

How to Improve Your Search Engine Ranking: Myths and Reality

AO-JAN SU, Northwestern University
Y. CHARLIE HU, Purdue University
ALEKSANDAR KUZMANOVIC, Northwestern University
CHENG-KOK KOH, Purdue University

Search engines have greatly influenced the way people access information on the Internet, as such engines provide the preferred entry point to billions of pages on the Web. Therefore, highly ranked Web pages generally have higher visibility to people and pushing the ranking higher has become the top priority for Web masters. As a matter of fact, Search Engine Optimization (SEO) has become a sizeable business that attempts to improve their clients' ranking. Still, the lack of ways to validate SEO's methods has created numerous myths and fallacies associated with ranking algorithms.

In this article, we focus on two ranking algorithms, Google's and Bing's, and design, implement, and evaluate a ranking system to systematically validate assumptions others have made about these popular ranking algorithms. We demonstrate that linear learning models, coupled with a recursive partitioning ranking scheme, are capable of predicting ranking results with high accuracy. As an example, we manage to correctly predict 7 out of the top 10 pages for 78% of evaluated keywords. Moreover, for content-only ranking, our system can correctly predict 9 or more pages out of the top 10 ones for 77% of search terms. We show how our ranking system can be used to reveal the relative importance of ranking features in a search engine's ranking function, provide guidelines for SEOs and Web masters to optimize their Web pages, validate or disprove new ranking features, and evaluate search engine ranking results for possible ranking bias.

Categories and Subject Descriptors: H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval

General Terms: Algorithms, Design, Measurement

Additional Key Words and Phrases: Search engine, ranking algorithm, learning, search engine optimization

ACM Reference Format:

Su, A.-J., Hu, Y. C., Kuzmanovic, A., and Koh, C.-K. 2014. How to improve your search engine ranking: Myths and reality. *ACM Trans. Web* 8, 2, Article 8 (March 2014), 25 pages.
DOI: <http://dx.doi.org/10.1145/2579990>

1. INTRODUCTION

Search engines have become generic knowledge retrieval platforms used by millions of Internet users on a daily basis. As such, they have become important vehicles that drive users towards Web pages highly ranked by them (e.g., [Cho and Roy 2004; Moran and Hunt 2005]). Consequently, finding ways to improve ranking at popular search engines is an important goal of all Web sites that care about attracting clients.

This project is supported by the National Science Foundation (NSF) via grant CNS-1064595.

The subject of this work appears in *Proceedings of IEEE/ACM International Conference on Web Intelligence 2010* [Su et al. 2010].

Authors' addresses: A.-J. Su and A. Kuzmanovic, Electrical Engineering and Computer Science Department, Northwestern University, 2145 Sheridan Road, Evanston, IL 60208; Y. C. Hu (corresponding author) and C.-K. Koh, School of Electrical and Computer Engineering, Purdue University, West Lafayette, IN 47906; email: ychu@purdue.edu.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2014 ACM 1559-1131/2014/03-ART8 \$15.00

DOI: <http://dx.doi.org/10.1145/2579990>

Ways to improve a Web page's search engine ranking are different. On one side, SPAM farms are a well-known approach to attempt to boost a Web site's ranking. This is achieved by artificially inflating a site's popularity, that is, by increasing the number of nepotistic links [Davison 2000] pointing to it. Luckily, ways to detect and contain such approaches appear to be quite successful [Benczúr et al. 2005; Gyongyi et al. 2006; Wu and Davison 2005]. On the other side, the entire industry of Search Engine Optimization (SEO) is booming. Such companies (e.g., [TOPSEOs 2014; SeoPros 2014]) and experts claim to be capable of improving a Web page's rank by understanding which page design choices and factors are valued by the ranking algorithms.

Unfortunately, the lack of any knowledge or independent validation of the SEO methodologies, and the ever-lasting interest on this topic, has opened the doors to various theories and claims, myths and folklore about which particular factor is influential, such as Aubuchon [2010], SEOMoz [2007], Kontopoulos [2007], Patel [2006], Li [2008], AccuraCast [2007], and Marshall [2009]. To the best of our knowledge, none of these claims is backed by any published scientific evidence. At the same time, the problem of predicting a search engine's overall ranking results is widely considered a close-to-impossible task due to its inherent complexity. It should also be noted that this article presents an academic, not a commercial study. Thus we have no preference for, or bias towards, any commercial search engine.

The key contribution of our article is that we demonstrate that simple *linear* learning models, accompanied by a recursive partitioning ranking scheme, are capable of predicting a search engine's ranking results with high accuracy. As an example, in Google's case, we show that when non-page-content factors are isolated our ranking system manages to correctly predict 8 pages within the top 10 for 92% of a set of randomly selected 60 keywords. In the more general scenarios, we manage to correctly predict 7 or more pages within the top 10 for 78% of the set of keywords searched.

In this work, we start by describing several initial unsuccessful attempts at predicting a search engine's ranking results, to illustrate that reproducing such ranking results is not a straightforward task. In our initial attempts, we set up cloned and artificial Web sites trying to decouple ranking factors in a search engine's ranking algorithm. These attempts failed because we could not obtain sufficient data points due to limited pages being indexed and infrequent visits of a search engine bot.

We then revised our approach to consider multiple ranking factors together. We developed an automated ranking system that directly queries a search engine, collects search results, and feeds the results to its ranking engine for learning. Our ranking engine incorporates several learning algorithms, based on training both linear and polynomial models. Using our ranking system, we show that a linear model trained with linear programming and accompanied with a recursive partitioning algorithm is able to closely approximate a search engine's ranking algorithm. In addition, we use our ranking system to analyze the importance of different ranking features to provide guidelines for SEOs and Web masters to improve a search engine's ranking of Web pages. Furthermore, we present case studies on how our ranking system can facilitate in validating and disproving potential new ranking features reported in the Internet. More specifically, we confirm that the particular search engine (Google) imposes negative bias toward blogs, and that HTML syntax errors have little to no impact on the search engine's ranking.

To study the potentially different ranking algorithms used in different search engines, we also use our ranking system to analyze the ranking results of the Bing search engine. Our experimental results show that the recursive partitioning algorithm is also beneficial for improving the ranking system's prediction accuracy for Bing. In addition, we compare the relative importance of ranking features of the two search engines' ranking algorithms and show the disagreements between the two.

In particular, our case studies for Bing show that the search engine favors incoming links' quality over quantity and Bing shows no bias toward the overall traffic of a Web site.

In addition to the study we presented in Su et al. [2010], this article adds the following new contributions: in Section 2.3, we describe our unsuccessful initial attempts in learning a search engine's ranking system. We expect that our lessons learned will be valuable for others who will attempt to explore similar problems. In Section 6, we present results for using our ranking system to analyze the Bing search engine and compare them with our findings in Su et al. [2010]. Furthermore, we provide insights into the ranking algorithm of the Bing search engine and present new case studies for Bing. Thus, while in this work we necessarily limit our analysis to two search engines, namely Google and Bing, our approach is far more generic than reverse engineering the two search engines.

The rest of the article is structured as follows. In Section 2, we define the problem and outline the folklore around a search engine's ranking system. In Section 3, we present the detailed design of our ranking system. We present the evaluation results of using our ranking system to analyze the Google search engine in Section 4 and several case studies in Section 5. We present the evaluation results of using our ranking system to analyze the Bing search engine and associated case studies in Section 6. We discuss related issues in Section 7. Finally, we conclude in Section 8.

2. PROBLEM STATEMENT

2.1. Goals

There have been numerous efforts that attempt to reveal the importance of ranking factors to a search engine [Aubuchon 2010; Kontopoulos 2007; Patel 2006; SEOMoz 2007]. While some of them are guess-works by Web masters [Kontopoulos 2007; Patel 2006], others are based on experience of Search Engine Optimization (SEO) experts [Aubuchon 2010; SEOMoz 2007]. While we recognize that guessing and experience might indeed be vehicles for revealing search engine internals, we strive for more systematic and scientific avenues to achieve this task. The goal of our study is to understand the important factors that affect the ranking of a Web page as viewed by popular search engines. In doing so, we validate some folklore and popular beliefs advertised by Web masters and the SEO industry.

2.1.1. Ranking Features. In this work, we aim to study the relative importance of Web page features that potentially affect the ranking of a Web page, as listed in Table I. For each Web page (URL), we collect 17 ranking features. These ranking features can further be divided into 7 groups. The page group represents characteristics associated with the Web page, including Google's Page Rank score (PR) and the age of the Web page in a search engine's index (AGE). The URL group represents features associated with the URL of the Web page. Parameter HOST counts the number of occurrences of the keyword that appear in the hostname and PATH counts the number of occurrences of the keyword in the page segment of the URL.

