# Relevance criteria identified by health information users during Web searches

Abe Crystal  (corresponding author)
School of Information and Library Science
University of North Carolina at Chapel Hill

100 Manning Hall, CB #3360

Chapel Hill, NC 27599-3360

919.593.6129 (voice)

919.962.8071 (fax)
abe@unc.edu


Jane Greenberg
School of Information and Library Science
University of North Carolina at Chapel Hill

100 Manning Hall, CB #3360

Chapel Hill, NC 27599-3360

919.962.8066  (voice)

919.962.8071 (fax)

janeg@ils.unc.edu

## Abstract

This study focused on the relevance judgments made by health information users using the Web.  Health information users were conceptualized as motivated information users concerned about how an environmental issue affects their health.  Users identified their own environmental health interests, and conducted a Web search of a particular environmental health Web site.  Users were asked to identify (by highlighting with a mouse) the criteria they use to assess relevance in both Web search engine surrogates and full-text Web documents. Content analysis of document criteria highlighted by users identified the criteria these users relied on most often.  Key criteria identified included (in order of frequency of appearance): research, topic, scope, data, influence, affiliation, Web characteristics, and authority/person. A power-law distribution of criteria was observed (a few criteria represented most of the highlighted regions, with a long tail of occasionally-used criteria).  Implications of this work are that IR systems should be tailored based on users' tendencies to rely on certain document criteria, and that relevance research should combine methods to gather richer, contextualized data. Metadata for IR systems, such as that used in search engine surrogates, could be improved by taking into account actual usage of relevance criteria.  Such metadata should be user-centered (based on data from users, as in this study) and context-appropriate (fit to users' situations and tasks).

## *Introduction*

John, a thirty four-year old father of two, is concerned about the impact of car exhaust on his young son's asthma. On his lunch break one day, he opens his Web browser and loads the Website of a well-known nonprofit organization focused on healthy environments. He inputs the query "car exhaust asthma" into the site's search engine, which returns a long list of results. Among the results are position papers urging Congress to adopt stricter pollution standards, scientific papers on various environmental topics, reviews of popular books on asthma, and brochures for school officials to distribute to parents. How will John determine which of these many documents are relevant to his interests? And how could the search engine better support his needs?

Scenarios like this one motivated us to explore the relevance judgments of Web users. Our goal is to support the design of more effective information retrieval (IR) systems for the Web by providing a better understanding of what document criteria users employ to assess documents in the context of a Web search. Specifically, we examine Web relevance evaluation in the framework of user-centered IR research (Sugar, 1995), with an emphasis on document criteria (particularly those criteria that can be represented via structured metadata).

A key purpose of IR systems is providing relevant documents to users (Borlund, 2003). Because relevance is of such central importance in IR, many user studies have examined how users make relevance judgments (Mizzaro, 1997). In particular, researchers have examined the specific *criteria* users employ to evaluate documents and assess relevance (Barry, 1994; Carlyle, 1999). This research is important because it can ground the design of IR systems in users' actual needs and abilities. For example, if user studies find that people consistently use the *currency* of documents to make relevance judgments (Schamber,

Eisenberg & Nilan, 1990), IR systems can be designed to return more recent documents, or to include surrogates that clearly highlight publication dates.

Our study extends this line of research to Web IR. Previous studies have focused on users' relevance judgments when interacting with OPACs, bibliographic databases, or simply sets of documents, abstracted from any system (e.g., Janes, 1991; Barry, 1994). The Web provides a notably different context for IR interaction: it contains diverse types of documents incorporating hypertext and multimedia, information is provided by many different authors and publishers, and documents are organized in numerous heterogeneous systems, from simple lists to complex hyperlinked networks. Research is needed to address how relevance judgments on the Web may differ from those in earlier IR environments, and then to incorporate these more nuanced conceptions into IR systems (Kekäläinen & Järvelin, 2002). Our study addresses this need by exploring the criteria employed by Web users to make relevance judgments. Our user study complements recent work based on transaction logs (Spink, 2003), as well as other user studies with different foci (Rieh, 2002; Tombros, Ruthven & Jose, 2005). These empirical findings can, in turn, be used to support improved Web IR systems through the more effective use of metadata, surrogates, retrieval algorithms, and filtering mechanisms.

Our study builds upon recent advances in relevance research by incorporating four key approaches:

1. A two-stage model of relevance judgment, in which users evaluate surrogates and then full-text documents (Tang & Solomon, 2001; Rieh, 2002).

2. Consideration of both individual criteria and groups of criteria (Barry, 1998; Tang & Solomon, 2001).

3. A naturalistic focus on motivated users with real information needs (Barry, 1994; Tang & Solomon, 2001).

4. A focus on document criteria that can be effectively exploited in IR system design and metadata development (as emphasized by Borlund (2003, p. 918)).

We focus on a particular type of Web user: "health information users," who are concerned about the impact of the environment on their health (as in the scenario outlined above). Our study is therefore particularly relevant to the design of IR systems that provide health information. Our research method, on the other hand, is intended to be widely applicable, and could be used in other contexts (personal finance, for example). We hope that our methods and results will provide a foundation both for future empirical research and for innovations in IR system and interface designs.

## Background

The conceptual framework guiding this study was derived from previous research on information-seeking behavior and relevance judgments (Marchionini, 1995; Tang & Solomon, 2001; Rieh, 2002). The framework models individuals who undertake information retrieval as part of particular tasks and work situations. They seek information relevant to their interests—information they can make sense of and integrate into their work. IR systems support this information-seeking behavior by returning document surrogates and documents to users based on the queries they pose to the system (see Figure 1).

[ Figure 1 appears about here ]

**Figure 1. Conceptual framework: relevance judgment as part of information-seeking process.**

Relevance judgments are embedded in larger information-seeking processes, but it is when results are available to examine that users focus on assessing them. This activity is represented in a Marchionini's (1995) model of information-seeking behavior, which appears on the left-hand side of Figure 1. Interaction with the results provided by the IR system is a two-phase process (Tang & Solomon, 2001). In Phase 1, users examine *surrogates* provided by the system. The surrogates consist of specific *metadata elements*, such as the title of a document, the date it was created or published, a description of the document, or subject keywords. Phase 1 allows users to choose among the set of documents returned by the system. Since a surrogate provides only a fraction of the information in the full document, Phase 1 can be considered a process of "predictive judgment," in which users try to predict which documents will ultimately be useful (Rieh, 2002). Phase 2 encompasses evaluation of specific *documents* to assess whether they are relevant and useful in the context of the users' task and situation. So Phase 2 represents a process of "evaluative judgment" (Rieh, 2002).

The challenge for IR system design is to effectively support users in both of these phases. Empirical research can provide a foundation for this design challenge by establishing usage patterns and user preferences. For each phase, one can ask, "What characteristics of the presented information do people find useful?"

*Surrogates, metadata, and relevance judgments*

The challenge for Phase 1 of IR interaction is to construct *surrogates* (or more generally, *document representations*) that enable users to efficiently and accurately assess the relevance of a whole document (Barry, 1998). Good surrogates should provide metadata that enables users to quickly and accurately predict the relevance of the document. To

design such surrogates, "it seems clear that we should first understand what metadata elements are important or useful to users" (Lan, 2002, p. 13). Two key questions are:

- *How much* metadata is useful?

- *Which* metadata elements are particularly valuable?

These questions can be examined by studying how people interact with surrogates.

