# Ranking Search Results in Library Information Systems — Considering Ranking Approaches Adapted From Web Search Engines

Christiane Behnert *, Dirk Lewandowski

*Hamburg University of Applied Sciences, Faculty of Design, Media & Information, Department of Information, Finkenau 35, 22081 Hamburg, Germany*

## ABSTRACT

For an information retrieval system to be successful, it must have the ability to rank search results. As web search engines are the most often used and — in terms of ranking functionality — the most advanced existing systems, the principles they are based on and the strategies they use can be advantageous when applied to the library context. We categorize ranking factors into six different groups: 1. text statistics, 2. popularity, 3. freshness, 4. locality and availability, 5. content properties and 6. user background. We discuss the basic concepts and assumptions these ranking factors involve and offer potential implementations in the library context. The practice recommended here is for libraries to not only apply selected ranking factors — as existing library information systems already do — but to systematically test for the ranking factors best suited to their systems. We argue for a user-centric view on ranking, because in the end, ranking should be for the benefit of the user, and user preferences may vary across different contexts.

© 2015 Elsevier Inc. All rights reserved.

## INTRODUCTION

In order to understand the factors behind relevance ranking, this article surveys conceptual approaches behind web search engine ranking and how ranking factors can be adopted to library information systems. The exemplary search results ranking performed by web search engines can be a useful model for other information systems providers, especially libraries, to emulate. Since people are now used to web search interfaces and relevancy-ranked results lists, they expect searching in library catalogs to be as easy, and the presentation of results to be as good, as when they search the web. The aim of this article is to provide librarians and system developers with an overview of suitable ranking factors as used in web search engines from an academic perspective and offer recommendations for applying these factors (or their underlying principles) to library information systems.

The reason search results are ranked in an information retrieval (IR) system derives from the assumption that information-seeking users should get *all* the information relevant to their search query and *only* that information. In order to help the user judge the relevance of a single search result, the results are presented in a certain way — the most relevant documents are presented first, with less relevant documents beneath them. This raises the question: How does the IR system "know" which documents are (most) relevant to satisfying an information need? A clear definition of the term *relevance* is problematic, and differing views on the meaning of *relevance* can lead to misunderstandings (Bade,

2007; Mizzaro, 1997; Saracevic, 2015), as it is highly subjective and understood intuitively (Saracevic, 1996, 2006).

Although mathematical and statistical methods of varying complexity do exist to determine the relevance of a search result, such methods use algorithms to integrate *assumptions of relevance*. But it is the subjective relevance of a result that matters to the user in the end (Bade, 2007), "because an information-retrieval system exists only to serve its users" (Swanson, 1986, p. 390). This concept of subjective relevance can be referred to as pertinence, defined as the user's cognitive ability to understand the knowledge obtainable from a search result (Stock & Stock, 2013). We can regard any effort put into systems determining the relevance of search results as an effort to find ways to algorithmically model the users' views on relevance. In this article, we therefore discuss ranking factors as basic ideas of how we can technically simulate users' relevance judgments.

The need to rank search results derives from the behavior of the typical user, who is either unwilling or unable to assess *all* the results shown in response to a given query. There are two general reasons for this. The first is that there may simply be too many items in the database, and only some of them are needed. The second is that a users' query may be too general or ambiguous, generating a large number of results. Relevance ranking can at least partially compensate for a user's inability to construct queries that lead to a well-defined number of hits.

Research investigating web search engine user behavior offers us a good general impression of how users search. As we will see, many of the characteristics of search engine usage are also applicable to library information systems. Several studies have been conducted to analyze search behavior in the context of web IR (i.e., methods of information retrieval in the context of the World Wide Web) and the findings are

* Corresponding author. Tel.: + 49 40 42875 3653.
*E-mail addresses:* christiane.behnert@haw-hamburg.de (C. Behnert),
dirk.lewandowski@haw-hamburg.de (D. Lewandowski).

that the majority of queries consists only of one or two words, whereas according to Bendersky & Croft (2009), long queries, i.e., queries consisting of 5 or more terms, represent only 10% of the query volume. Usually with short queries, Boolean operators are rarely or only implicitly used (Höchstötter & Koch, 2008). Furthermore, users only look at the first result page and consider mainly the top-ranked hits (Barry & Lardner, 2011; Goel, Broder, Gabrilovich, & Pang, 2010; Jansen & Spink, 2006; Pan et al., 2007; Schmidt-Mänz & Koch, 2005; Spink, Wolfram, Jansen, & Saracevic, 2001).

Studies also showed that web search often acts as a starting point in the information seeking process (Rowlands et al., 2008): before users start searching in library catalogs, they tend to obtain information on the desired materials via web search, and then carry on searching in the Online Public Access Catalog (OPAC) (Pera, Lund, & Ng, 2009) or the library's website (De Rosa et al., 2005, 2010). Thus, one major implication for library systems is that they "need to look and function more like search engines" (Connaway & Dickey, 2010, p. 5).

When searching the library OPAC, generally the same search and browsing behavior as in search engines could be observed (Hennies & Dressler, 2006): users consider the top results on the first result page to be most relevant (Antelman, Lynema, & Pace, 2006, p. 135). Queries also usually consist of only a few words, i.e. one, two, or three words (Niu & Hemminger, 2010; Schneider, 2009). Studies also show that users rely on default settings (Asher, Duke, & Wilson, 2013; Jones, Cunningham, McNab, & Boddie, 2000) and, more importantly, that they expect a library catalog to have the same search capabilities and options for displaying results as they are accustomed to from web search engines (Yu & Young, 2004).