The domain group consists of features related to the domain of a Web site. D_SIZE reports the number of Web pages indexed by Google in the domain and D_AGE reports the age of the first page index by archive.org in the domain. Groups header, body, heading, and link are features extracted from the content of the Web page. TITLE counts the number of occurrences of the keyword in the title tag. M_KEY counts the number of occurrences of the keyword in the metakeyword tag and M_DES counts the number of occurrences of the keyword in the meta-description tag. DENS is the keyword density of a Web page, which is calculated as the number of occurrences of the keyword divided by the number of words in the Web page. H1 through H5 is the

Table I. Ranking Features

Group	Feature	Detail
Page	PR	pagerank score
	AGE	age of the web page in a search engine's index
URL	HOST	keyword appear in hostname
	PATH	keyword in the path segment of url
Domain	D_SIZE	size of the web site's domain
	D_AGE	age of the web site's domain
Header	TITLE	keyword in the title tag of HTML header
	M_KEY	keyword in meta-keyword tag
	M_DES	keyword in meta-description tag
Body	DENS	keyword density
Heading	H1	keyword in h1 tag
	H2	keyword in h2 tag
	H3	keyword in h3 tag
	H4	keyword in h4 tag
	H5	keyword in h5 tag
Link	ANCH	keyword in anchor text
	IMG	keyword in image tag

number of occurrences of the keyword in all the headings H1 to H5, respectively. ANCH counts the number of occurrences of the keyword in the anchor text of an outgoing link and IMG counts the number of occurrences of the keyword in an image tag.

Google claims to use more than 200 parameters in its ranking system. Necessarily, we explore only a subset of all possible features. Still, we demonstrate that ranking features listed in Table I are adequate for providing high ranking prediction accuracy. Moreover, we are capable of establishing important relationships among the explored features.

2.2. State-of-the-Affairs (the Folklore)

The great popularity and the impact that search engines have in shaping users' browsing behavior have created significant interests and attempts to understand how their ranking algorithms work. Still, there is no consensus on the set of the most important features. Different people express quite different opinions, as we illustrate shortly.

To get a glimpse of how little agreement there is about the relative importance of ranking features, we collect different opinions for a search engine's ranking features and summarize in Table II. We acknowledge this is not a comprehensive survey; rather it is a collection of anecdotal evidence from several diverse sources. In particular, the second column labeled by *SEOmoz'07* [SEOmoz 2007] is a list of top 10 ranking factors created by surveying 37 SEO experts by SEOmoz.org in 2007. This column represents observations from knowledgeable experts. The third column labeled by *Survey* [Kontopoulos 2007] is a list of top 10 ranking features rated by a poll of Internet users interested in this topic. This column represents the perception of the search engines' ranking algorithm from general Internet users. While the Internet users may not be page ranking experts, they are the ultimate end-users who tune their Web pages according to their beliefs and thus who are affected the most. Finally, the fourth column labeled as *Idv* [Marshall 2009] is the top 10 ranking feature list posted by an Internet marketing expert on his personal Web page. This column represents an individual investigator that studies this topic. We note the SEOmoz and Survey opinions are for search ranking in general, and the LDV opinion is about Google search.

Due to the lack of systematic measuring and evaluating guidelines, it is not surprising to see significant differences in ranking between the three lists. The SEO experts

Table II. Various Ranking Feature Opinions

#	SEOMoz'07	Survey	Idv
1	Keyword use in title tag	Keywords in title	Keyword in URL
2	Anchor text of inbound link	Keywords in domain name	Keyword in domain name
3	Global link popularity of site	Anchor text of inbound links	Keyword in title tag
4	Age of site	Keywords in heading tags	Keyword in H1, H2 and H3
5	Link popularity within the site's internal link structure	Keywords in URL	Page Rank
6	Topical relevance of inbound links to site	Anchor text from within the site	Anchor text of inbound link to you
7	Link popularity of site in topical community	Internal links	Site listed in DMOZ Directory
8	Keyword use in body text	Keywords in Alt attribute of images	Site listed in Yahoo Directory
9	Global link popularity of linking site	Relevance of external links	Rank Manipulation by Competitor Attack
10	Topical relationship of linking page	Keywords in body	Site Age

obviously favor ranking features associated with hyperlinks as they rated 7 out of the top 10 ranking features in this category. On the contrary, the other two opinions have only 3 and 1 top 10 ranking features associated with links, respectively. In addition, the Web site's age feature is in the top 4 features among SEOs, but is absent in the second list and is at the bottom of the third list.

From time to time, Internet users will see rumors that spread about a search engine's ranking algorithm. However, to the best of our knowledge, there does not exist a systematic approach to validate or disprove these assumptions. This motivates us to perform research in this topic and build a system to facilitate the necessary evaluation process.

2.3. Initial Attempts

Our path towards achieving the preceding goal was not straightforward. In this section, we present our initial unsuccessful attempts in approaching this problem. We hope that the lessons learned here could save time and hence be useful for others who attempt to explore this problem.

Reverse engineering a search engine's ranking algorithm of its organic search engine result pages is a difficult task because it is a multi-variable optimization problem in a highly uncertain environment. Hence, our initial attempts were focused on designing several experiments that aim to *isolate* a subset of ranking features and evaluate their relative importance in isolation. We describe these attempts next.

2.3.1. Cloned Web Sites. Our first attempt was trying to decouple Web content from other ranking factors such as page rank score, domain reputation, and page freshness. Our technique involves cloning Web sites from two reputable domains (slashdot.org and coffeegeek.com) to our new established domain. We retain each Web item in the original Web site except hosting the content in our own Web server and change the domain name in each URL. We submit our new Web sites to Google via its Web-master tool [Google 2014c]. In addition, we inform Google bot of every Web page in the new Web sites by generating a sitemap for Google bot to crawl. However, this attempt could not live up to our expectations due to the following reason: Google bot crawled and

indexed only a very limited number of pages in our cloned Web site. Our understanding is that the Google bot conserves its resources in this way. In addition, our cloned Web site might be detected as a SPAM Web site by Google and thus prevent Google's crawler to visit our Web site again. This left us with very few data points to evaluate the search engine's content score ranking algorithm.

2.3.2. Artificial Web Sites. In another attempt to isolate the impact of the features, we use our isolated domain and set up artificial Web sites with fabricated Web pages. Our idea is that even if we cannot make a search engine index a large number of pages in our Web sites, we can modify the content of the Web pages and still get more data points. We carefully design our Web pages to contain different numbers of targeted keywords and we place these keywords in different HTML tags. Similarly, we submit these new Web sites to Google and wait for Google bot to visit. Once the Google bot crawls and indexes a Web page, we measure and record the ranking result and change the content of the indexed page in an attempt to obtain more data points. However, after several weeks of measurements, we still could not get enough data points to analyze the results. This is because the frequency of Google bot's visits was very low (once per week) and the number of indexed pages was small (less than 20 per domain).

2.3.3. Summary. In our initial attempts, we tried to set up new Web sites (cloned or fabricated) in order to evaluate a search engine's ranking algorithm by decoupling the ranking factors. These attempts failed because the search engine bot indexed very limited number of Web pages for newly established Web sites. In addition, the frequency of the search engine bot's visit is also too low, which makes the turn-around time of content change to ranking change very long. In short, it may take a very long time to collect sufficient data points to evaluate a search engine's ranking algorithm by the preceding methods. Nonetheless, we find that our original idea of separating different ranking factors to be quite useful. We demonstrate that it is possible to apply this approach using advanced searching querying methods, such as those available in the Google API, later in Section 5.1.

3. METHODOLOGY

In this section, we discuss the design of our ranking system that analyzes a search engine's ranking algorithm.

3.1. Design Goals

There are several design goals that such a ranking system should meet.

- The system should be automated. In particular, the data collection (from the Web) and the offline analysis should all be automated.
- The system should output human-readable results. In particular, the output of the system should give intuitive explanations to why some Web pages are ranked higher than others, and provide guidelines for how to improve the ranking of Web pages.
- The system should be able to handle a large amount of data, as the training set and the training of the ranking model should converge in a reasonable amount of time.
- The ranking system should be extensible. For instance, it should be able to accommodate new ranking features when they are available.
- The system should rely only on publicly available information and it should avoid any intrusive operations to the targeted search engine.

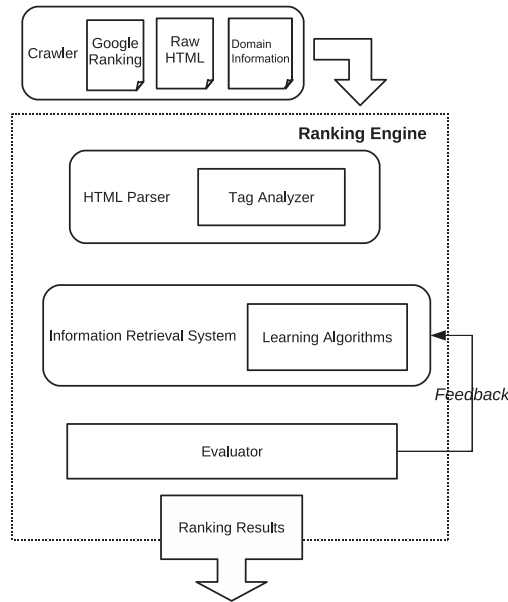


Fig. 1. System architecture.

3.2. A Practical Approach

We have designed and implemented a ranking system that meets the aforesaid design goals. The architecture of our system is depicted in Figure 1. The two major components are the crawler and the ranking engine.

Our initial unsuccessful attempts (discussed in Section 2.3) suggest that it is difficult and impractical to isolate the impact of individual features that may affect the ranking, for example, one at a time. Therefore, we resort to considering the impact of multiple features together. This boils down to issuing search queries to a search engine and collecting and analyzing the search results. The data collection is performed by the *crawler* which queries a search engine and receives the ranked search results. In addition, it downloads HTML Web pages from their original Web sites and queries domain information as described in Section 3.3.