*How much metadata is useful?* Some researchers have argued that more metadata is better. According to the "length hypothesis," longer surrogates, by providing more information to users, lead to more accurate relevance judgments (see Mizzaro, 1997 for a review). One limitation of these studies is that they were conducted in a restricted IR environment based on retrieval of bibliographic records, rather than retrieval of full text, as is common today. In addition, these studies focused strictly on binary relevance judgments rather than the larger context of IR and support for the user's task. User-centered IR and human-computer interaction (HCI) researchers (Sugar, 1995; Shneiderman, Byrd & Croft, 1998) argue for examining the *usability* of IR systems in terms of how well they facilitate task completion.

Some recent studies have taken this user-centered perspective when evaluating the length hypothesis. Drori (2000a, 2000b) augmented Web search surrogates by adding metadata elements. He then studied the usability of these surrogates. He found that in some cases the augmented surrogates improved user performance and satisfaction. In particular, he found that adding keywords and categories (i.e., subject headings) to surrogates improved users' sense of comfort and confidence when using the search results. In some cases (particularly when category metadata was employed) time-on-task was also reduced. Chen and Dumais (2000) performed similar usability studies, and also found that adding category (i.e., subject) metadata greatly aided users (see also (Dumais *et al.*, 2001)).

6

Participants in their studies preferred an interface with the additional metadata to one without, and reported much greater subjective satisfaction with the enhanced interface. In addition, task efficiency (the time required to find answers to set questions) was improved when using the enhanced interface.

In short, it appears that providing additional information in surrogates does enable users to make better predictions about document relevance. It seems likely, however, that all metadata elements are not created equal—some may be more effective than others at supporting users' predictive judgments within the context of their situation and task. This perspective leads to the second question:

*Which metadata elements are particularly valuable?*

A long tradition in relevance research is studying which *criteria* people use to assess the relevance of documents. This research has raised questions such as, Do people only care whether the document is "on topic," or are there other factors at play? Can the criteria people use to assess relevance be enumerated and classified? Numerous studies have shown that a finite, relatively *universal set of relevance criteria* exists, and that a *few elements account for the preponderance of use* (see Mizzaro, 1997 and Borlund, 2003 for reviews).

Relevance criteria are important for surrogate design because they can be mapped to metadata elements. For example, an empirical finding that users often refer to the publication date of articles could motivate designers to include publication date as a metadata element in an IR system and display it in surrogates. An example of this approach is Hufford's (1991) study of reference librarians at major university libraries, using both card catalogs and OPACs. He counted the number of times specific elements were used and found that seven elements accounted for 90.7% of total uses. Based on these results, he argued that OPACs should have minimal surrogates with a limited number of elements.

Results of Hufford's study and other related work (e.g. Janes, 1991) are difficult to extend to the Web context, since the Web is a radically different environment. In particular, the Web has a huge and diverse scope, contains full-text documents of greatly varying type and design, and is widely used by people with average education and searching skills (as opposed to search experts such as librarians). Researchers have begun to address this issue by studying how users interact with surrogates provided by Web search engines (Lan, 2002). The research presented here extends this line of inquiry.

*Document characteristics and information usage*

In Phase 2 of IR interaction, users directly evaluate documents. Valuable insights for IR can be gained by studying how users evaluate documents. Because surrogates contain such limited information, they restrict the scope of relevance criteria available to users. Researchers have extended taxonomies of relevance criteria by examining how users interact with full-text documents.

Barry (1994) had students and faculty from various disciplines (e.g., history, English) examine both surrogates and documents from a mediated online search related to a current information need. Participants worked with hard copies, and were asked to circle any part of the document that would lead them to pursue it. This technique of directly eliciting relevance criteria from end users proved effective and has motivated much subsequent research, including this study. Barry's taxonomy of relevance criteria included 23 categories, such as "depth/scope," "affectiveness," and "tangibility." These categories were combined into seven larger groups, such as "criteria pertaining to user's belief and preferences." The results show the wide range of "beyond topical" criteria that users consider when evaluating documents.

Numerous subsequent studies (see, for example, Bateman, 1998; Barry & Schamber, 1998) using a variety of methods have identified additional criteria beyond those in Barry's study. Generally these criteria seem to be similar and relatively universal (Barry & Schamber, 1998; Mizzaro, 1997). Researchers have classified these criteria in different ways. Tang and Solomon (2001) argue for a dual classification: Objective (related to the document) vs. Subjective (related to the user), and Primary (essential for relevance judgment) vs. Secondary (non-essential).

Pre-Web studies generally used fixed document sets based on a predefined query. Interaction in Web IR is typically quite different than this—users' generate their own queries on the fly and enter them directly into search engines. Furthermore, Web documents are much more diverse and complex (including multimedia and hypertextual elements) than the textual articles retrieved from databases in older studies. More recent studies (Lan, 2002; Tombros, Ruthven & Jose, 2005) have extended the relevance criteria approach to Web search. Lan (2002) studied graduate students searching for research related to their dissertations. Analyzing highlighted regions in the documents and interviews, he found users employed a wide range of document characteristics and argued for creating descriptive metadata based on empirical findings. Tombros et al (2005) asked "what features make a Web document useful for information seeking?" and conducted a think-aloud study to discover what features users identified as useful to their searches. Their taxonomy of features included general categories, such as Text, Structure, and Non-textual. They also found a wide range of features used.

Both of these studies found rather different criteria than those reported in earlier studies. Lan (2002) emphasized the importance of characteristics than enable users to filter or reconsider documents when topicality alone is insufficient. Tombros et al (2005)

identified very general categories (different from those in any previous study) that would likely be more useful for IR algorithms than metadata development or surrogate design. These findings suggest that relevance judgment in Web IR interaction is different than in OPAC or database interaction. In order to make effective use of metadata in Web IR, relevance criteria should be refined and classified to reflect users' behavior in Web IR interaction. The research presented here extends this line of inquiry by observing motivated information users evaluating Web documents.

Research in this area is important because it can link IR system design, which is explicit and objective, with users' evaluations of documents, which are tacit and subjective. As Barry (1998) puts it: "clues to relevance are not always explicitly stated by document characteristics" (p. 1301). For example, while a user may seek a document that is easy to understand, there may be no query option for "easy to understand documents," and the surrogates returned by the system may not indicate difficulty either. The user must assess difficulty by reading the document. By directly examining users' tacit and subjective evaluations of documents, empirical research in this area can provide a foundation for more sophisticated system design and better retrieval results.

## Research Questions

The study was designed to extend existing research on relevance judgment and document usage. The following research questions guided the design and analysis of this study:

RQ1) What types of document criteria do Web users (specifically, "health information users") use to assess the relevance of Web documents?

RQ2) When evaluating typical Web search engine (textual) surrogates…

    a. Which document criteria do Web users use to determine whether the documents are relevant to them?

    b. In which parts of the surrogates do Web users most commonly find these criteria?

RQ3) When evaluating complete Web documents…

    c. Which document criteria do Web users use to determine whether the documents are relevant to them?

    d. In which parts of the documents do Web users most commonly find these criteria?

## Methods

### Overview

We designed a combined *user study* (Henderson et al., 1995) and *content analysis* (Barry, 1994; Lan, 2002) to explore what criteria Web users employ to assess the relevance of documents during a Web search. Our goal was to explore how motivated information users—specifically, "health information users"—use various characteristics of documents and surrogates to assess the relevance of Web documents. This user group was chosen for two reasons. First, online health information is one of the most important domains to study, as users are increasingly seeking health information and making health decisions based on health Web sites (Anderson, 2004). Second, our ongoing research partnership with the National Institute of Environmental Health Sciences (NIEHS) gave us access to expertise within this content/task domain.