Academic researchers often use specialized scientific web search engines such as Google Scholar to find journal articles and other sources of information. In the library context, scholarly articles have not been as easily searchable nor have they been directly available (Lewandowski, 2010a). Traditional OPACs with "second-generation" features (e.g. cross-references and exact match Boolean search) (Antelman et al., 2006) still lack a single search interface that allows searching across multiple databases (Luther, 2003), which users expect, having grown accustomed to it from searching on the web. Instead, articles are searchable in separate databases or portals. End users are frequently uncertain which database to choose.

When comparing search functionality and how results are displayed in web search engines vs. library information systems, we should also note that there are certain characteristics of the contents in the respective databases that make library materials somewhat more difficult to rank. Results presented by library catalogs are bibliographic records, i.e., metadata. We have (1) the metadata of printed and other physically tangible materials, for example books, periodicals, CDs, DVDs, and maps, and (2) the metadata of digital contents, for example licensed e-journals and even links to other external content such as audio and video files. Library materials increasingly comprise more than just printed monographs and journal articles. Now, "web content" such as links to licensed e-journals, e-books, research data and infographics are also included.

Traditional IR techniques alone are insufficient for these types of library content. Because of the change in user behavior when submitting search queries and the expectation that result quality will be indicated by means of a ranking, it is important to implement ranking factors in library information systems inspired by web IR. Traditional OPACs lack relevance ranking, despite the fact that "[a]lphabetizing makes for easy lookups, but ranking is better for human interest" (White, 2007, p. 600). As a consequence, the integration of search engine technology into library catalogs via discovery software is an essential component of solving OPAC ranking problems (Lewandowski, 2009, 2010b; Schneider, 2006).

Ranking features have already been implemented in *next-generation catalogs* and *discovery tools*, which enable users to not only find but also access licensed materials. Along with enriched content, faceted navigation and spell-checking, one of the defining features of discovery

systems is relevance ranking (Yang & Hofmann, 2011; Yang & Wagner, 2010). Discovery tools such as Serial Solutions' *Summon* or ExLibris' *Primo* provide ranked search result lists using web technology that corresponds more closely to user expectations than traditional catalogs (Breeding, 2006, 2007). With open source software such as *VuFind* and *Blacklight*, libraries can take things one step further. These applications give libraries control over the technology and the ability to set up their own relevance rankings (Oberhauser, 2010; Parry, 2010). Whichever approach is chosen, what current systems have in common is that they apply *some* ranking factors, but lack a systematic review of possible factors to decide from.

Below, we discuss ranking factors used by web search engines and their potential adaptation for use in library information systems. In contrary to the web search industry's perspective of improving web search systems or the search engine optimization (SEO) community's perspectives in terms of increasing the visibility of websites, we aim for showing in which regard ranking concepts from web search and from library information systems relate to each other. For this purpose, we avoid going into details of (technical) *ranking signals* or website design elements, as, for instance, can be categorized into on-the-page and off-the-page factors (Sullivan, 2015) or into positive and negative website elements (Weideman, 2009). Instead, we focus on the basic concepts of relevance ranking and categorize the ranking factors into six groups, being modified after Lewandowski (2009). Each group is illustrated with an overview of the individual factors. The first group, *text statistics*, comprises factors which are primarily derived from traditional IR methods. Text statistics include the fundamental ranking factors for all text-based retrieval systems, because there always has to be a query text that can be matched with the documents' representation if any search results at all are to be obtained. Since such ranking factors alone cannot lead to a *quality-induced* ranking, there are other factors building on this first group, as shown in Fig. 1. These factors consider the "wisdom of crowds" and rank results based on a document's *popularity*. Another group is *freshness*. The up-to-dateness of a document is not only important in web IR, it is also the standard ranking concept used in traditional library catalogs since their inception. Within the group *locality & availability*, ranking factors consider the physical location of both the user and the document, since mobile data connectivity now enables access independent of physical location. Apart from these four major ranking groups, we introduce two others which provide additional valuable information for relevance ranking. The group *content properties* includes characteristics of the document content, while the factors contained within the last group, *user background*, derive from characteristics of the user. In the last section of this article, we summarize the discussed ranking factors and offer suggestions for the development of future ranking functions.
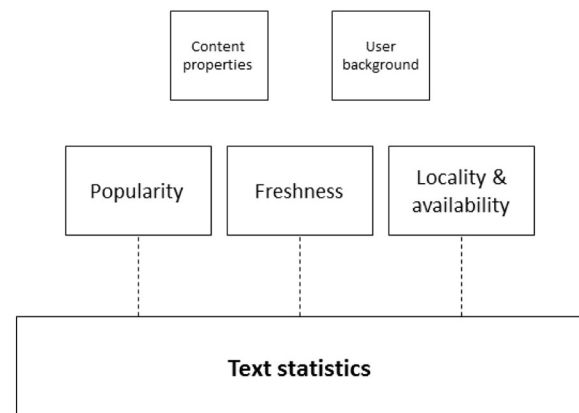


**Fig. 1.** Overview of ranking factor groups.