Second, since multiple features can affect the ranking of Web pages in complicated ways, the ranking engine extracts features under study from raw Web pages and performs learning to train several ranking models to approximate the ranking results by a search engine. In this part, we make several contributions: (1) We confirm that a search engine's ranking function is not a simple linear function of all the features, by showing a nonlinear model can outperform, that is, approximates a search engine's ranking better than, a simple linear model. However, a nonlinear model is difficult for humans to digest. (2) We present a simple recursive ranking procedure based on a simple linear model and show that it can achieve comparable accuracies to the nonlinear model. The theoretical underpinning for such a procedure is that recursive application of a linear model (function) can effectively approximate a nonlinear function. In addition, the linear model converges more efficiently and outputs more human-readable results.

3.3. The Crawler

The crawler submits queries to a search engine and obtains top 100 Web pages (URL) for each keyword. Without losing generality, we limited our queries to HTML files to

avoid Web pages in different file format such as PDFs and DOCs that could create unnecessary complications in our experiments. In addition, we focus on Web pages composed in English in our experiments. For example, the Google API syntax we use for the previous two features is `as_filetype:html&lr=lang_en`. Moreover, to obtain the date that a search engine indexed the Web pages, we submit our queries with an additional parameter `qdr:y10`. By doing so, the date that a search engine indexed the Web page will be returned in the search result page for us to extract the age of the page ranking feature. Finally, for each Web page, the crawler does the following.

- (1) It downloads the Web page from the original Web site.
- (2) In case of the Google search engine, it queries the URL's page rank score by Google toolbar's API [Google 2014a].
- (3) It obtains the age of a page (the latest date a search engine crawled the Web page) by parsing the search result page.
- (4) It obtains the size (the total number of pages) of the domain by querying a search engine with `site: [domain]`.
- (5) It queries `archive.org` and fetches the age of the Web site (the date when the first Web page was created on this Web site).

3.4. The Ranking Engine

There are three components in the ranking engine. The HTML parser [PHP HTML Parser 2014] converts Web pages into the Document Object Model (DOM) for the tag analyzer to examine the number of keywords that appear in different HTML tags such as anchor text. The ranking engine trains the ranking model by combining features obtained from the Web page contents, page rank scores, and domain information. After the model is created, the ranking engine evaluates the testing sets by applying the model. The evaluator then analyzes the results and provides feedback to the ranking engine which is used to adjust parameters in the learning algorithms such as error threshold. In the following sections, we describe the two ranking models we experimented in this article: linear programming and SVM. We use ranking features listed in Table I to train our ranking models.

3.4.1. Linear Programming Ranking Model. In this section, we describe our linear programming ranking model. Given a set of documents $I = (i_1, i_2, \dots, i_n)$, predefined search engine ranking $G = (1, 2, \dots, n)$, and a ranking algorithm A , the goal is to find a set of weights $W = (w_1, w_2, \dots, w_m)$ that makes the ranking algorithm reproduce the search engine ranking with minimum errors. The objective function of the linear programming algorithm attempts to minimize errors (the sum of penalties) of the ranking of a document set. Eq. (1) defines the objective function which is a pairwise comparison between two documents in a given dataset.

$$\Phi(W) = \sum_{i=1}^n \sum_{j=i+1}^n c_i \cdot |i - j| \cdot D(i, j). \quad (1)$$

In Eq. (1), c_i is a factor that weights the importance of the i_{th} document (e.g., a top 5th page is more important than a top 50th page). $|i - j|$ is the distance (ranking difference) between the i_{th} and the j_{th} page. Finally, $D(i, j)$ is a decision function we define as

$$D(i, j) = \begin{cases} 0 & \text{if } f(A, W, i) \geq f(A, W, j), \\ 1 & \text{if } f(A, W, i) < f(A, W, j). \end{cases} \quad (2)$$

where $f(A, W, i)$ is the score produced by algorithm A with a set of weights W for the i_{th} page in the given dataset. Page X is ranked higher than page Y if it receives a higher

score than page Y. The decision function denotes that if the ranking of the two pages preserves the order as the search engine's ranking, the penalty is zero. Otherwise, the penalty will be counted in the error function which is denoted in Eq. (1).

Since we cannot import conditional functions (e.g., $D(i, j)$) into a linear programming solver, we transform the decision function into the following form

$$f(A, W, i) + D_{ij}F_{\max} \geq f(A, W, j), \quad (3)$$

where F_{\max} is the maximum value to which $f(A, W, \cdot)$ would evaluate, and $D_{ij} \in \{0, 1\}$. When $D_{ij} = 0$, the preceding inequality is satisfied only if

$$f(A, W, i) \geq f(A, W, j). \quad (4)$$

When $D_{ij} = 1$, the inequality is always satisfied. Therefore, we have effectively converted the original minimization problem into the problem of minimizing the D_{ij} . Hence, we can now replace $D(i, j)$ in the objective function with D_{ij} . Finally, the score function $f(A, W, i)$ can be represented by a dot-product of ranking parameters $X = (x_1, x_2, \dots, x_m)$ and weights $W = (w_1, w_2, \dots, w_m)$ denoted as $f(A, W, i) = f_A(w_i \cdot x_i)$. For example, the parameter x_i can be the number of keywords that occur in the title tag and w_i is the weight associated with x_i .

In addition to the objective function, we set constraints to our linear programming model. For each pair of pages (i, j) where $i < j$ (i is ranked higher than j by the search engine), we have a constraint

$$f(A, W, j) - f(A, W, i) \leq \tau, \quad (5)$$

where τ is the maximum allowed error which is set to a predefined constant τ . The constant τ is adjusted by the feedback from the evaluator to refine the ranking results. For example, when linear programming solver cannot find a feasible solution, we relax the maximum allowed error τ . The linear programming solver we use in our experiments is ILOG CPLEX [CPLEX 2014]. In addition, we apply a recursive partitioning algorithm as we describe in Section 3.4.3 shortly.

3.4.2. Support Vector Machines' Ranking Model. Support Vector Machines (SVMs) are a set of supervised learning methods used for classification, regression, and learning ranking functions [Vapnik 2000]. In an SVM, data points are viewed as n-dimensional vectors (n equals to the number of ranking features in our case). An SVM constructs a hyperplane or a set of hyperplanes in a high-dimensional space, which is used as a classifier to separate data points. In our experiments, we use the SVM-rank [Joachims 1999, 2009] implementation with linear and polynomial kernels to train the ranking functions. The ranking features we used in our SVM experiments are the same as the linear programming model as discussed in Section 2. The parameter c in SVM-rank controls the trade-off between training error and margin. The ranking engine adjusts the value of parameter c according to the feedback provided by the evaluator in order to find the best value for prediction accuracy. Finally, we perform a recursive partitioning algorithm as we describe in the following section.

3.4.3. Recursive Partitioning Ranking Algorithm. It is common that a search engine keeps several layers of indices in practice. For example, the first layer of indices may serve as a cache and it is able to answer queries for top 20 pages. When the first-layer query fails, it is then sent to subsequent indices. Additionally, the search engine's internal ranking algorithm can be nonlinear. To capture such a nonlinear and/or nonequational behavior for a search engine's ranking function, we developed a recursive partitioning algorithm to approximate this ranking behavior. We apply this algorithm to both our linear and SVM models. We evaluate the power of this recursive partitioning ranking

algorithm in Section 4.3. First, we describe our recursive partitioning algorithm with pseudocode shown here.

ALGORITHM 1: Recursive partitioning ranking algorithm

```

1: procedure PARTITION( $S, X$ )           // S: a set of pages, X: top X
2:   Rank( $S$ )                           // Train or apply ranking models
3:   while  $|S| > 2.5 * X$  do
4:      $N = \text{Max}(2 * X, \frac{|S|}{2})$ 
5:      $S \leftarrow \text{Top}(S, N)$          // Return top N pages
6:     return Partition( $S, X$ )
7:   end while
8:   return Top( $S, X$ )                 // Return top X pages
9: end procedure

```

In Algorithm 1, S denotes a set of pages in a dataset and X denotes the target top X pages to be evaluated (e.g., top 10 pages). Algorithm 1 can be explained by giving an example of how to train recursive ranking models for selecting top 10 pages ($X = 10$) out of 100 Web pages ($|S| = 100$). In the first round of recursion, the learning algorithm produces a set of weights by training all 100 pages (in line no. 2 of Algorithm 1). Next, line no. 4 calculates $N = 50$. In line no. 5, the function returns the top 50 pages out of 100 using the model learned previously. Next, the algorithm moves on to the second recursion with a set of 50 pages in S . Similarly, in the second recursion, a new set of weights for ranking is learned and the variable N in line no. 4 becomes 25. The partitioning algorithm further extracts top 25 pages from the 50 pages (in line no. 5) and proceeds to the third round. In the third round, an additional new set of weights is learned in line no. 2 of Algorithm 1. The condition statement in line no. 3 is not met in this round. Therefore, the algorithm escapes the recursive while-loop. Algorithm 1 then proceeds to line no. 8 and returns the top 10 pages by applying the ranking model learned in the third round.

The process of evaluating the testing sets using our recursive partitioning ranking algorithm is similar to the steps we described previously. The input of this evaluation process contains a set of pages to evaluate S , a variable X , and ranking models learned in the previous training procedure. For example, to evaluate top 10 pages out of 100 pages in a dataset, the recursive partitioning algorithm first obtains top 50 pages using the weights learned in the first round of the training process. The top 50 pages are sent to the second round and the algorithm extracts top 25 pages using the set of weights learned in the second round. Finally, the third round evaluates top 10 pages out of the 25 pages using the third set of weights learned in the third round of the training process.