Health information users were defined as individuals with a strong interest in their health and the ability to use online information retrieval tools to learn more about health-related topics. A health information user was identified as a person comfortable with the Web and highly motivated, but not expert in technical subject areas. Relevance was conceptualized as usefulness to an individual's information need in the context of their background and interests (i.e., "situational relevance" (Borlund, 2003)). This concept was operationalized as users' self-reported judgments of document usefulness.

To meet the goals of the research it was necessary to create a realistic, but controlled, environment in which motivated users could evaluate information as it applied to real tasks and situations (Barry, 1998). In previous studies, researchers have gathered relevance criteria from users via content analysis of highlighted regions (Barry, 1994) and post-search interviews (Rieh, 2002). In this study, naturalism and control were balanced by combining these methods. A user study was first conducted to elicit users' natural usage of document and surrogate characteristics. Unobtrusive content analysis of users' responses was then performed to discover larger patterns and trends.

*Participants*

Participants were recruited from the university and local community via listserv postings, paper flyers, and word of mouth (participants were offered a chance to win an MP3 player as compensation for their time). In order to meet the research goal of studying health information users, recruitment messages requested individuals with a personal interest in a particular environmental health topic. Twelve individuals (nine female) participated. Participants were students and professionals from a variety of fields. The topics selected by participants included an array of environmental health issues, such as environmental justice,

the effects of lead paint on health, and the health consequences of living near power lines

(see Appendix A).

*Procedure*

Participants met individually with a researcher in a private office. First, an informed

consent form was given. Then, the participant completed a profile questionnaire which

recorded demographic information and the environmental health topic of interest. The

researcher then conducted a brief interview to clarify the scope of the participant's

information need, and help the participant transform that need into a query suitable for

execution by a search engine. For example, one participant was concerned about possible

causes of asthma. The interview established that her primary concern was car exhaust, so a

query of "car exhaust asthma" was created. The purpose of the interview and assistance

with query formulation was to help ensure that users entered a query which would return

meaningful results within the specified domain. Pilot testing had shown that users often

came to the study with only vague ideas of what to search for and little sense of the scope of

information provided by NIEHS. Some assistance was therefore necessary to help them

formulate a query which would return useful results.

[ Figure 2 appears about here ]

**Figure 2.  A user's highlighting within a surrogate.**

An example query was used to illustrate the experimental procedure. Participants

were asked to evaluate the first document summary (surrogate) and highlight, using their

mouse, any terms of interest (see Figure 2 and Figure 3).

[ Figure 3 appears about here ]

**Figure 3.  A user's highlighting with a document.**

Then participants were to load the document itself and evaluate it, again highlighting any terms of interest (Lan, 2002). Participants were asked to proceed as they would during an ordinary search, and not to spend more time on evaluation than they would normally. The researcher instructed participants to perform this process for each document in the first page of search results (ten documents)[1]. Video capture software was used to record the entire interaction for later discussion and analysis. After the procedure was explained, the participant was shown the search engine[2] and asked to enter the query and proceed through the results as instructed. The researcher left the room so as to allow participants to work undisturbed.

After participants completed their evaluation of the search results, the researcher conducted a post-search interview (Rieh, 2002), using a subsequent think-aloud protocol (Henderson et al., 1995). As the participant and researcher viewed the video record of the search session together, the researcher prompted the participant to think aloud, describing any reactions to the document or summary. Next, the researcher asked four questions used by (Lan, 2002) about each search result to clarify the participant's relevance judgments and metadata use:

1. Based on the document summary, what was your initial judgment of the document's usefulness to you?

2. Based on the document itself, what was your final judgment of its usefulness to you?

3. What were the factors that led you to that decision?

---

[1] This constraint is consistent with research showing that Web users rarely look beyond the first page or two of search results (Spink, 2003).

[2] For the first three participants, the search engine (Infoseek) on NIEHS website (http://www.niehs.nih.gov/) was used. For the remaining participants, Google (restricted to the NIEHS domain) was used. Though they use different ranking algorithms, the two engines provide essentially identical surrogates with similar user interfaces.

4.  Was there any additional information you would have liked to have had in the
    surrogate?

This conversation was recorded on audiotape and the researcher took notes on the
participant's behaviors, manner, and comments.  After the think-aloud protocol was
completed, the researcher conducted a brief post-interaction interview.  In the interview,
participants were asked for their overall impressions, their satisfaction with the search
experience, and their suggestions for additional information or affordances that could
improve the process for them.  Participants also completed a questionnaire asking for their
ratings of the document summaries, and of the NIEHS Web site as a whole.

*Content analysis*

Video logs were analyzed to determine where in the document (or surrogate) users identified
useful information, as well as what types of information they found most helpful in
evaluating documents (see Figure 4).

[ Figure 4 appears about here ]

**Figure 4.  Example coding of a highlighted region.**

For each highlighted word or phrase, the appropriate document *location* was selected
from the list of possible locations (see Table 1) and coded in a database.  After all
highlighted text had been coded for location, each highlight was coded for the document
*criterion* it represented.  A *criterion* was defined as a clearly distinguishable quality or
characteristic of a document or surrogate.

One or multiple criteria (for as many criteria as were contained within the highlight) were identified and coded in a database. In cases where the nature of the criterion being highlighted was unclear, the think-aloud data and the researcher's notes were consulted to clarify the participant's intent.

[ Table 1 appears about here ]

**Table 1. Locations of highlighted criteria in surrogates and documents.**

[ Table 2 appears about here ]

**Table 2. Document criteria identified by content analysis of users' highlighted regions.**

## *Results*

The content analysis yielded data on participants' usage of criteria and locations within both surrogates and documents. These data address the research questions.

### *Criteria*

RQ1 asked, "What types of document criteria do Web users (specifically, "health information users") use to assess the relevance of Web documents?" The complete set of criteria and categories used for analysis is shown in Table 2.

The list of criteria used for coding was developed based on both previous relevance research (Barry, 1998; Lan, 2002) and on inductive analysis during the coding process. Previous research (Barry, 1998) distinguished between "document criteria" and "source criteria." *Document criteria* include the publication date of the document, the pagination or length of the document and the document type (e.g., a journal article, a book, a dissertation, etc.). *Source criteria* related to the intellectual source of the document, and include authors and/or editors of documents, any organization involved with the creation and/or publication of the document (e.g., universities, institutes, professional organizations that appeared as the author/editor affiliation, the sponsoring agency, or the publisher), and the larger source or publication in which the document appeared (e.g., the journal in which an article appeared, the book in which a chapter appeared).

In this study, "criteria" encompassed both of these aspects. In addition, coding was not limited to previously developed schemes—additional criteria were developed inductively over the course of the analysis. Finally, criteria were clustered into categories based on similarity. This grouping was done only to facilitate interpretation of the results and did not affect the coding process in any way.

*Criterion selection*

[ Figure 5 appears about here ]

**Figure 5. Distribution of highlighted criteria in surrogates and documents, ordered by criterion.**

RQ2a asked, "When evaluating typical Web search engine (textual) surrogates, which document criteria do Web users use to determine whether the documents are relevant to them?" RQ3a asked, "When evaluating Web documents, which document criteria do Web

users use to determine whether the documents are relevant to them?" [ Figure 5 appears

about here ]

Figure 5 shows the highlighted criteria and the number of times each was highlighted, in

both the surrogates and the documents. [ Figure 6 appears about here ]

Figure 6 shows the highlighted criteria grouped by category, and the number of times

each was highlighted, in both the surrogates and the documents.

