

Statistical methods used in search engine evaluations

Erzsébet Tóth^a, Gábor Fazekas^b

^aCollege of Nyíregyháza
Institute of Mathematics and Informatics
tothe@zeus.nyf.hu

^bUniversity of Debrecen, Institute of Informatics
fazekasg@inf.unideb.hu

Abstract

This paper provides a brief overview of the different statistical methods and approaches utilized in search engine evaluations. Six significant research projects are described with their fulfilled activities and statistical results.

Introduction

In evaluations researchers measure the quality of Web search services that can be determined on the basis of several measures. However, it is extremely difficult to find reliable measures that reflect this quality appropriately. Oppenheim [6] suggested that benchmarking tests should include the following measures at the very minimum:

1. precision;
2. relative recall;
3. speed of response;
4. consistency of results over an extended period;
5. proportion of dead or out of date links;
6. proportion of duplicate hits;
7. overall quality of results as rated by users;
8. evaluation of GUI for user friendliness;
9. helpfulness of help and variations of software for new and experienced users;
10. options for display of results;
11. presence of adverts;
12. coverage;
13. estimated/expected search length;

14. length and readability of abstracts;
15. search engine effectiveness.

The lack of a standard set of measures causes a great problem in evaluations. Because of this deficiency research on search engine evaluation is inconsistent in methods and approaches. So there is a real need for working out a standard set of measures to evaluate search engines properly. From a statistical point of view it would be also interesting to examine if there is any relationship among these measures.

Evaluations have been carried out mostly on robot engines, but in principle any of the search engines can be measured. In most cases the results of evaluations remain valid only for a short period of time and indicate only the performance of search engines at that time. Although several difficulties associated with search engine evaluations, we have to make an effort to measure search engines currently in use. However, so far standardized evaluation approaches have not been applied for this purpose. In general experiments report their own idiosyncratic approaches and they mostly avoid the use of standardized evaluation methods.

1. Research projects focusing on search engine evaluation

1.1. Sroka [10] has evaluated the performance of Polish versions of English language search engines and homegrown Polish search engines (Polish AltaVista, Polish Infoseek, Virtual Poland, NEToskop, Onet.pl, WOW). In his evaluation precision was emphasized that he determined on the basis of topical relevance judgements. However, he did not consider the authoritativeness of retrieved web pages. He analysed the first 10 hits retrieved from each search engines with this method. He also studied the overlap of the retrieved results and the response time of each search engine. The number of retrieved hits from each search engine was recorded, but he omitted recall as a relevancy criterion. He formulated 10 queries about various topics, four of them required Boolean logic (AND, AND NOT). He conducted searches with and without Polish diacritical marks. He concluded that he retrieved the largest number of documents for each query by using diacritical marks. The average number of relevant documents out of the first ten hits was calculated for each search engine. He determined a mean precision score for each query to find which queries were the most difficult for search engines to handle.

1.2. Clarke and Willett [2] have carried out 30 searches on three different search engines, such as AltaVista, Excite and Lycos. The relevance of the first 10 retrieved hits for each query was determined on the basis of a three-point scale: a score of 1 was assigned to relevant documents, 0.5 to partially relevant documents and 0 to non-relevant hits. Mean values for precision, recall and coverage were calculated. The significance of the differences in performance between the three search engines were evaluated by using the Friedman two-way analysis of variance test. Leighton and Srivastava used this test for the analysis of search-engine results.

The purpose of the applied statistical method is to test the null hypothesis that k different samples (corresponding to the three search engines) have been produced from populations having the same median. The data can be seen in a table consisting of N rows (here rows are equivalent to the 27 queries) and k columns. The data in each row are ranked from 1 to k . The Friedman statistic, F_r , is built on the sum of the ranks for each search engine. [8]

1.3. *Leighton and Srivastava* [3] compared the performance of five search engines, such as AltaVista, Excite, Infoseek, Hotbot and Lycos. They constructed a test suite containing 15 queries that were submitted to search engines. They measured the precision of the first 20 results by taking into account the percentage of relevant hits within the first 20 retrieved. The relevance assessment of the hits was based on six different relevance categories. They conducted five experiments for the first 20 precision.