4. EVALUATION

In this section, we evaluate the accuracy of our ranking system in predicting Google's ranking results. We first explain our experimental methodology. Next, we evaluate our system's overall ranking accuracy. Then, we evaluate the importance of the recursive partitioning ranking algorithm and demonstrate its effectiveness in dealing with the underlying nonlinearities. Finally, we evaluate the relative importance of various Web page features.

4.1. Experimental Setup

Using our crawler, we collect search results from Google for 60 keywords in four categories, shown in Table III. The four categories are Linux commands, chemical

Table III. Query Keywords

Type	Keywords
Linux commands	tcpdump, modprobe, egrep, chmod, dhclient, dmesg, netstat, nslookup, traceroute, rsync, crontab, iconv, telnet, iptables, xtail
Chemicals	potassium, boron, chlorine, manganese, lithium, magnesium, phosphorus, sulfur, fluorine, iodine, helium, zinc, platinum, cobalt, uranium,
Music Terms	adagio, arpeggio, baroque, cadence, crescendo, diminuendo, fortissimo, legato, moderato, pianissimo, pizzicato, ritardando, rubato, sforzando, staccato,
Astronomy Terms	aphelion, apogee, chromosphere, ecliptic, equinox, photon, pulsar, supernova, zodiac, galaxy, polaris, neutron, perihelion, magnetosphere, perigee.

elements, as well as music and astronomy terms. We select these terms and the corresponding keywords to keep our experiments simple and at the same time not lose generality. In particular, the goal is to avoid plural, similar words that Google might take into account, etc. Later, in Section 5.1, we demonstrate that our approach is applicable to other popular keywords as well. In this article, we focus on search terms that consist of a single word. Search engines could have proprietary processes to preprocess multiwords query terms and is out of scope of this work.

In each experiment, we randomly select 15 keywords to form the *training set*. We then run these keywords through our ranking system and develop a ranking model. Then, we use the remaining 45 keywords, which we term the *testing set*, to evaluate the accuracy of our model.

4.2. Overall Ranking Accuracy

Here, we evaluate the effectiveness of our ranking system by comparing its ranking results to those of Google. In the experiments, we use linear programming (Section 3.4.1) as well as the two SVM learning algorithms (Section 3.4.2), that is, linear and polynomial. In all cases, we apply three rounds of recursion (Section 3.4.3). While predicting the top 10 pages is not the same as ranking pages correctly, we compute the fraction of correctly returned top 10 pages (also known as “recall @ 10”) as the measure of ranking accuracy. We select this measure simply because from the SEO perspective, predicting if a page will reach the top 10 is what really matters. Indeed, statistics have shown that most Internet search engine users (as many as 90%) never click past the first page of search results [Hopkins 2012], and tend to only consider the first few search results [Joachims et al. 2005].

Figure 2 shows the ranking results of our system under the LP model and the two SVM models, using the ranking features listed in Table I as input variables. In this figure, the x-axis represents the hit rate of the top 10 pages, that is, the percent of the top 10 pages ranked by Google that are also captured in the top 10 pages output by our system. The y-axis shows the percent of pages, distributed over all keywords from the test set, that satisfy the given hit rate. For example, the point (60,100) in the figure for the SVM polynomial algorithm means that for all 45 keywords from the test set, our algorithm managed to always (100% of time) correctly predict at least 6 pages out of the top ten ones reported by Google (hence, 60%). In general, the closer the curve is to the upper right corner, the better the result.

Figure 2 shows that the linear programming model achieves accuracy comparable to SVM-linear, and has only slightly lower accuracy compared to the SVM-polynomial model. In particular, for the LP model, 78% of the keywords experienced a hit rate greater than 70%, and 54% of the keywords experienced a hit rate greater than 80%.

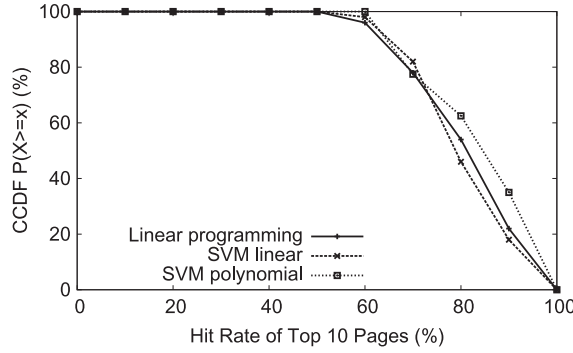


Fig. 2. Comparison of ranking models.

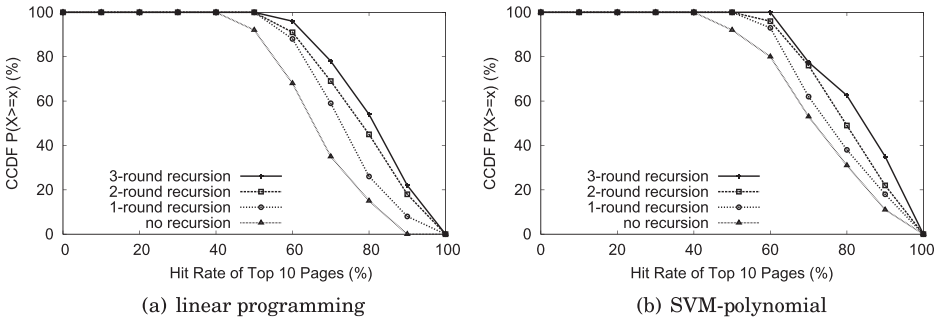


Fig. 3. Power of recursive partitioning.

Despite the slight lag behind the polynomial model, the linear model is much more practical and convenient because it provides human-readable insights (feature weights in our case). Hence, we use linear programming models in the rest of our experiments. An exception is the next section, where our goal is to understand the role of recursive partitioning ranking. Hence, we explore this issue in the context of the nonlinear algorithm as well.

4.3. The Power of Recursive Partitioning Ranking

Several factors potentially introduce nonlinear effects in Google’s ranking results. First, it is more than likely that Google uses nonlinear functions in their ranking algorithms. Second, it is preferable for a search engine to keep multiple indices (e.g., one for the most accessed top 20 pages, another one for pages 20–40, etc.), instead of having a huge index. These indices can be ranked independently and then merged together. Hence, this can add another level of nonlinearity that can further complicate the ranking results prediction problem.

Our goal here is to understand the role of the recursive partitioning and its ability to “smooth out” these nonlinearities. Hence, we compare the performance of recursive and nonrecursive ranking approaches. We study both linear and polynomial models. Our hypothesis is that recursion should be much more effective in the linear case, because the polynomial model should be able to track nonlinearities more effectively. We validate this hypothesis next.

Figure 3 compares the results of the LP model and the SVM-polynomial model without recursive partitioning and with 1, 2, and 3 rounds of recursive partitioning. We make the following observations: (1) Without recursive partitioning ranking, the LP

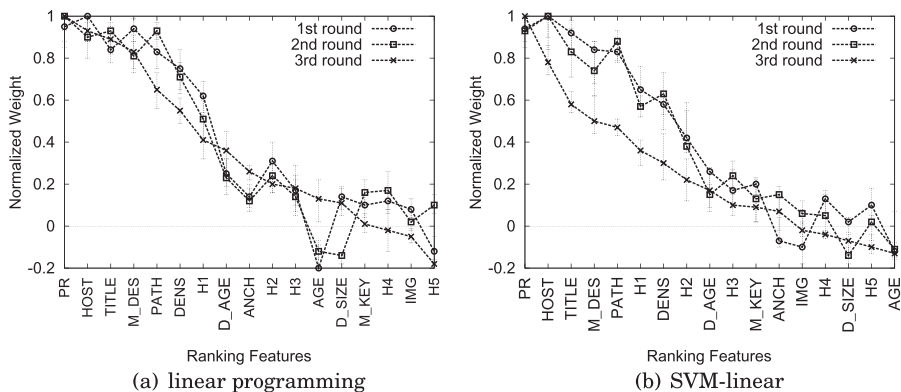


Fig. 4. Weights in different rounds.

model achieves lower prediction accuracy than the SVM-polynomial model. For example, the former predicted 8 out of the top 10 pages by Google for 18% of the keywords, while the latter predicted 8 out of the top 10 pages by Google for 32% of the keywords. (2) Recursive partitioning significantly improves the ranking accuracy for both models. The improvement reaches diminishing returns after 3 rounds. The average hit rate after 3 rounds of recursion improves by 20% over the no-recursion algorithm. The recursive partitioning algorithm successfully removes noise contributors (the lower half) of the dataset and the ranking on the remaining data is improved. (3) After 3 rounds of recursive partitioning, the LP model achieves similar accuracy as the SVM-polynomial model. Hence, it is sufficient to analyze the linear model which gives human-readable outputs, as the weights of the linear model directly reflect the relative importance of different features. We next investigate the weights of the features in the linear model.

4.4. Relative Importance of Features

We next analyze the relative importance of different ranking features towards contributing to the overall ranking of a page. Since in the LP model the ranking function is simply a linear combination of all the ranking features, their relative importance boils down to the relative values of the weights in the linear function.

