κ .
- \Box Scott's π is the one exemplified.
- \Box Fleiss' κ is a generalization of Scott's π to arbitrary numbers of annotators and categories. It also does not presume that all cases have been annotated by the same group of people.
- Presuming that raters A and B work independently, the probability $P(yes)^2$ ($P(no)^2$) denotes the probability of both voting yes (no) by chance. Another way of computing p_e is to sum the mutplication of the rater-specific probabilities of each rater voting yes (no).
- \square Some assign the following interpretations to κ values measured (disputed):

κ	Agreement
<0	poor
0.01-0.20	slight
0.21-0.40	fair
0.41-0.60	moderate
0.61-0.80	substantial
0.81-1.00	almost perfect

[Wikipedia]

figure Within TREC evaluations, typically a "substantial" agreement ($\kappa \approx [0.67, 0.8]$) is achieved.

[Manning 2008]