A total of 555 highlighted regions (see [ Figure 5 appears about here ]

Figure 5 and [ Figure 6 appears about here ]

Figure 6) were captured in the video logs of the 120 documents examined by users. Content

analysis of these highlighted regions (see Figure 4 for an example) identified 788 criteria.

There were more criteria than regions because the coding procedure allowed for multiple

criteria to be assigned to a single highlighted region. Participants highlighted topical criteria

most often, but a range of criteria were also commonly identified. Researchers often test

whether distributions of information creation or information use can be characterized by

summary laws or formulae (e.g., Newby, Greenberg & Jones, 2003; Bates, 1996). Studying

criteria this way helps to place users' behavior in context, and helps to raise questions for

further research and analysis. To examine whether the distribution of criteria observed here

could be similarly characterized, we fit the data to an exponential curve. The data

approximately fit this exponential distribution ($y = 134.67e^{-0.1568x}$; $R^2 = 0.90$). This result

reflects the relative preponderance of a few widely-used criteria in conjunction with a long

tail of occasionally-used criteria. In addition, there were different patterns of highlighted

criteria in the surrogates and the documents (see [ Figure 7 appears about here ]

Figure 7 and [ Figure 8 appears about here ]

Figure 8). A broader range of criteria were highlighted in the documents than in the surrogates—but while some criteria (e.g. purpose or goal) were identified repeatedly in both the surrogate and the document, others (e.g. investigator) were identified almost exclusively in either surrogates or documents.

[ Figure 6 appears about here ]

**Figure 6. Distribution of highlighted criteria in surrogates and documents,**

**ordered by criterion category.**

[ Figure 7 appears about here ]

**Figure 7. Relative appearance of highlighted criteria in surrogate vs. document.**

[ Figure 8 appears about here ]

**Figure 8. Relative appearance of highlighted categories of criteria in surrogates vs. documents.**

Considering the distribution of criteria by categories (see [ Figure 6 appears about here ] Figure 6) further illustrates the diversity of criterion usage—six of the eight categories represented at least 5% of the total highlighted regions. There were likewise differences in appearance of criterion categories in the documents versus the surrogates (see [ Figure 8 appears about here ] Figure 8). Criteria in the topical relevance category were highlighted often in both document and surrogates, while criteria in the authority/person category were almost always highlighted in the document.

*Criterion location*

RQ2b asked, "When evaluating typical Web search engine (textual) surrogates, in which parts of the surrogates do Web users most commonly find these criteria?" RQ3b asked, "When evaluating Web documents, in which parts of the documents do Web users most commonly find these criteria?" [ Figure 9 appears about here ]

Figure 9 shows the distribution of criterion locations in the surrogates; Figure 10 shows the distribution of criterion locations in the documents.

Of the 555 highlighted regions, 408 were in the documents, while 147 were in the surrogates.

In the surrogates, virtually all criteria identified in the content analysis were highlighted within the description (see [ Figure 9 appears about here ]

Figure 9). In the documents, criteria were identified in multiple locations (see Figure 10). The most common location, paragraph text, accounted for the bulk of the highlights. Secondary aspects of document structure, such as titles, section headings, lists and hyperlinks were also used. While specific structural locations were only highlighted occasionally, when considered in the aggregate, structural features in general were commonly used.

[ Figure 9 appears about here ]

**Figure 9. Distribution of criterion locations in the surrogates.**

[ Figure 10 appears about here ]

**Figure 10. Distribution of criterion locations in the documents.**

[Table 3 appears about here]

**Table 3.  Consistency of relevance judgments.**

*Relevance judgments*

Table 3 shows initial relevance judgments (based on the document surrogates) and final relevance judgments (based on the full documents).  There were 99 unique judgments from 11 participants[3].  Participants often changed their initial judgment—59.6% of the initial judgments were consistent with the final judgment, as opposed to 40.4% that were inconsistent.  The null hypothesis (that the observed proportion of 59.6% consistent judgments could have occurred by chance alone) cannot be rejected ($p > 0.07$).  In other words, the observed proportion of shifts in relevance judgment is consistent with a model of relevance judgment as a random binomial process—participants' initial judgments were no better than chance at predicting their final judgments.  Most (65.0%) of the changes in assessed relevance were shifts from relevant to not relevant, while 35.0% were shifts from not relevant to relevant.

## Discussion

The results indicate that relevance judgments of Web users are complex and multifaceted, drawing on a range of document criteria and locations.  These findings have important implications for IR system design.

*Document characteristics and metadata usage*

**Criterion types and categories.**  As [ Figure 5 appears about here ]

---

[3] One participant's relevance data had to be dropped to difficulty that participant had in completing the procedure properly.  A few judgments from other participants were dropped due to unclear judgments or data collection problems.

Figure 5 illustrates, participants employed diverse criteria in their document evaluation processes. As anticipated, the most frequently-identified criteria were topical ones—but these did not constitute even a majority of identified criteria. Consistent with previous research (Barry, 1994, 1998; Lan, 2002), this study found that Web users rely on an array of 'extra-topical' document characteristics when making relevance judgments. Not all characteristics were identified equally, however—a decaying exponential or power-law path was observed, in which a few criteria were heavily used, followed by long tail of occasionally-used criteria.

This distribution is consistent with well-known long-tail distributions in large information spaces, such as the Zipf and Bradford distribution (Bates, 1996, 2002). This finding implies that designers of retrieval systems should focus on the most important criteria or aspects of documents, and prioritize these in the systems' information architecture and user interface (see Brinck et al, 2002). According to the polyrepresentation model (Ingwersen, 1996), IR systems should provide a diverse set of "roads to information" (p. 23). The data presented here suggest a corollary to this principle—the information elements used to facilitate retrieval and to allow users to evaluate documents should be structured and prioritized to match users' tasks, knowledge, and cognitive abilities.

For example, the health information users studied here were shown to rely heavily on "scoping" criteria, such as the intended audience of a document. An IR system intended to support this user group should take this into account. This might be done by extracting scope criteria automatically, or using explicit scope metadata. In either case, the scope criteria could be incorporated into the retrieval mechanism (weighting consumer-friendly documents more highly, for example), or into the user interface (allowing users to sort retrieved documents by their intended audience). The research method presented here can

be applied in specific organizational and situational contexts to determine which criteria are most critical to the target user group—an approach consistent with the "socio-cognitive" conception of relevance (Cosjin & Ingwersen, 2000).

Considering the distribution of highlighted criteria by categories (see [ Figure 6

appears about here ]

Figure 6), further patterns are apparent. Topical criteria were commonly used—but not overwhelmingly so. This suggests that users relied on subject information, but needed additional information. In some cases, this information might have served to narrow or restrict overly broad search results. The large number of highlighted criteria in the "scope" category supports this view—users were evidently looking for cues that could connect a document from a general search to their particular problem or interest. In other cases, the "extra-topical" information might have augmented or complicated a particular topic.

"Research" criteria, relating to research design, variables studied, and so forth, were the most commonly-highlighted criteria. The importance of these criteria reflects the user group ("health information users") that were studied, but also a general increase in sophistication on the part of Web users. Even though the users in this study were not environmental health experts, they still sought detailed information about data collection and research methods. This strategy implies that they wanted to analyze the assumptions behind documents, rather than taking the available conclusions and recommendations at face value.

At the same time, users were not seen to be interested in assessing the credibility or authority of retrieved documents. None of the highlighted regions directly reflected credibility issues, nor did any users comment on credibility, and the authority/person category was the least commonly-highlighted group. This tension between detailed

23

investigation of research methods and limitations on the one hand, and implicit acceptance

of the credibility and authority of Web documents on the other, deserves further attention.