Figure 4 shows the weights and the 90% confidence intervals (i.e., each coefficient falls into its corresponding bar range in 90% of the trials) of the linear equation from the LP and SVM ranking models, sorted in decreasing order, for the three recursion rounds. The first insight is that the dominant features (i.e., the first 7 on the x-axis) carry larger weights in the first two recursion rounds than in the last. This effect is particularly pronounced in Figure 4(b), for the SVM-linear model. This suggests that these features are dominant in “pushing” Web pages to the top 50 or 20 pages, yet in order to get into the top 10 ones, other factors (including those that are in general less valued) become relatively more important as well.

We observe that, despite some disagreement in the two linear models, the feature weights in the two models are highly correlated. In particular, the first 5 features are in the same order for both algorithms. Not surprisingly, page rank is the dominant factor in both cases. Nevertheless, HOST (keyword appearing in the hostname), TITLE (keyword appearing in the title tag of the HTML header), M_DES (keyword appearing in the meta-description tag), and PATH (keyword appearing in the path segment of the URL) are the other leading factors, respectively.

Figure 4 also shows a high correlation among the factors that have little to no impact (the bottom of the x-axis). In particular, AGE (the age of the Web page in a search engine's index) is at the bottom of the list; on the contrary, D_AGE, the *domain* age, is an important parameter valued by the ranking system. This is not a surprise since the domain age in general increases its reputation. Continuing with the factors at the bottom of the list, both algorithms show that M_KEY (keyword in the metakeyword tag) has low impact; on the contrary, M_DES (keyword in the meta-description tag) has a higher impact, as we explained earlier. Further, D_SIZE (domain size) and IMG (keyword in the image tag) do not have much influence.

5. CASE STUDIES

Here, we perform several case studies. In particular, we first focus on isolating the content score and explore the general effects of dealing with a smaller number of ranking features. Next, we evaluate the ability of our approach to add new ranking features in a methodical way; we explore whether a given category, such as blogs, is a factor considered by the search engine ranking algorithm. Finally, we explore whether HTML syntax errors affect ranking results or not.

5.1. Isolating Subsets of Ranking Features

In this case study, we explore how well our ranking system works, that is, approximates Google, when focusing on a *subset* of Web pages for which a subset of parameters are known to matter. Such a study allows us to understand the strengths and weaknesses of our approach in ranking different types of Web content. In particular, we focus on all but page features shown in Table I. The set of features include the features that belong to the URL, domain, header, body, heading, and link groups. To conduct the study, we need to decouple the subset of page ranking features from others. To exclude the impact of the page features, we crawl only “young” Web pages (i.e., past 24 hours) using Google search API `tbs=qdr:d`. By looking only at “young” Web pages, we manage to effectively remove the age factor, since all pages crawled are of the same age. Moreover, we manage to remove the page rank factor, because it takes much more than one day for a Web page to obtain its page rank score, as we also evaluate experimentally.

In order to find sufficiently many new Web pages (at least 100 per keyword) that are generated in the most recent 24 hours, we necessarily abandon the keywords shown in Table III. This is simply because not enough new Web pages that contain these keywords (and are indexed by Google) appear daily. Hence, we turn our attention to more popular keywords. In particular, we query the popular Google Trends Web site [Google 2014b] and obtain a list of 60 popular keywords for our experiment. We then query the aforesaid advanced Google API (`tbs=qdr:d`), which gives us only new Web pages, using the 60 keywords. If a given keyword cannot return 100 or more links, we abandon such a keyword. Finally, we randomly select 15 keywords as the training set and the remaining 45 as the testing set.

Figure 5 shows the ranking results using the LP model and three iterations of recursive partitioning and ranking. The x-axis shows our ranking system's hit rate of top 10 pages returned by Google. We observe that the hit rate is about 80% for 92% of keywords, and about 90% for 77% of keywords. These results are much better than those in Figure 2, suggesting that our ranking system is more accurate, that is, approximates Google's ranking results much better, when only content-related features affect the ranking. This case study suggests when the parameters are more specific, that is, fewer and closely related, our ranking engine performs better.

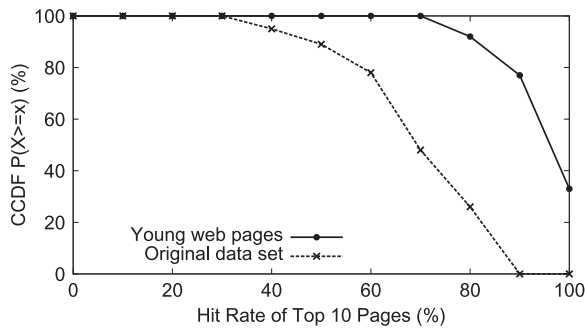


Fig. 5. Ranking results for “young” pages.

For horizontal comparison, we apply the same subset of ranking features (i.e., without PR and AGE) to the original dataset we evaluated in Section 4. We plot the ranking results for the younger pages and for the original dataset in Figure 5. The figure clearly shows that the accuracy of our ranking system is much worse after removing the PR and AGE ranking features when ranking the original dataset. This shows PageRank is an important ranking factor, when it is available, as we have discussed in Section 4.4. But when it is not available, which is the case for young pages, the PageRank ranking feature does not affect the ranking accuracy.

5.2. Negative Bias toward Blogs

In this case study, we show that our methodology can be used to study if there exists any positive or negative bias towards certain categories or types of Web pages. Uncovering the existence of such bias has clear implications for Web site designers who wish to improve their Web site ranking.

In this particular case study, we introduce a new ranking feature that represents a Web page *category*, for example, news, music, blogs, etc., and we try to discover if there exists positive or negative bias toward a certain category blog. Recently, there have been rumors on the Internet saying that blogs rank lower than regular Web pages by Google [Blog1 2013; Blog2 2013]. We validate this belief using our ranking system.

To characterize Web sites into different categories, we apply a simple keyword approach. For example, in case of the blog category, we looked at keyword “blog” in the URL, title tag, and Google’s snippet, in a case-insensitive manner. If the keyword appears in any of the three places (regardless of how many times it appears), we categorized the page as a blog. We apply a similar approach to all other categories using different keywords.

Next, we add a new feature to our ranking algorithm, named CAT (which stands for category) and explore each of the categories *in isolation*. For example, when we explore the news category, we test both potential positive and negative bias towards the given category (by assigning a positive or negative unit value to the CAT feature parameter in the linear model for pages from the given category); at the same time, all other non-news-related Web pages that are returned by the search engine have the CAT feature turned off. We then run our ranking system with all features (including CAT) using the LP model and three iterations of recursive partitioning and ranking.

For all explored categories, except for blogs, the assigned category feature (either positive or negative) has *no* influence whatsoever on the ranking prediction (the result was the same as when we used the original 17 features). Figure 6 compares the ranking results with and without testing the negative bias hypothesis for the blogs feature. The

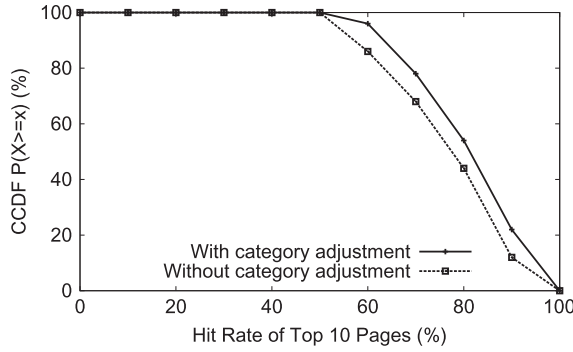


Fig. 6. Category analysis.

figure clearly shows that the negative hypothesis indeed shows true for blogs. Indeed, the prediction results improve when we adopt the negative bias hypothesis for blogs. The weight for the negative blog bias is as high as 0.2, which puts it in the top 10 features explored in our article.

For example, in our iptables dataset, a blog article titled “*Iptables dependency: why we got there and how we got out*” [Palaniswamy 2005] has high scores in all of our top 3 ranking factors. Specifically, it has a high page rank score of 6, and the target search term appears in the title and URL. However, this blog is ranked 62nd out of 100 by Google and it would have been ranked 22nd without the negative bias toward blogs. Another example is a blog titled “*Rsync Version 3 Alpha Out - O’Reilly ONLamp Blog*” [Gift 2007] in our rsync dataset that has everything it needs to be ranked as a top 10 Web page. It has a good page rank score of 4 and the keyword appears in the title, URL, and meta-description tag. In addition, it also has a good keyword density in its content. However, this blog is ranked 32nd by Google while it should have been ranked 8th without the bias.

Thus, we show that our system is capable of validating new conjectures about Google’s ranking algorithm and that Google imposes negative bias toward blogs.

5.3. HTML Syntax Errors Do Not Matter

In this final case study, we show our approach can be used to debunk certain myths about Web page ranking. In particular, some SEO experts hypothesised that Google estimates the quality of a Web page which includes the correctness of HTML syntax [Anderson 2007; RSS Pieces 2007]. Whether this is true or not has implications for how important it is for Web site developers to pay attention to the correctness of HTML syntax. Our study described next conducted by adding this new feature into our ranking system, demonstrates that the hypothesis is false.

In this experiment, we use the HTML tidy library program [HTML Tidy 2014] to analyze and count HTML errors of each Web page in our dataset. This new feature is then added to our ranking system to train new ranking models. In the training set, the number of Web pages with one or more syntax errors is 275 out of 1,500 pages (18.33%). In the testing set, the number of pages with syntax errors is 972 out of 4,500 (21.60%).

