

Evaluation of Information Retrieval Performance of Turkish Search Engines

SUMMARY

This is an investigation on the information retrieval performances of search engines based on various measures. We searched 17 queries of differing types on four Turkish search engines, namely Arabul, Arama, Netbul and Superonline. We classified each document/Web site contained in the retrieval results as being “relevant” or “non-relevant”. Based on this classification, we calculated the precision and normalized ranking ratios in various cut-off points for each query run on each search engine. We checked the “dead” or “broken” links among the retrieval results to determine how often the crawlers of search engines visit the sites they index and how often they update their indexes, if needed. We found out the coverage and novelty ratios of each search engine by searching five keywords that have been the most frequently submitted queries to the Turkish search engines. Those keywords are “mp3”, “oyun” (game), “sex”, “erotik” (erotica) and “porno” (porn). By means of two modest experiments, we tested to see if Turkish search engines make use of index terms that are assigned by the authors of Web pages and included under the “keywords” and “description” meta tags of HTML documents. Using Kruskal-Wallis and Mann-Whitney statistics, we tested if up-to-dateness, precision, normalized ranking, coverage and novelty ratios of each search engine differ significantly from each other.

Major findings of our research are as follows: On the average, one in six documents retrieved by search engines was not available due to dead or broken links. Netbul retrieved fewer documents with dead or broken links than other search engines did. Some search engines retrieved no documents (so called “zero retrievals”) or no relevant documents for some queries. On the average, five in six documents retrieved were not relevant. Average precision ratios of search engines ranged between 11% (Netbul) and 28% (Arama) (Superonline being 20% and Arabul 15%). Arama retrieved more relevant documents than that of Arabul and Netbul in the first five documents retrieved. Search engines do not seem to make every efforts to retrieve and display the relevant documents in higher ranks of retrieval results. Average normalized ranking ratios of search engines ranged between 20% (Arabul) and 54% (Arama) (Superonline being 37% and Netbul 30%). Arama retrieved the relevant documents in higher ranks than that of Arabul and Netbul. The strong positive correlation between the precision and normalized ranking ratios got weakened as the number of documents that we evaluated increased. Search engines were less successful in finding relevant documents for specific queries or queries that contained broad terms. Although non-relevant documents were higher in number, search engines were more successful in single-term queries or queries with Boolean “OR” operator. The success rate was lower for queries with Boolean “AND” operator. Search engines seemingly do not use stemming algorithms to better analyze queries and to increase retrieval performance. The use of Turkish characters such as “ç”, “ö”, and “ş” in queries still creates problems for Turkish search engines as retrieval results differed for such queries. Superonline’s coverage rate was much higher than that of other search engines for the most frequently searched queries on the Turkish search engines. Except Arama, search engines index fewer documents/sites with domain names ending with “.tr”. Arama is the indisputable leader in covering documents with Turkish addresses. Almost all search engines scored high in novelty ratios for the most frequently searched queries. Different search engines tend to retrieve different relevant documents for

the same queries. For retrieval purposes, Netbul and Superonline seem to index and make use of metadata fields that are contained in HTML documents under “keywords” and “description” meta tags.

The research report concludes with some recommendations to improve the information retrieval performances of Turkish search engines.