**Criterion location.** Analysis of the location of highlighted document characteristics

revealed that participants relied primarily on plain text. In the surrogate, they

overwhelmingly used the document description (automatic extract). In the document, they

generally used the text itself. Still, roughly one-third of the highlighted locations in the

document included some structural feature, such as a section heading or hyperlink. The

finding indicates that users did find structural features useful in some cases.

The extent to which users employ such structural features is highly dependent on

document design. Ideally, Web documents should be highly structured and designed to

support easy scannability (Morkes & Nielsen, 1997), so users don't have to rely on raw text.

Highly scannable documents might enable users to assess documents more rapidly and

accurately, improving relevance judgments (see below). In addition, the structural features

could be exploited by retrieval systems. Link analysis (Yang, 2005) and tag analysis (as in

Google's exploitation of HTML structure), for example, have proven useful in retrieval of

unstructured Web documents. Researchers have also begun to investigate the use of

structured documents such as XML in retrieval, as in the INEX project (Fuhr *et al.*, 2004).

Findings from user-defined Web relevance research can inform these efforts by suggesting

what types of document structure would be most helpful to users, and how systems can

present this structure in user interfaces.

The surrogate data show a greatly restricted range of highlighted criteria in comparison to

the documents (see [ Figure 7 appears about here ]

Figure 7 and [ Figure 8 appears about here ]

Figure 8).  For example, some criteria (e.g., investigator, hyperlink) were only highlighted in documents, never in surrogates.  The relative paucity of information in the surrogates may be one reason why surrogates were of limited effectiveness in supporting relevance judgments.  If the criteria that users rely on to augment or restrict topical relevance are not present in a surrogate, it is unsurprising that such a surrogate is of little use for a user looking beyond topical relevance.  Other studies (Drori 2000a, b; Dumais, Cutrell & Chen, 2001) have demonstrated that adding subject metadata such as categories and keywords to surrogates improves the usability of Web search engines.  Future research should seek to extend this work by exploring how extra-topical criteria can be effectively integrated into surrogates.

**Relevance on the web.**  Participants' initial relevance judgments (based only on the surrogate) were found to be no better than chance.   Users were sometimes disappointed by the promise of a potentially relevant document that turned out to be not relevant.  At other times, users made serendipitous discoveries—some users initially rated a document as not relevant, then found that it was actually relevant.  In contrast, earlier (pre-Web) work on surrogates found relevance judgments to be accurate 70% - 90% of the time (Barry, 1998).  In addition, participants gave the surrogates a mean score of 3.6 on the "informative" scale (a 7-point semantic differential scale, where 1 indicated "informative" and 7 indicated "uninformative") and a mean score of 3.8 on the "clear" scale (a 7-point semantic differential scale, where 1 indicated "clear" and 7 indicated "confusing").  These ratings suggest that participants were neutral regarding the informativeness and clarity of the surrogates they used.

These findings imply that Web surrogates have substantial room for improvement. Designers should make better use of specific document elements (based on empirical studies and user suggestions).  They also need better *quality* metadata, particularly descriptions, since users rely on these heavily.

## Conclusion

This study focused on the relevance judgments of health information users using the Web to search for information on how an environmental issue affects their health.  Users were asked to identify the criteria they use to assess relevance in both Web search surrogates and full-text documents.  Analysis of users' responses provided a summary of the criteria

these users rely on, with implications for theories of relevance behavior, the methodology of related studies, and IR practice.

*Theoretical implications*

This research found that Web users employ a wide range of criteria in documents and surrogates as they assess the relevance of retrieved documents. These findings from the Web environment support a major hypothesis of earlier relevance research—that users rely on many criteria in addition to topicality (Mizzaro, 1997). It was noted that users employ these criteria both to *augment* topical relevance and to *filter* on-topic documents. Augmenting topical relevance involves using additional relevance criteria to assess the relevance of a document when topicality is complex or unclear. For example, a user might be unsure that a particular document is on topic, but notes that it is very recent and easy to read, and so decides to save it for further evaluation. Filtering on-topic documents involves using relevance criteria to eliminate documents that are not pertinent or useful (Lan, 2002). For example, a user might find a research article on her topic that was published twenty years ago, and has been superseded by more recent findings.

This study also supports Barry's (1998) argument that the criteria identified by users are closely tied to users' situations and information needs. The user group studied here— "health information users"—differs from previous studies, which largely focused on academic users. The different perspective provided by these users enabled the identification of additional relevance criteria not found in previous studies. For example, users were often concerned with the practical implications or impact of the information in a document. On the other hand, criteria that were important in other studies, such as "affectiveness" or "access" (Barry, 1998) were not identified in this research. So the theoretical implication of

studying additional user groups is *not* continual expansion of relevance criteria taxonomies—this search for a "universal" set of criteria will never be complete. Rather, it should be acknowledged that users will employ criteria closely tied to their situations and information needs. In the end, relevance evaluation is a subjective and individual process (Harter, 1992), so the goal of research should be to identify clear patterns than can inform system design in specific contexts, rather than to seek generalities. In this study, for example, a clear pattern was that health information users sought practical and actionable information about their health concerns. Providers of health information should develop IR systems that better deliver such information (e.g., a summarization system targeted at lay users (Kan & Klavans, 2002)).

A further theoretical consideration is the extent to which the "length hypothesis" is relevant on the Web. That is, do longer Web surrogates, containing more information, enable users to make better relevance judgments? Our study does not directly address this question. It is apparent, though, that both the breadth (exhaustivity) and depth (specificity) of relevance criteria available to users should be considered. Do Web users need a *wider range* of clues (breadth) or more *detailed and comprehensive* clues (depth)? This question, too, is context-specific, but empirical studies can illuminate it. Our data suggest, for example, that health information users would benefit from a greater breadth of clues in Web surrogates.

*Methodological implications*

Relevance judgments are individual, subjective, and closely bound with contextual factors, including the user's task, environment, and existing knowledge. Instead of seeking general sets of relevance criteria based on psychological models, researchers should study relevance behavior in context and with system design in mind. Numerous methods for

studying relevance behavior and criteria have been developed in previous studies. The methodological contribution of this study is to effectively combine approaches that have been used in previous studies and successfully adapted by multiple independent researchers:

- elicitation of relevance criteria using highlighting (Barry, 1994)

- direct interaction with Web pages using a mouse to highlight regions (Lan, 2002)

- content analysis of user-highlighted regions (Barry, 1994, 1998)

- post-search interviews using a video record (Rieh, 2002)

- examination of relevance judgment as a two-phase process (Tang & Solomon, 2001)

Future research can gather richer data by combining these techniques, as presented here.

*Practical implications*

The key goal of this research was to support the design of more effective IR systems for the Web by providing a better understanding of how users assess documents when searching a Web site. Our data show that health information users rely on a wide range of criteria when evaluating the relevance of Web documents. In contrast, Web IR systems tend to incorporate only a few criteria. The algorithms used to rank documents are based on term and hyperlink frequencies. The surrogates provided to users show little more than a title (which may contribute little to relevance judgment) and an automatically extracted snippet of the document text.