In the calculation of precision a weighing factor was used to increase value for ranking effectiveness. Precision and ranking effectiveness were combined into one metric. This metric incorporates several qualities: first, we have to take into consideration that a link either fulfils the relevance criteria under examination, or it does not. A binary scale of relevance was applied, because the relevance categories were different definitions of relevance. Second, more weight is given to effective ranking of relevant documents. Third, the statistic has to reflect the fact that if the search service retrieves fewer results with the same number of good results, it is easier for the user to find relevant links. At last they had to decide how to handle inactive and duplicate links. In the first three tests duplicates were only deleted from the numerator of the precision ratio, search services were penalized. In the last two tests duplicates were not penalized. In all five experiments the retrieved inactive links were penalized. However, they did not consider two other types of duplicates, such as mirror pages and clusters of pages.

A formula calculating the performance of the service on a query is between zero and one. The first 20 hits for each query have been grouped by their status and type of relevance. We begin the formula for this metric with converting the relevance categories into binary values of zero or one. We assign a value to each position on a 20 position linear scale of value in order to weigh ranking. On this scale the first position represents the greatest value and the last position represents the smallest value. In the formula the first 20 hits are divided into three groups. In each group the links receive an equal weight, the weight that the first link in the group would have obtained in a 20 position linear scale of value. The first three links have a weight of 20, the next seven have a weight of 17 and the last ten have a weight of 10. They added up these weighted values to produce the numerator of the metric.

For example, if a service retrieved five good links, it would calculate the score of $(3*20)+(2*17)=60+34=94$ if these hits were the first five ranks.

If these hits were all between ranks 11 and 20, it would score only $(5*10)=50$.

They calculated the denominator of the metric from the number (up to 20) of links retrieved by the search service. If the service retrieved 20 or more links, then

the sum of all weights to 20 was applied.

$$(3*20)+(7*17)+(10*10)=279.$$

The denominator is altered if we have fewer than 20 links returned. If the denominator were not altered, it would always be 279. If there are no links returned, the denominator would be zero, with the metric undefined. Because of this boundary condition, they calculated the denominator by adding up all of the weights to 20, 279, and then subtracting 10 for each link less than 20 retrieved.

For example, if a service retrieves 15 links, the denominator is $279-(5*10)=229$, but if it retrieves one link, the denominator is $279-(19*10)=89$.

Then they divided the numerator by the denominator to calculate the final metric. They used a Rube Goldberg machine as a function that could be described with this complete formula:

$$\frac{(\text{Links } 1-3 * 20) + (\text{Links } 4-10 * 17) + (\text{Links } 11-20 * 10)}{279 - [(20 - \min(\text{No. of links retrieved}, 20)) * 10]}.$$

In the tests the Friedmann's randomized block design was utilized, because the normality assumption needed for the ANOVA model was not fulfilled. In the Friedman test the blocks were the queries and the treatment effects were the search services. The Friedman test analyses population medians rather than mean values because of the skewness that is present. In all five experiments they refused the null hypothesis that the search service's medians could be equal. Because they refused null hypotheses, they could make pairwise multiple comparisons between individual services.

They have conducted a correspondence analysis of the queries by search services. For this purpose they used scores from experiment two as weights. They could examine how a search service or a query corresponded to the composite score of the whole by using a correspondence analysis. They displayed all of the information in the correspondence relationship by presenting services and queries on a graph in higher dimensional space.

They concluded that Alta Vista, Excite and Infoseek provided more relevant hits than Hotbot and Lycos. The first 20 hits for the top three services included all of the words from the search expression more frequently than the first 20 hits for the lower two services. This metric formula could be tested against other weighing and scoring schemes, such as the traditional Coefficient of Ranking Effectiveness. [5] In the future a study should be made where the test suite is large enough to compare structured search expressions versus unstructured ones.

1.4. *Chignell, Gwizdka and Bodner* [1] have carried out two experiments for studying the performance of commercial search engines, such as Excite, Hotbot and Infoseek. In the first experiment they analysed the effect of time of day and the effect of query strategy on query processing time for each of three search engines. They used nine prespecified queries that could be divided into three different categories: general, higher precision and higher recall. Document relevance was measured by using a 'consensus peer review' procedure. So they chose six other search engines on which the same queries were conducted. They obtained binary

judgment of relevance from the hits of six different search engines (AltaVista, Lycos, Northern Light, Search.com, Web Crawler, Yahoo) to the same query. A hit from one of the search engines was considered to be relevant if it was also retrieved by at least one of the six referee search engines in response to the query. The usage of this method is questioned in relevance assessment, because there is a minimal overlap between the set of hits of the search engines. They also measured the number of broken and duplicate links for each search.

