

Web Search Engine Comparison

The exercise is about comparing the search results from Google versus Bing, the two leading US search engines. Many search engine comparison studies have been done. All of them use samples of data, some small and some large, so no general conclusions can be drawn. But it is always instructive to see how the two search engines match up, even on a small data set.

The process you will follow is to issue a set of queries and to evaluate the returned results for relevance. These studies do not seek to answer the ultimate question of which search engine is “best”. Rather we stick to more modest research questions which are:

RQ1: Which search engine performs best when considering the first five results for a given query?

RQ2: Is there a difference in relevance between the search engines when considering informational queries and navigational queries, respectively?

To begin the class is divided across the set of Schools at USC. Students are pre-assigned according to their USC ID number, as given in the table below.

USC ID ends with	School to crawl	Root URL
01~20	Dornsife (College)	http://dornsife.usc.edu/
21~40	Gould (Law)	http://gould.usc.edu/
41~60	Keck (Medicine)	http://keck.usc.edu/
61~70	Marshall (Business)	http://marshall.usc.edu/
71~80	Viterbi (Engineering)	http://viterbi.usc.edu/
81~00	Price (Public Policy)	http://priceschool.usc.edu/

Now that you have been assigned a USC School, below are the queries you will submit. There are a total of nine (3+3+1+1+1) navigational queries, three informational queries, and one final query.

Input Navigational Queries: Devise a set of queries for your USC School as follows;

- Choose 3 Faculty names from your school and enter the following query using the names from your school, e.g. “Ellis Horowitz Viterbi” or “David Cruz Gould” or “Tara Blanc Price” (do NOT use quotes in your query; include only the faculty name and the school name. Your query should be exactly as shown above.)

Determine relevance (see below for how to determine relevance) for each individual faculty name; do not average over the three names;

- **Choose 3 Faculty departments**, e.g. “Computer Science Viterbi”, or if there is no department use a division name, e.g. “Director of Admissions, Gould”. If there are no departments or divisions, come up with a suitable categorization on your own. Your query should be exactly as shown above, that is it includes ONLY the department or division name followed by the school name.

Determine relevance for each individual department name; do not average over the three names;

- **Determine School Location**, a map, e.g. “Viterbi USC map” or “Price USC map”. Your query should be exactly as shown, the school name, USC followed by the word “map”.
- **Determine the Founder**: The USC School of Engineering is named after Andrew Viterbi, the USC School of Business is named for Gordon S. Marshall; the USC School of Public Policy is named for Sol Price, etc. Issue a query to find a web page describing the individual who has named the school, e.g. “Andrew Viterbi”, “Gordon Marshall”, “Sol Price”; the web page can be a USC page, or if not, a Wikipedia entry. **Your query should contain ONLY the name of the founder of the school.**
- **Determine School Alumni News** web page, e.g. “USC Viterbi Alumni” or “USC Gould Alumni”. Your query should only contain “USC” followed by the name of the school, followed by the word “Alumni”.

Input Informational Queries: Devise a set of queries for your USC School as follows

- **Requirements for an undergraduate degree in a given department** or if there are no departments than simply the requirements for an undergraduate degree, e.g. “USC Computer Science Undergraduate degree requirements”
 - **Requirements for a Masters degree in a given department** or if there are no departments than simply the requirements for a Masters degree , e.g. “USC Computer Science Masters degree requirements”
 - **Requirements for a Ph.D. degree** in a given department or if there are no departments than simply the requirements for a Ph.D. degree or whatever the most advanced degree that is offered , e.g. “USC Computer Science Ph.D. degree requirements”
- If your School does not offer an undergraduate, Masters, or Ph.D. degree, devise a query for whatever degree(s) are offered.

Query 13: Attempt to create a query for your USC school that Google includes in its top five results, but Bing does not include in its top five results.

Note: do not alter the above queries so more relevant results are returned; use only the queries as specified above since they are typical of what a casual user might enter.