System designers should consider customizing IR systems to particular contexts and user groups—such as the health information users discussed here—by mapping identified user needs to IR algorithms and user interfaces. In this case, health information users were seen to have frequently used criteria related to "methodology" and "scope." System design could reflect this user preference by more explicitly incorporating these criteria, in a variety of

ways. First, *surrogates* could be expanded to include more clues to these criteria. For example, a surrogate for a document discussing a study might state the sample and types of data collected. Second, *filtering mechanisms* could be provided to enable users to more effectively specify queries. For example, users could limit their searches to documents written for the general public, rather than for scientists. This functionality could be provided either in the query interface (parametric searching) or in the search results interface (as in "facet-based" systems such as FLAMENCO (Yee et al, 2003) and RB++ (Zhang & Marchionini, 2004)). Third, *relevance ranking algorithms* could be designed to rank documents with "relevant" methodology or scope more highly. This might require developing a user profile that would allow matching of document methodology or scope with a user's preference.

All of these approaches would entail the creation of additional metadata, compared to what is currently used by Web search engines. This might be done manually, by having resource authors or catalogers assign appropriate terms from a controlled vocabulary to appropriate metadata elements (e.g., "technical" or "scientific" to indicate a document aimed at scientists). Moreover, systems incorporating automatic metadata creation functionalities might be employed (Greenberg, Spurgin & Crystal, 2005). In any case, creating metadata is costly, so organizations should seek to prioritize its development based on user needs. By using empirical research to identify the most important areas for metadata support, organizations can more *efficiently allocate limited metadata creation resources.*

*Limitations and future research*

As discussed above, this research emphasized the intensive study of a particular class of users—"health information users"—so as to understand the behavior of motivated

30

participants with realistic information needs. Generalizability was not a key a goal of this research—rather, our emphasis was on domain focus and contextualization.

In addition, the procedure was somewhat rigid and therefore may not fully reflect individuals' natural searching behavior. A useful direction for future work would be the development of fully unobtrusive techniques that could provide insight into users' natural behavior (ideally, over longer periods of time as well). A subtler issue is that the method used here only captures the *frequency* of criteria highlighted by users. Frequency cannot be directly equated with importance (Barry, 1994)—it is conceivable that some criteria are used only occasionally, but are extremely important when they are used. Still, it seems logical that the criteria users identify most frequently are of some utility (Tombros, Ruthven & Jose, 2005). Another opportunity for methodological innovation in this area, then, is the development of techniques for measuring and comparing the importance of different relevance criteria.

Future work could extend this study by studying additional classes of users in different domains. The studies should seek to map particular user groups and their tasks to specific metadata requirements, which can improve system and interface design. This *user-centered approach to metadata* has the potential to improve IR systems and interfaces through the three techniques outlined above: surrogates, filtering mechanisms, and ranking algorithms.

### Acknowledgments

## *Appendix A*

Queries entered by users:

- environmental justice

- "environmental justice"

- household cleaning products

- health effects power lines

- cell phone cancer

- pesticide risk developing countries

- paint consumer safety

- lead paint health effects

- "wood dust" cancer

- mold -asthma

- "lead poisoning"

- asthma car exhaust

## *References*

Anderson, J. G. (2004). Consumers of e-health: Patterns of use and barriers. *Social Science Computer Review, 22*(2), 242-248.

Barry, C. L. (1994). User-defined relevance criteria: an exploratory study. *Journal of the American Society for Information Science, 45*, 149 - 159.

Barry, C. L. (1998). Document representations and clues to document relevance. *Journal of the American Society for Information Science, 49*(14), 1293-1303.

Bateman, J. (1998). Changes in relevance criteria: a longitudinal study. *Proceedings of the American Society for Information Science Annual Meeting*, 35, 23-32.

Bates, M. J. (1996). Document familiarity in relation to relevance, information retrieval theory, and Bradford's law: The Getty online searching project report No. 5. *Information Processing & Management, 32*, 697 – 707.

Bates, M. J. (2002). After the dot-bomb: Getting Web information retrieval right this time. *First Monday, 7*(7), retrieved from: http://firstmonday.dk/issues/issue7_7/bates/index.html .

Borlund, P. (2003). The concept of relevance in IR. *Journal of the American Society for Information Science and Technology*, *54*(10), 913 - 925.

Carlyle, A. (1999). User categorisation of works: Toward improved organisation of online catalogue displays. *Journal of Documentation, 55*(2), 184 - 208.

Chen, H., & Dumais, S. (2000). Bringing order to the web: Automatically categorizing search results. In *Proceedings of the ACM SIGCHI Conference on Human Factors in Computing Systems (CHI '00)*, The Hague, The Netherlands, April 1 – 6. New York: ACM Press pp. 145 - 152.

Cosjin, E., & Ingwersen, P. (2000). Dimensions of relevance. *Information Processing & Management, 36*, 533 - 550.

Drori, O. (2000a). The benefits of displaying additional internal document information on

    textual database search result lists. In J. L. Borbinha & T. Baker (Eds.), *Proceedings of*

    *Research and Advanced Technology for Digital Libraries, 4th European Conference, ECDL*

    *2000* (pp. 69 - 82). Lisbon, Portugal: Springer.

Drori, O. (2000b). Improving display of search results in information retrieval systems –

    users' study.  (Technical Report No. 200034). Jerusalem, Israel: Leibniz Center for

    Research in Computer Science.  Retrieved December 7, 2003 from

    http://shum.huji.ac.il/~offerd/papers/drori082000.htm .

Dumais, S., Cutrell, E., & Chen, H. (2001). Optimizing search by showing results in context.

    In *Proceedings of the ACM SIGCHI Conference on Human Factors in Computing Systems*

    *(CHI '01)*, Seattle, WA, USA, March 31 - April 5.  New York: ACM Press, pp. 277 –

    284.

Fuhr, N., Malik, S., & Lalmas, M. (2004). Overview of the INitiative for the Evaluation of

    XML retrieval (INEX) 2003. In *Proceedings of the second INEX workshop, Dagstuhl,*

    *Germany, December 15--17, 2003* (pp. 1 - 11).

Greenberg, J., Spurgin, K. & Crystal, A. (2005). Functionalities for automatic metadata

    generation applications: a survey of metadata experts' opinions. *International Journal of*

    *Metadata, Semantics & Ontologies, 1*(1).

Harter, S. (1992). Psychological relevance and information science. *Journal of the American*

    *Society for Information Science, 43*(9), 602 - 615.

Henderson, R., Podd, J., Smith, M., & Varela-Alvarez, H. (1995). An examination of four

    user-based software evaluation methods. *Interacting with Computers, 7*(4), 412 - 432.

Hufford, J. (1991). Elements of the bibliographic record used by reference staff members.

    *College & Research Libraries, 52*(1), 54-64.

Janes, J. W. (1991). Relevance judgments and the incremental presentation of document

representations. *Information Processing & Management, 27*(6), 629 - 646.

Kan, M. & Klavans, J. (2002).  Using librarian techniques in automatic text summarization

for information retrieval.  In *Proceedings of ACM/IEEE Joint Conference on Digital

Libraries (JCDL '02)*, Portland, Oregon, USA, June 14-18.  New York: ACM Press.

Kekäläinen, J., & Järvelin, K. (2002). Using graded relevance assessments in IR evaluation.

*Journal of the American Society for Information Science and Technology, 53*(13), 1120 - 1129.

Lan, W.-C. (2002). *From document clues to descriptive metadata.* Unpublished Dissertation,

University of North Carolina at Chapel Hill, Chapel Hill, NC.

Marchionini, G. (1995). *Information seeking in electronic environments.* Cambridge; New York:

Cambridge University Press.

Mizzaro, S. (1997). Relevance: the whole history. *Journal of the American Society for Information

Science, 48*(9), 810 - 832.

Morkes, J. & Nielsen, J. (1997). Concise, scannable, and objective: How to write for the

Web. Unpublished report, retrieved July 20, 2005, from

http://www.useit.com/papers/webwriting/writing.html .