Figure 7 compares the results from the original and the new ranking models. The figure shows that the performance of the new model is very close to the original one. In addition, the weights of the new ranking feature in each of the 3 rounds are very close to zero (0.068, 0.056, and 0.033, respectively). This indicates that HTML syntax errors have very little to no impact on a Web page’s Google ranking.

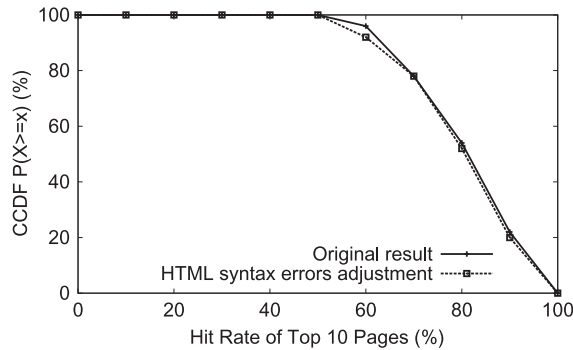


Fig. 7. HTML syntax analysis.

5.4. Debunking the Folklore

Here, we revisit the folklore listed in Section 2.2 by using the findings of our ranking system. We make the following observations.

- *PageRank is overlooked.* PageRank scores have been one of the key ranking factors since the birth of Google. However, this factor tended to be overlooked by SEO experts and search engine users. Only one out of the three lists in Table II explicitly includes PageRank score as one of the top 10 ranking factors. However, our results shown in Figure 4 demonstrate that page rank score is still one of the top ranking factors in Google’s ranking algorithm.
- *Keyword in URL plays a big role.* Based on our results, the following suggestions can be made to SEOs and Web masters: choosing an appropriate hostname and file-name can be as important as selecting a good title for a Web page. In our findings, we discover that the appearance of the keyword in the URL (including hostname and path) plays an important role in Google’s ranking algorithm. However, as described in Table II, only the *Survery* and *Idv* lists recognize this ranking factor, whereas the *SEOmox’07* list composed by SEO experts did not address the importance of this factor at all.
- *Meta-description tag counts.* The importance of metatags in a search engine’s ranking function has long been disputed. Some believe that Web masters can abuse metatags by stuffing irrelevant keywords into metatags, which could be the reason why some search engines hesitate to use metatags as ranking factors. Our measurement results show that this is indeed the case, as the metakeyword tag has little to no influence on the overall ranking. However, we find that meta-description is still one of the dominating ranking factors in Google’s ranking algorithm. Unfortunately, none of the three lists in Table II includes this ranking factor.

6. EXTENDED EVALUATION ON BING AND YAHOO! SEARCH ENGINES

Search engines have become the primary medium for Internet users to acquire knowledge from the Internet. However, these search engines often provide vastly differing search results on any particular search term. The cause of these differences between search engines includes the size of their indexed pages, that is, their indexed corpus, and more importantly, how a search engine presents its search results, that is, its ranking algorithm. In the previous sections, we have evaluated Google’s search results and made the following observations: (i) The page rank is the dominant factor that influences Google’s ranking results; (ii) Google’s ranking algorithm emphasizes search keywords appearing in the hostname, the title tag, and the meta-description tag.

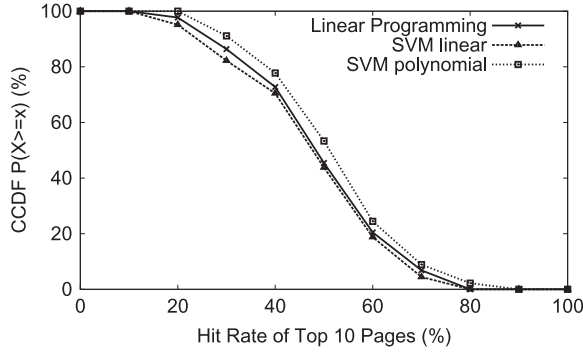


Fig. 8. Comparison of ranking models for Bing.

In this section, we evaluate the Bing search engine and explore: (i) whether our ranking system can predict Bing search engine’s ranking results; (ii) whether Bing and Google have different views on the importance of ranking features; and (iii) whether our ranking system can validate or disprove new ranking features for Bing?

Yahoo! is the third major company that shares a significant portion of the Web search market that we intend to evaluate. However, during the time of writing this article, Yahoo! has adopted Bing as the provider of its Web search service [Seth 2010], and we have confirmed this by obtaining identical search results from Bing and Yahoo!. Therefore, for the rest of the article we focus on Bing and our evaluation results apply to both search engines.

6.1. Experiment Setup

We first present our experiment setup for evaluating the Bing search engine. To evaluate Bing and compare its ranking model with Google’s, we use the crawler of our ranking system to collect ranking results from Bing for the same 60 keywords in Table III. We also intend to extract all the ranking features listed in Table I for Bing. However, Bing neither has page rank scores nor provides a search API for querying the date it indexed a Web page. Therefore, for the page rank score ranking feature (PR), we query Google using Google toolbar’s API to obtain the scores for each Web page in Bing’s search results. In Section 6.4.1, we use our ranking system to demonstrate that page rank scores from Google are beneficial for improving the accuracy of ranking prediction for Bing. As for the Web page (AGE) ranking feature, we have to exclude it from our evaluations since we have no means to obtain this information from Bing. The rest of the experiment setup is identical to that for evaluating Google, as described in Section 4.1.

6.2. Ranking Accuracy Analysis

We now evaluate the prediction accuracy of our ranking system for the Bing search engine. First, we compare the ranking accuracy of different ranking models, all operating under our recursive partitioning algorithm.

Figure 8 shows the results of our ranking system under the LP model and the two SVM models for Bing. The results show that the two linear programming ranking models have similar prediction accuracy while the SVM-polynomial model performs slightly better than the other two. For the LP model, for 44% of the keywords we evaluated, our system correctly predicts more than 50% of the top 10 pages. Compared to our evaluation results for Google presented in Section 4, our ranking system performs better in approximating Google’s ranking algorithm than Bing’s. We believe it is because

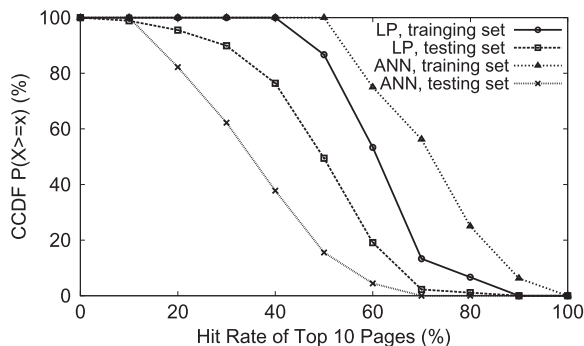


Fig. 9. Artificial neural network ranking model.

Bing’s ranking algorithm is based on a machine learning ranking model [RankNet 2009], and hence is not a good “fit” for the numerical ranking models used in our experiments. In addition, Bing might use a different set of features in its ranking system that we were not able to measure. Shortly, we further investigate the machine learning ranking methodology by developing a new ranking model based on an Artificial Neural Network (ANN) [RankNet 2009] and integrating it into our ranking system.

ANNs are based on a machine learning methodology inspired by the functional aspect of biological neural networks. The ANN learning methods are robust to noise in the dataset but the models or functions learned by an ANN are often difficult for humans to interpret [Mitchell 1997]. Here, we examine whether we can build an ANN ranking model to better approximate Bing’s ranking results than the numerical models we have studied. In particular, the ANN ranking model we developed is based on the research in Cheng et al. [2008] and is implemented on top of the fast artificial neural network library [FANN 2014].

Figure 9 plots the ranking results of our LP and the new ANN ranking models with our recursive partitioning algorithm. For better understanding of the characteristics of the two ranking models, we use the learned models to evaluate both the training and the testing datasets. First, we compare the prediction accuracy of the two models for the training set. The figure shows that the ANN ranking model performs better than the LP model for predicting ranking results on the training set. This is because ANN can better adapt to Bing’s nonequational ranking model [RankNet 2009].

However, the accuracy of an ANN is also prone to overfitting problems because uncertainty (e.g., missing ranking features) and noise (e.g., approximated ranking scores) prevent ANN from obtaining a more generalized ranking model. As shown in the figure, the learned ranking model by ANN performs worse than the LP model in predicting the testing set. Next, we evaluate the effectiveness of our recursive partitioning algorithm for predicting Bing’s ranking results.

Figure 10 compares the results of the LP model and SVM-polynomial model with different rounds of recursive partitioning. The figure shows that the recursive partitioning algorithm is able to improve the ranking accuracy for both models. For example, for the LP model, the average hit rate after three rounds of recursion improves by 16% over the case without using the recursion algorithm. This shows that the recursive partitioning algorithm is valuable for improving the accuracy of our ranking system for both Bing and Google.