Place your queries in a separate text file which you will submit with the rest of your assignment.

Each of your thirteen queries should be run on both Google and Bing. You should capture the top five results (the URL) for each query. For each of the top 5 results for each query you should compute a relevance score as follows:

For faculty names relevance = 1 for a search result to the faculty's home page¹; relevance = 0.5 for course page taught by the faculty member, and relevance = 0.25 for a page with only a little information about the faculty member, and otherwise relevance = 0;

For faculty departments or divisions relevance = 1 for a search result to the department's home page, relevance = 0.5 for an page that is internal to the department and otherwise relevance = 0;

For school location, relevance = 1 for a search result containing map and/or directions, otherwise relevance = 0; note that a Google map that provides the exact building location is as relevant as a USC campus map.

For school founder's name relevance = 1 for a search result that describes the individual, relevance = 0.5 for a page that gives the history of the school and mentions the individual, and otherwise relevance = 0;

For alumni news web page relevance = 1 for a result that points to an alumni news page; if one exists and is not returned, then relevance is 0. A returned page that talks about the school's alumni get a relevance of 1. A page describing a specific alum gets a relevance of 0.25.

For the informational queries relevance = 1 if the page describes the requirements, relevance = 0.5 if it contains a link to the actual requirements, and otherwise relevance = 0.

Note: in the event that your Google account enables personalized search, please turn this off before performing your tests.

Output

Once you score all of the search results for all of the queries you should produce the following statistics.

1. A text file containing the list of queries that you used and for each query the top five URLs produced as results, and for each URL the relevance score that you assigned. The data should include the results for both Google and Bing.

2. An Excel or Google docs spreadsheet showing the following:

2.1 For each query a bar graph with Y-axis from 0 to 1 and X-axis the top five results; the value for each result is the relevance score for Google and Bing; so your bar graph should have ten bars

2.2 A bar graph whose Y-axis is 0 to 5 and whose X-axis is query 1, query 2, . . . , query 5 and whose value for each query is the number of overlapping search results for that query. Results are assumed to overlap if the identical link is contained in the top 5 results².

¹ Notes on special cases: a professor may have more than one home page, perhaps one created by him and one created by his department; both may receive a relevance score of 1; to receive a relevance score of 1, the homepage must have a usc.edu domain; links to external sites such as a LinkedIn entry for a professor is not considered a home page, though it can be recorded with relevance 0.5; a resume or CV is not considered a home page, but may get relevance = 0.25

² If Google and Bing show different URLs, but they point to the identical page, this should be considered as an overlap

2.3 For the nine navigational queries, a bar graph showing the ratio of relevant vs. irrelevant pages. A page is considered relevant if its relevance score is greater than 0. A page is considered irrelevant if its score is 0. Each column of the bar graph represents the ratio of the number of relevant results for Google (Bing) in the top five divided by the number of irrelevant results for Google (Bing) in the top five results. For example, see pp. 7, figure 1 of the Lewandowski paper

2.4 For the three informational queries, a bar graph showing the ratio of relevant vs. irrelevant pages in the top five results. See pp 8, Figure 2 of the Lewandowski paper

Note1: you should use a spreadsheet program to produce the above graphs.

Note2: place all of your results on a single sheet of the spreadsheet

Note3: do not reformulate your queries in such a way that the search engine produces more relevant results; the point of the exercise is to examine the results when a “normal” query (as defined above) is entered

References

Evaluating the retrieval effectiveness of Web search engines using a representative query sample by Dirk Lewandowski, Hamburg University of Applied Sciences, Journal of the American Society for Information Science and Technology, 2013

<http://arxiv.org/ftp/arxiv/papers/1405/1405.2210.pdf>

Submission

You are required to submit your results electronically to the csci572 account on SCF so that it can be graded. To submit your file electronically, enter the following command from your Unix prompt:

```
submit -user csci572 -tag hw1 MYFILE1 MYFILE2
```

where MYFILE1 contains your queries and MYFILE2 contains your results.