Newby G., Greenberg J., & Jones P. (2003).  Open source software development and

Lotka's Law: Bibliometric patterns in programming.  *Journal Of The American Society

For Information Science And Technology 54*(2), 169 – 178.

Rieh, S. Y. (2002). Judgment of information quality and cognitive authority in the web.

*Journal of the American Society for Information Science and Technology, 53*(2), 145-161.

Schamber, L. (1994). Relevance and information behavior. *Annual Review of Information Science

and Technology, 29*, 3-48.

Schamber, L., Eisenberg, M. & Nilan, M. (1990). A re-examination of relevance: Toward a

dynamic, situational definition. *Information Processing & Management, 26*(6), 755 – 776.

Spink, A. (2003). Web search: emerging patterns. *Library Trends, 52*(2), 299 - 306.

Sugar, W. (1995). User-centered perspective of information retrieval research and analysis

methods. *Annual Review of Information Science and Technology, 30.*

Tang, P. and Solomon, P. (2001). Use of relevance criteria across stages of document

evaluation: On the complementarity of experimental and naturalistic studies. *Journal*

*of the American Society for Information Science and Technology, 52*(8), 676 – 685.

Tombros, A., Ruthven, I., & Jose, J. M. (2005). How users assess web pages for information

seeking. *Journal of the American Society for Information Science and Technology, 56*(4), 327 -

344.


Yang, K. (2005). Information retrieval on the Web. *Annual Review of Information Science and*

*Technology, 39.*

Yee, K. P, Swearingen, K., Li, K., Hearst, M. (2003). Faceted metadata for image search and

browsing. In *Proceedings of the ACM SIGCHI Conference on Human Factors in Computing*

*Systems (CHI '03)*, Fort Lauderdale, FL, USA, April 5 – 10. New York: ACM Press

pp. 401 – 408.

Zhang, J. & Marchionini, G. (2004). Coupling browse and search in highly interactive user

interfaces: a study of the relation browser++. In *Proceedings of the ACM/IEEE Joint*

*Conference on Digital Libraries (JCDL '04)*, Tuscon, AZ, USA, June 7-11. New York: ACM

Press, p. 384.

| Location | Description |
|---|---|
| In document summary | |
| Title | The title of the page as displayed in the search results (extracted from the <TITLE> tag of the document). |
| Description | The summary (automatic extract) provided by the search engine. |
| Date | The date the page was last updated. |
| Uniform Resource Locator (URL) | The URL (Web address) of the page. |
| Format | The type of document, such as HTML, PDF, or Word. |
| In the document itself | |
| Title | The title of the page, as written in the actual text of the document (not in the <TITLE> tag). |
| Section heading | A heading indicating a major or minor section of the document, as in an outline. |
| Paragraph text | Ordinary text with no special attributes. |
| Emphasized text | Text that was highlighted, bolded, italicized or otherwise emphasized. |
| Image | Any figure, photograph or illustration. |
| Hyperlink | A link (usually underlined, though this varied |

| | by document) to another Web page. |
|---|---|
| Navigation | Links in a specific area of the page dedicated to providing navigation through a larger Web site. |
| List item | Text in a bulleted or numbered list. |
| Citation or reference | Text that referred to another source, such as a journal article or book. |

Table 1.

| Criterion | Description |
|---|---|
| *Affiliation category* | *Criteria relating a document to a particular organization, event or initiative.* |
| Event or meeting | Reference to a specific event, such as a scientific conference. |
| Institution | Reference to a specific institution, such as the one that conducted a particular study. |
| Program or function | Reference to a specific project or other group within an institution. |
| *Authority/person category* | *Criteria relating a document to a particular authority, such as a person or work.* |
| Author or editor | Name or other identifier of an author or editor of the document. |
| Investigator | The scientists responsible for a study. |
| Reference or citation | Explicit reference to another source, such as a journal article or book. |
| *Data category* | *Criteria that constitute specific pieces of raw data, without interpretation or analysis.* |
| Data | Specific results from a scientific study. |
| Fact | A statement of fact, typically based on scientific research. |
| *Influence category* | *Criteria having to do with the larger effects of a* |

| | |
|---|---|
| | *research project or other initiative.* |
| Importance or impact | The salience of a given topic; for example, its importance to human health. |
| Practical advice or implication | Recommendations or discussion relevant to practical problems or everyday living. |
| Purpose or goal | Statement of the objectives of the research or other activity discussed in the document. |
| *Methodology category* | *Criteria related to the research design or methodological position of a study.* |
| General approach | The overall perspective the document represented, such as laboratory experimentation, or field interviews. |
| Method or design | Discussion of a specific research design. |
| Theory or model | Specific reference to a theoretical position that guided research. |
| Variable or factor | Specific reference to a particular (typically causal) variable in a study, or factor in a discussion of health implications. |
| *Scope category* | *Criteria representing a restriction of document focus to a specific context.* |
| Geographic coverage | Discussion of a specific region or area. |
| Level or audience | Who the document was written for (e.g. elementary-school children; policymakers). |

| | |
|---|---|
| Time period | Mention of a specific time period, typically related to the content of the document (e.g. when certain research was performed). |
| *Topical relevance category* | *Criteria used to determine if a document is on- or off-topic for a particular user and query.* |
| Irrelevance filter | A criterion used specifically to determine that the document was not relevant to the user. |
| Topic or subject | Mention of the general topic discussed in the document. |
| *Web characteristics category* | *Criteria related to the structure and format of the document as a Web page.* |
| Document type | Description of what type of document the Web page was, such as a research report, or an informational brochure. |
| Hyperlink | A link to another Web page. |

Table 2.