Multivariate analysis of variance (MANOVA) was applied for the analysis of the data. They found a multivariate effect for the two-way interaction of Query Type and Search Engines ($F(20, 604.6)=10.6$, $p < 0.001$). A significant univariate interaction for precision caused this effect ($F(4, 186)=6.9$, $p < 0.001$). The univariate interactions for the other three dependent variables (corresponding to the three search engines) were not important.

They realized a significant main effect of search engine query time ($F(2, 186) = 65.5$, $p < 0.001$). A significant difference was shown between the query processing times of Excite and Infoseek ($p < 0.001$), and Hotbot and Infoseek ($p < 0.001$) by post hoc Tukey testing.

They found a significant main effect of search engine on the number of broken links, on the number of duplicate links and on precision. There was a significant main effect of time of day on query processing time. They realized a significant main effect of query type on the precision scores. The three search engines executed best the general queries that were followed by high precision and high recall queries. At this stage of the experiment they have analysed only the effects of the independent variables on a single dependent variable.

We may define a user-oriented composite measure of performance on the basis of four dependent variables. They received the ranking of the three search engines for each dependent variable by performing post hoc Tukey tests. A simple formula was devised by using ranking information: the number of first, second and third place rankings were summed for each search engine. The first ranks were multiplied by a coefficient of 3, the second ranks by a coefficient of 2, and the third rank by a coefficient of 1.

For instance: the composite measure of performance for Excite is:

$$3*(1/4)+2*(2/4)+(1/4)=2/3= 66.7\%$$

This measure has two boundaries. First, it does not take into account the case where there is no statistical significant difference between two search engines. Second, all dependent variables are treated as being equal.

The second experiment analysed the influence of geographical coverage and Internet domains on search engine performance. They used three search engines (Altavista, HotBot, Infoseek), six Internet domains and four queries in a fully factorial design of the 72 observations. The four queries were translated to three different languages. They measured the precision of the first 20 hits on the basis of human relevance judgments.

Full factorial multivariate analysis was done, in which search services and domain names were used as independent factors, with 14 dependent measures. They found a significant multivariate interaction between search engines and domains ($F(221, 425.24)=1.56$, ($p<0.001$)). Interaction between search engines and domains had significant univariate effects on the number of unique hits, on total number of hits, on the proportion of retrieved hits to each search engine collection size, on the quality of returned results, and it also had a borderline significant effect on search length 1. The domain name and the search engine had a significant multivariate effect separately. They conducted univariate analyses to determine the source of these effects. The search engine had a significant univariate effect on the following measures: differential objective precision, best and full precisions, and search length 1.

1.5. Analyses of search engine transaction logs [9] have shown that the average users formulate term-based queries including approximately two terms and they sometimes use operators. In contrast search experts apply more search terms and advanced search operators than the average searchers. Lucas and Topi [4] conducted a comprehensive study to examine the influence of query operators and term selection on the relevancy of search results. 87 participants involved in the survey formulated queries on eight topics that were used on a search engine of their choice. Besides this search experts constructed queries on each of the topics that were submitted to the eight preferred search engines of participants (AltaVista, AOL, Excite, Go, Google, iWon, Lycos, Yahoo!).

All of the queries were executed and analysed during a 1-day period. A cutoff value of 10 was used in judging the relevancy of pages, because relevant links appearing on the first page of search hits were most likely to be viewed. The relevancy of the first 10 web pages was judged on the basis of a four-category ordinal scale for relevancy. The relevancy criteria associated with each of the relevancy scores for the given search topic were determined independently.

The most important results of the study were associated with the research model and examined by using two different regression analyses, such as full multiple regression and step-wise hierarchical multiple regression. Average standardized relevancy was a dependent variable in the regression model that indicated the search engine performance. In the research model seven variables were connected to operator usage, and four variables were connected to term usage. All of them were treated as independent variables. If we take together these variables they will explain 31,8% of the variance in the dependent variable. We consider this result both statistically significant and a relevant amount of variance. The first-order Pearson correlations between the research variables were presented in a matrix.

It was found that search term selection and usage are more important than the selection and usage of operators. The two independent variables closely related to dependent variable are term variables. One of them assesses the number of terms (e.g. taking into account the absolute difference between the number of terms in a subject's query and a corresponding expert query). Another assesses the number of subject terms that match with terms in the corresponding expert query. If we

take together these two variables they will explain more than 25% of the variance in the dependent variable. The number of misspelled terms is also important which forecasts the search performance.

We find only one operator variable among the four most significant independent variables, which is namely the number of other nonsupported operators. It actually assesses the number of NOT, OR and (-) operators in contexts where their usage is not supported. However, the usage of nonsupported ANDs and nonsupported (+) operators is positively linked to search performance.