6.3. Relative Importance of Ranking Features for Bing

Next, we evaluate the relative importance of different ranking features for the Bing search engine. Figure 11 shows the weights of different ranking features for the LP and

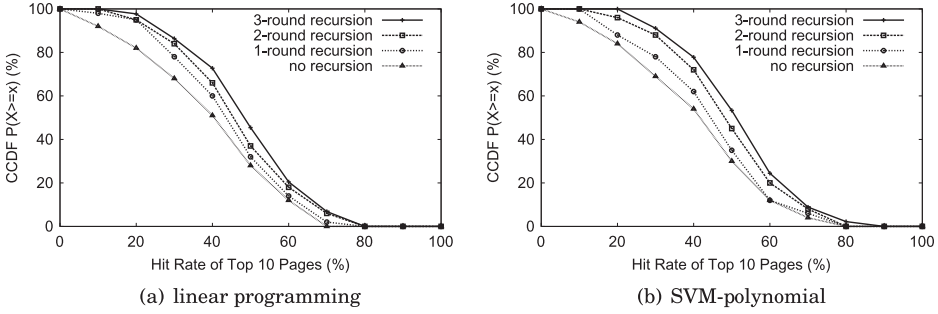


Fig. 10. Evaluating recursive partitioning for Bing.

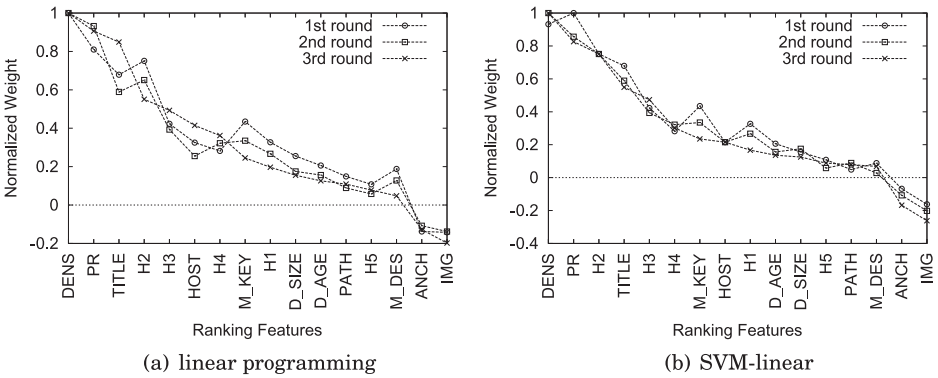


Fig. 11. Weights in different rounds for Bing.

SVM-linear ranking models with our recursive partitioning algorithm. Our first observation is that the most important ranking feature for Bing is keyword density (DENS) in a Web page. This indicates that Bing emphasises on content-based (i.e., in-page) ranking features. Next, the results show that Bing favors metakeyword rather than meta-description, while our results for Google show the opposite relationship. Moreover, to our surprise, the H1 tag ranking feature weighs less than the other header tags in the results. We believe it is because the H1 tag, which uses the largest font, is not as commonly used as the other header tags in the Web pages. Therefore, Bing's ANN learning algorithm attributes less weight to this ranking feature.

6.4. Case Studies

Next, we present several case studies on the Bing search engine. First, we analyze the impact of *Google's* page rank and the number of incoming links to Bing's ranking results. Next, we explore whether the overall Web traffic of a Web site affects Bing's ranking algorithm.

6.4.1. Quantity versus Quality of Incoming Links. In this case study, we analyze whether the quantity or the quality of incoming links has more impact on Bing's ranking results than the other. Incoming links are URL links leading to the target Web page and the number of incoming links is one indication of the popularity of that Web page. On the other hand, *Google's* page rank score measures the quality of a Web page by assigning weighted numerical value to each incoming link based on the page rank algorithm. In this experiment, we obtain the page rank score from Google and the number of incoming links from Yahoo! site explorer for every Web page in our dataset. We use our

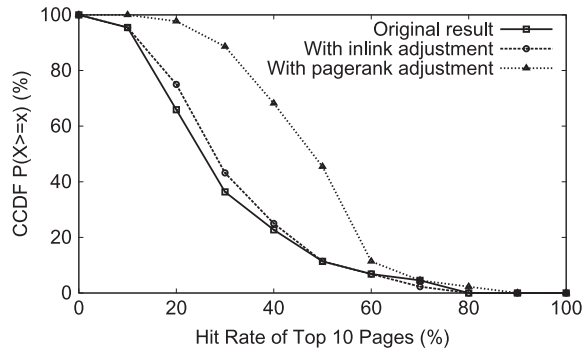


Fig. 12. Incoming links and page rank analysis.

ranking system with the recursive partitioning algorithm in evaluating the impact on these two ranking features.

Figure 12 compares the results of the two ranking features. The curve of original result shown in the figure plots the accuracy of our ranking system without the in-link and page rank ranking features. The in-link adjustment curve plots the ranking accuracy by adding only the in-link ranking feature into the ranking system. Similarly, the page rank adjustment curve plots the ranking accuracy by adding only the page rank ranking feature to our ranking system. The figure shows that *Google's* page rank ranking feature significantly improves the accuracy of our ranking engine while the impact of the number of incoming links ranking feature is negligible. This case study clearly shows that the quality of incoming links is much more important than the quantity. Furthermore, the results show that Bing also measures the quality of incoming links and include the results into its ranking model. Although its underlying methodology is not publicly available, *Google's* page rank score can provide an approximation to the measure.

6.4.2. Impact of Popular Web Sites. Next, we evaluate whether the overall Web traffic of a Web site impacts Bing's ranking model. In particular, we want to explore whether a popular Web site would obtain better placements for its Web pages in Bing search. If that is the case, a positive feedback effect will be formed because the popular Web sites will get even more popular by attracting more Web traffic from the search engine. In this case study, we obtain the Web traffic ranking from Alexa [2014] and add it as a new ranking feature to our ranking engine. We evaluate the impact of the newly added ranking feature by comparing the new ranking results to the original ones.

Figure 13 shows the results from this experiment. The figure shows that the impact of the Web traffic ranking feature is negligible in the overall prediction accuracy. This indicates that Bing's ranking model does not favor popular Web sites, and hence the Web pages from small or young Web sites are equally evaluated in the placement on Bing's ranking results.

6.5. Summary

In this section, we have evaluated the Bing search engine and made the following observations: (i) our evaluation results suggest that page rank scores obtained from *Google* can help to improve the accuracy of our ranking engine in predicting Bing's ranking results. We hypothesize that Bing's ranking algorithm incorporates a link analysis mechanism that is similar to *Google's* page rank algorithm. (ii) The major disagreement on the importance of ranking features between *Google's* and Bing's ranking algorithm includes: Bing puts much more emphasis on the keyword density

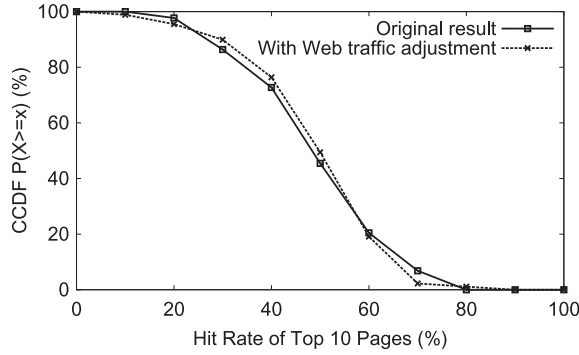


Fig. 13. Web traffic analysis.

ranking feature than Google, and Bing weighs the metakeyword tag more than the meta-description tag while Google weighs them oppositely. (iii) Our case study shows that Bing’s ranking algorithm has no bias towards popular Web sites.

7. DISCUSSION

Limited set of features. While we have demonstrated that the set of features shown in Table I is sufficient to predict search engine ranking results with high accuracy, this set of features is necessarily not comprehensive. In the following we explain why a limited set of features is still capable of providing high accuracy and why a noncomprehensive set of features can still have a practical application.

For example, it is likely that a search engine may use features such as user click data mined from search engine logs as one of the ranking features. Even if we knew for sure that this is the case, it would be challenging to accurately measure or estimate such information from an endpoint. However, by carefully probing the system, that is, selecting the appropriate set of keywords, it is possible to avoid the effects of such a parameter. For example, as we have shown in Section 4.1 (keywords shown in Table III), by selecting appropriate classes of keywords, it is possible to minimize the effects of parameters such as user click data, simply because the corresponding landing Web pages likely have negligible click rates. As another example, in Section 5.1 we used another set of keywords and crawled only young Web pages, thus again indirectly avoiding the given parameter (i.e., young Web pages in general have low click rates).

The preceding experiments have demonstrated that, despite the fact that we do not consider the given parameter while the search engine likely does (as implied from the results shown in Figure 5), we are still able to provide meaningful feedback to the end-user regarding the parameters that the end-user can actually control. Given that end-users should in general worry only about the page features that they can actually control, our methodology still provides practical feedback to an SEO user.

Recursive partitioning. Models such as gradient boosted regression trees [Friedman 2001, 2003] are known to perform well for purposes such as the one treated in this article, particularly in terms of scalability [Panda et al. 2009]. Nonetheless, we have demonstrated that linear learning models, coupled with a recursive partitioning ranking algorithm, are also capable of performing well in terms of predicting ranking results with high accuracy.

Discussion of the findings. Our results demonstrate that it is possible to estimate relative importance among different parameters used by a search engine. The key idea in resolving this inherently complex problem lies in the ability to effectively isolate the

corresponding space that reveals the relationship among the given parameters. For example, choosing unpopular keywords whose search results land at unpopular Web pages that attract negligible user traffic can help avoid the impact of page popularity, and hence enable a more accurate estimation of the remaining parameters.