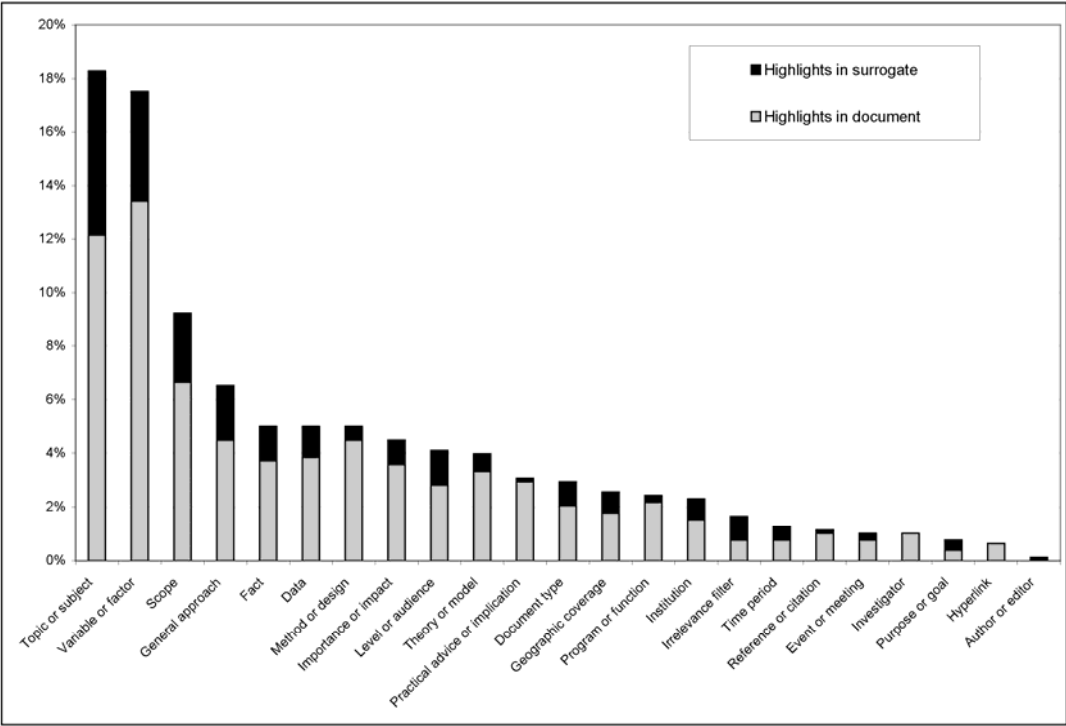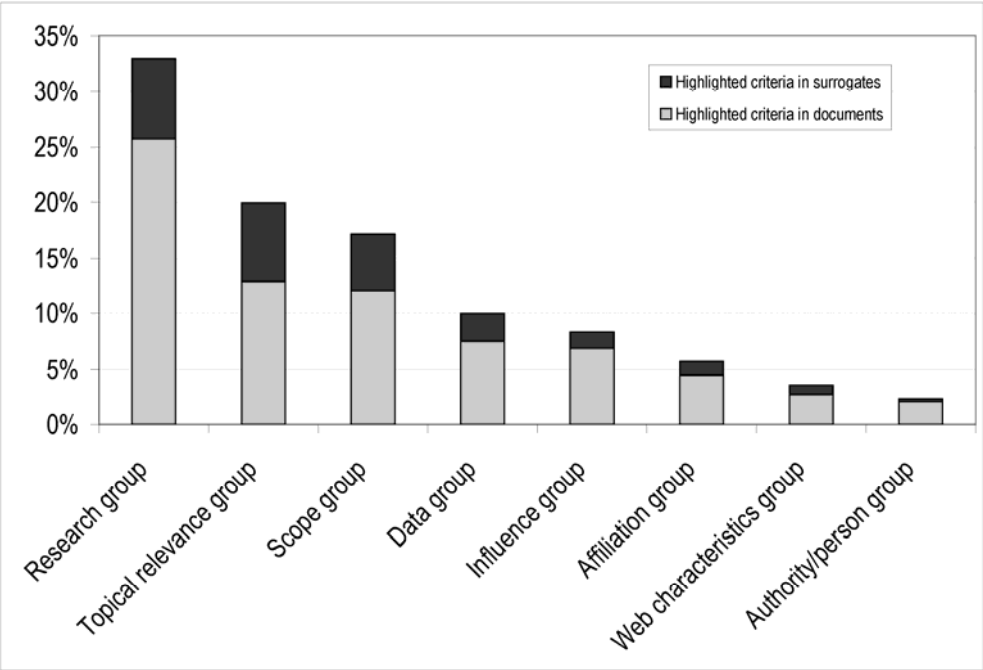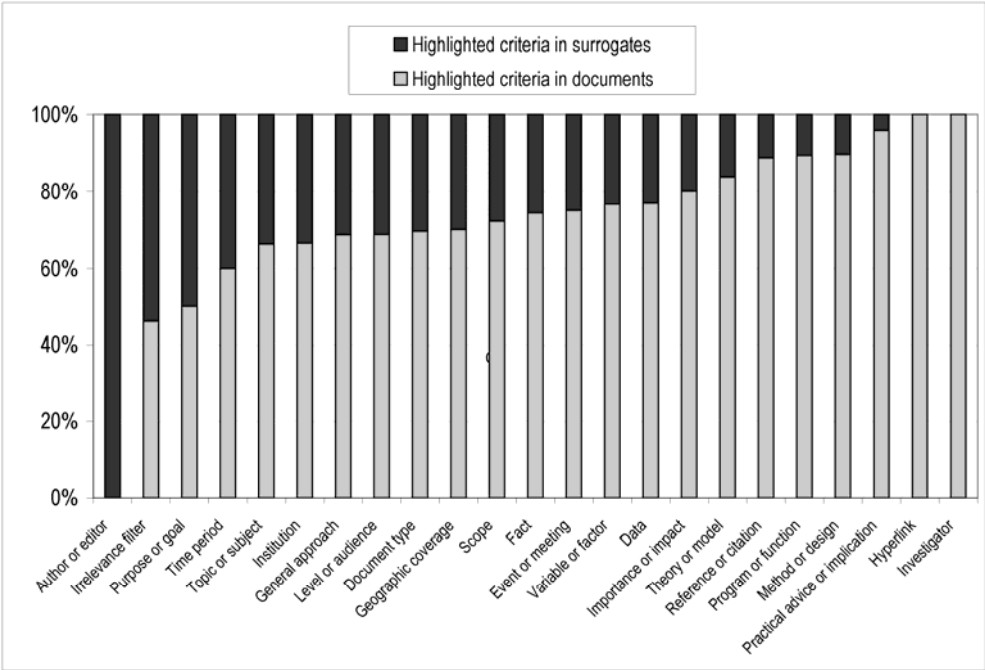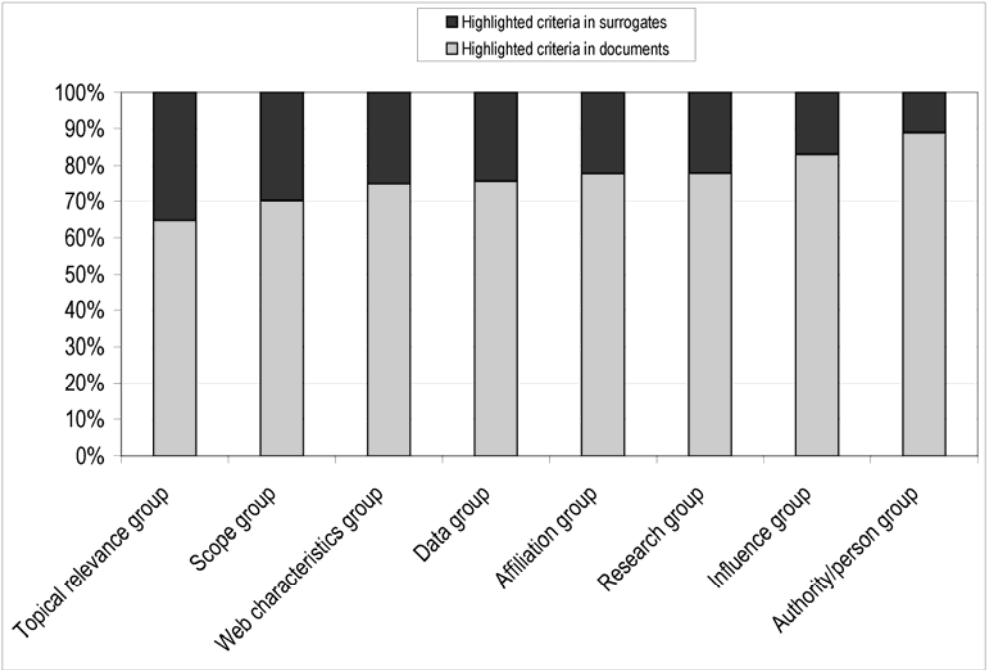| Initial (surrogate) judgment | | ➔ Final (document) judgment | |
|---|---|---|---|
| Relevant | 61.6% | Relevant (consistent) | 57.4% |
| | | Not relevant | 42.6% |
| Not relevant | 38.4% | Relevant | 63.2% |
| | | Not relevant (consistent) | 36.8% |

Table 3.

Figure 1.

Figure 2.

Figure 3.

Figure 4.

Figure 5.

Figure 6.

Figure 7.

Figure 8.

Figure 9.

Figure 10.