1.6. Radev, Libner and Fan [7] examined how successful the most popular search engines are at finding accurate answers to natural language questions. Altogether 700 queries were submitted to nine different search engines. They downloaded and stored the top 40 documents retrieved by each search engine. It was established that all search engines returned at least one correct answer on more than three-quarters of the factual natural language questions.

In the experiment they tested three hypotheses that were the following:

1. Search engines are effective at answering factual natural language questions.
2. Certain characteristics of questions forecast the likelihood of retrieving a correct answer across all search engines.

3. Questions with particular characteristics are more likely to draw the correct answer from specific search engines.

A score was assigned to each of the search engines as the sum of the reciprocal ranks of documents containing the correct answer. Then they calculated the mean score across all queries for each search engine to evaluate the first hypothesis mentioned above.

To evaluate hypotheses two and three, they coded the 700 queries on the following four factors:

1. type of answer required
2. presence of proper noun
3. time dependency of the query
4. number of words in query

Corresponding to the three hypotheses, they planned to use an analysis of variance (ANOVA) (1) to compare the general score of the nine search engines, (2) find out the importance of the above four factors in predicting score, and (3) measure differential performance of search engines on each of these factors.

The initial distribution of scores showed a positive skew, because there was a large proportion of zero-value scores. This skew would not fulfil the normality assumption of ANOVA. To solve this problem a two-step analysis was carried out.

1. Scores were converted to values of zero or nonzero. After that a binary logistic regression was conducted. A nonzero score was selected, when the search engine retrieved at least one correct answer in the top 40 documents for a given question. A zero score was chosen, when there was no correct answer retrieved. This analysis looks for a relationship between the four question characteristics and whether a correct answer is retrieved at all.

2. The second part of the analysis primarily dealt with nonzero values - cases

where at least one correct answer was returned. The distribution of this restricted dataset still had some positive skew, so square-root transformation was applied to it. Then an analysis of variance (ANOVA) was conducted.

On the basis of the findings they stated that all the search engines managed to return the correct answer in the top 40 documents 75% of the time or more.

Conclusion

Research on search engine evaluation is inconsistent in applied methods and approaches, for this reason there is a growing need for a set of benchmarking tests for search engines. In addition to this a standard set of measures should be worked out for monitoring the performance of search engines. This overview can serve as a guideline for choosing an appropriate statistical method for use in search engine evaluations. In the design of our experiment we have to decide which statistical method would correspond to our research purposes and create a statistical model with the appropriate variables

References

- [1] Chignell, M. H. - Gwizdka, J. - Bodner, R. C.: Discriminating meta-search: a framework for evaluation. In: *Information Processing and Management* vol. 35. 1999. p. 337-362.
- [2] Clarke, J. - Willett, P. : Estimating the recall performance of Web search engines. In: *Aslib Proceedings* vol. 49. no. 7. July/August 1997. p. 184-189.
- [3] Leighton, H. V. - Srivastava, J.: First 20 precision among World Wide Web search services (search engines). In: *Journal of the American Society for Information Science*, vol. 50. no. 10. 2000. p. 870-881.
- [4] Lucas, W. - Topi, H.: Form and function: the impact of query term and operator usage on web search results. In: *Journal of the American Society for Information Science* vol. 53. no. 2. 2002. p. 95-108.
- [5] Noreault, T. - Koll, M. - McGill, M.J.: Automatic ranked output from Boolean searches in SIRE. In: *Journal of the American Society for Information Science*, vol. 28. 1977. p. 333-339.
- [6] Oppenheim, C.- Morris, A. - McKnight, C. - Lowley, S.: The evaluation of WWW search engines. In: *Journal of Documentation* vol. 56, no. 2, March 2000, p. 190-211.
- [7] Radev, D. R. - Libner, K. - Fan, W.: Getting answers to natural language questions on the Web. In: *Journal of the American Society for Information Science* vol. 53. no. 5. 2002. p. 359-364.
- [8] Siegal, S. - Castellan, N. J.: *Nonparametric statistics for the Behavioral Sciences*. McGraw-Hill: Singapore, 1988.
- [9] Silverstein, C. - Henzinger, M. - Marais, J. - Moricz, M.: Analysis of a very large web search engine query log. In: *SIGIR Forum*, vol. 33. no. 1. 1999. p. 6-12.
- [10] Sroka, M.: Web search engines for Polish information retrieval: questions of search capabilities and retrieval performance. In: *International Information & Library Review* vol. 32. 2000. p. 87-98.