Practical implications. Our work implies that it is possible to add a level of formality and “sanity” in the area of search engine optimization. As a result, we expect that methods like ours will be used to help legitimately optimize Web pages. It should also be clear that such methods are unable to “push a Web page to the top of the list” simply because there are parameters that are not under the control of an SEO user (such as a Web page’s popularity). Another area where we anticipate a practical application of our methods is search engine auditing, which we discuss next.

Avenues for future work. One of the key contributions of our work is to demonstrate that analyzing a search engine’s performance, despite a significant complexity, is a feasible task. While traditionally the key application of such analysis is in SEO optimization, another relevant and growing area is auditing search engines, that is, detecting a potential search bias. Indeed, the most persistent critique of search engines is the fear that they create reality instead of reflecting it [Grimmelmann 2013]. Our results strongly imply that it is feasible to develop such auditing mechanisms for search engines. Necessarily, additional research as well as formal and standardized methods are needed to achieve this goal.

8. CONCLUSIONS

In this article we study the problem of predicting Google’s and Bing’s ranking results. While we necessarily limited our analysis to the two search engines, our approach is certainly more generic and applicable to other search engines. Even though Google’s internal ranking function can be very complex, we demonstrated that it is possible to approximate Google’s organic search results by adopting a linear model trained with a linear programming optimizer along with a recursive partitioning scheme. We performed large-scale experiments using over 6000 Web pages and showed that our ranking system is capable of predicting Google’s ranking results with high accuracy. Specifically, our system is able to correctly predict 7 out of the top 10 pages for 78% of evaluated keywords.

Using our ranking system, we revealed the relative importance of ranking features in Google’s ranking function. In particular, page rank is the dominant factor, followed by the search keyword appearing in the hostname, in the title tag of the HTML header, in the meta-description tag, in the path segment of the URL, as the other leading factors. Such revelation provides guidelines for SEOs and Web masters to optimize their Web pages and obtain a better position in Google’s search result pages. Moreover, we showed how to use our system to validate or disprove new ranking features and to evaluate search engine ranking results for possible ranking bias. In particular, we used our ranking system to confirm the rumors in the Internet that the Google’s ranking is biased against blogs.

Finally, we evaluated our ranking system for the Bing search engine. We showed that our recursive partitioning algorithm is also beneficial for improving the prediction accuracy for Bing. In addition, we demonstrated that the page rank scores obtained from Google can significantly improve the ranking accuracy of our ranking system in predicting Bing’s ranking results. Furthermore, our results also revealed several disagreements on the relative importance of ranking features between the two search engine giants. In particular, Bing places a larger weight on the metakeyword tag than the meta-description tag while Google treats the two tags oppositely.

REFERENCES

- AccuraCast. 2007. Google algorithm's top ranking factors. <http://www accuracast.com/seo-weekly/ranking-factors.php>.
- Alexa. 2014. Alexa, the web information company. <http://www.alexa.com/>.
- Anderson, S. 2007. Google seo test google prefers valid HTML and CSS. <http://www.hobo-web.co.uk/seo-blog/index.php/official-google-prefers-valid-html-css/>.
- Aubuchon, V. 2010. Google ranking factors. <http://www.vaughns-1-pagers.com/internet/google-ranking-factors.htm>.
- Benczur, A. A., Csalogany, K., Sarlos, T., and Uher, M. 2005. Spamrank—Fully automatic link spam detection. In *Proceedings of the 1st International Workshop on Adversarial Information Retrieval on the Web (AIRWeb'05)*. 25–38.
- Blog1. 2013. How google works: Why does crappy website rank higher than mine? <http://www.trafficgenerationcafe.com/how-google-works-relevance/>.
- Blog2. 2013. My actual blog post ranks lower than pages associated with it. <http://productforums.google.com/forum/#!topic/webmasters/QOv273CK07I>.
- Cheng, J., Wang, Z., and Pollastri, G. 2008. A neural network approach to ordinal regression. In *Proceedings of the IEEE International Joint Conference on Neural Networks (IJCNN'08)*. IEEE World Congress on Computational Intelligence, 1279–1284.
- Cho, J. and Roy, S. 2004. Impact of search engines on page popularity. In *Proceedings of the 13th International Conference on World Wide Web (WWW'04)*. ACM Press, New York, 20–29.
- CPLEX. 2014. Ilog cplex: High-performance software for mathematical programming and optimization. <http://www.ilog.com/products/cplex/>.
- Davison, B. D. 2000. Recognizing nepotistic links on the web. In *Proceedings of the AAAI Workshop on Artificial Intelligence for Web Search*.
- FANN. 2014. Fast artificial neural network library. <http://leenissen.dk/fann>.
- Friedman, J. 2001. Greedy function approximation: A gradient boosting machine. *Annals Statist.* 29, 5, 1182–1232.
- Friedman, J. 2003. Statistical gradient boosting. *Comput. Statist. Data Anal.* 38, 4, 367–378.
- Gift, N. 2007. RSYNC version 3 alpha out - O'reilly onlamp blog. <http://www.oreillynet.com/onlamp/blog/2007/10/rsync.version.3.alpha.out.html>.
- Google. 2014a. Pagerank on Google toolbar. <http://www.google.com/support/toolbar/bin/answer.py?hl=en&answer=79837>.
- Google. 2014b. Google trends. <http://www.google.com/trends>.
- Google. 2014c. Google webmaster tools. <http://www.google.com/webmasters/>.
- Grimmelmann, J. 2013. What to do about Google? *Comm. ACM* 56, 9, 28–30.
- Gyongyi, Z., Berkhin, P., Garcia-Molina, H., and Pedersen, J. 2006. Link spam detection based on mass estimation. In *Proceedings of the 32nd International Conference on Very Large Data Bases (VLDB'06)*. 439–450.
- Hopkins, L. 2012. Online reputation management: Why the first page of Google matters so much. <http://www.leehopkins.net/2012/08/30/online-reputation-management-why-the-first-page-of-google-matters-so-much/>.
- HTML Tidy. 2014. HTML tidy library project. <http://tidy.sourceforge.net/>.
- Joachims, T. 1999. Making large-scale SVM learning practical. In *Advances in Kernel Methods—Support Vector Learning*, B. Scholkopf, C. J. C. Burges, and A. J. Smola Eds., MIT Press, Cambridge, MA, 169–184.
- Joachims, T. 2009. SVM-rank support vector machine. http://www.cs.cornell.edu/People/tj/svm.light/svm_rank.html.
- Joachims, T., Granka, L., Pang, B., Hembrooke, H., and Gay, G. 2005. Accurately interpreting click-through data as implicit feedback. In *Proceedings of the ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR'05)*. 154–161.
- Kontopoulos, G. 2007. Top Google ranking factors. <http://www.squidoo.com/topGoogleRankingFactors>.
- Li, L. 2008. Google's top search engine ranking factors. <http://lornali.com/online-marketing/seo/googles-top-search-engine-ranking-factors>.
- Marshall, B. 2009. Top 10 most important Google ranking factors. <http://blogs.myspace.com/index.cfm?fuseaction=blog.view&friendId=21196&blogId=493022330>.
- Mitchell, T. 1997. *Machine Learning*. McGraw-Hill.

- Moran, M. and Hunt, B. 2005. *Search Engine Marketing, Inc.: Driving Search Traffic to Your Company's Web Site*. Prentice Hall PTR, Upper Saddle River, NJ.
- Palaniswamy, A. 2005. Iptables dependency: Why we got there and how we got out. <http://www.zimbrablog.com/blog/archives/2005/11/iptables-dependency-why-we-got-there-and-howwe-got-out.html>.
- Panda, B., Herbach, J. S., Basu, S., and Bayardo, R. J. 2009. Planet: Massively parallel learning of tree ensembles with mapreduce. In *Proceedings of the 35th International Conference on Very Large Data Bases (VLDB'09)*.
- Patel, N. 2006. A breakdown of Google's ranking factors. <http://www.pronetadvertising.com/articles/a-breakdown-of-googles-ranking-factors.html>.
- PHP HTML Parser. 2014. PHP simple HTML dom parser. <http://simplehtmldom.sourceforge.net>.
- RankNet. 2009. RankNet: How Bing works. <http://neotracks.blogspot.com/2009/06/ranknethow-bing-works.html>.
- RSS Pieces. 2007. HTML validation: The hidden key to seo. <http://www.rsspieces.com/html-validation-the-hidden-key-to-seo>.
- SEOMoz. 2007. Search engine ranking factors. <http://www.seomoz.org/article/search-ranking-factors>.
- SeoPros. 2014. Seopros.org. <http://www.seopros.org/>.
- Seth, S. 2010. Yahoo! Transitions organic search back-end to Microsoft platform. <http://www.ysearchblog.com/2010/08/24/yahoo-transitions-organic-search-back-end-to-microsoft-platform/>.
- Su, A.-J., Hu, Y. C., Kuzmanovic, A., and Koh, C.-K. 2010. How to improve your Google ranking: Myths and reality. In *Proceedings of the IEEE/WIC/ACM International Conference on Web Intelligence (WI'10)*.
- TOPSEOs. 2014. Topseos.com. <http://www.topseos.com/>.
- Vapnik, N. V. 2000. *The Nature of Statistical Learning Theory*. Springer.
- Wu, B. and Davison, B. D. 2005. Identifying link farm spam pages. In *Proceedings of the Special Interest Tracks and Posters of the 14th International Conference on World Wide Web (WWW'05)*. 820–829.

Received February 2011; revised September 2013; accepted October 2013