



## Marketing Science

Publication details, including instructions for authors and subscription information:  
<http://pubsonline.informs.org>

### The Little Engines That Could: Modeling the Performance of World Wide Web Search Engines

Eric T. Bradlow, David C. Schmittlein,

To cite this article:

Eric T. Bradlow, David C. Schmittlein, (2000) The Little Engines That Could: Modeling the Performance of World Wide Web Search Engines. Marketing Science 19(1):43-62. <https://doi.org/10.1287/mksc.19.1.43.15180>

Full terms and conditions of use: <https://pubsonline.informs.org/Publications/Librarians-Portal/PubsOnLine-Terms-and-Conditions>

This article may be used only for the purposes of research, teaching, and/or private study. Commercial use or systematic downloading (by robots or other automatic processes) is prohibited without explicit Publisher approval, unless otherwise noted. For more information, contact [permissions@informs.org](mailto:permissions@informs.org).

The Publisher does not warrant or guarantee the article's accuracy, completeness, merchantability, fitness for a particular purpose, or non-infringement. Descriptions of, or references to, products or publications, or inclusion of an advertisement in this article, neither constitutes nor implies a guarantee, endorsement, or support of claims made of that product, publication, or service.

© 2000 INFORMS

Please scroll down for article—it is on subsequent pages



With 12,500 members from nearly 90 countries, INFORMS is the largest international association of operations research (O.R.) and analytics professionals and students. INFORMS provides unique networking and learning opportunities for individual professionals, and organizations of all types and sizes, to better understand and use O.R. and analytics tools and methods to transform strategic visions and achieve better outcomes.

For more information on INFORMS, its publications, membership, or meetings visit <http://www.informs.org>

# The Little Engines That Could: Modeling the Performance of World Wide Web Search Engines

Eric T. Bradlow • David C. Schmittlein

*The Wharton School, University of Pennsylvania, Philadelphia, Pennsylvania 19104-6371*

*ebradlow@wharton.upenn.edu • schmittlein@wharton.upenn.edu*

## Abstract

This research examines the ability of six popular Web search engines, individually and collectively, to locate Web pages containing common marketing/management phrases. We propose and validate a model for search engine performance that is able to represent key patterns of coverage and overlap among the engines.

The model enables us to estimate the typical additional benefit of using multiple search engines, depending on the particular set of engines being considered. It also provides an estimate of the number of relevant Web pages *not* found by any of the engines. For a typical marketing/management

phrase we estimate that the “best” search engine locates about 50% of the pages, and all six engines together find about 90% of the total.

The model is also used to examine how properties of a Web page and characteristics of a phrase affect the probability that a given search engine will find a given page. For example, we find that the number of Web page links increases the prospect that each of the six search engines will find it. Finally, we summarize the relationship between major structural characteristics of a search engine and its performance in locating relevant Web pages.

*(Capture/Recapture; Hierarchical Bayes; Marketing Information; Probability Models; World Wide Web)*

## 1. Introduction

The World Wide Web (WWW) is important to managers in three rather different respects. First, managers use it to engage in electronic commercial transactions as sellers or as buyers (Alba et al. 1997, Hoffman et al. 1996). Second, they use it to disseminate information to customers or gather information as (business) customers, including both Web advertising (acquiring new customers) and after-sales support to retain customers (Bakos 1997, Burke 1996, Hoffman and Novak 1996). Third, the Web is emerging as a rich source of managerial information that assists in decision-making, e.g., competitive intelligence, demographic information, market forecasts, general economic information, sources of external expertise or training, innovative managerial tools, tactics and strategies, and regulatory and other governmental information. Providers of such information include news organizations, governments, educational institutions, corporations, and nonprofit organizations, etc. Web search engines are commonly used to help locate this kind of information, and it is this performance of such engines that interests us here.

Search engine performance has begun to attract attention by both researchers and managers. Selberg and Etzioni (1996) studied search queries and their results using various popular search engines for the period July through September 1995. In a more recent and comprehensive study published in *Science*, Lawrence and Giles (1998) examined the URLs returned for a large number of queries during December 1997. A follow-up to that study, using more comprehensive search methods, a greater number of engines, and a larger number of phrase queries, has recently appeared in *Nature* for queries collected in February 1999. They were particularly interested in the relative number of URLs returned by different search engines and in estimating the number of URLs *not* found by any (or all) search engines. Coverage of those findings in *The Wall Street Journal* (1998) showed both the managerial interest and also the controversy generated by the findings. With significant advertising revenue at stake, those responsible for the engines are sensitive to assessments of their relative performance. Indeed, such assessments have loomed large in the business press discussion of the vast sums paid to acquire search engine sites.

In this study we will offer the following contributions. First, we present and validate a model for the performance of multiple Web search engines in finding URLs. We also analyze some natural, relatively simple models (Rasch-type ability/difficulty model and capture/recapture model) and find that they fail to represent key aspects of search engine performance (which the proposed model does contain). Second, we analyze the performance of six popular Web search engines in finding marketing/management phrases. Selberg and Etzioni (1996) studied all queries submitted to MetaCrawler, and Lawrence and Giles (1998, 1999) examined queries from the scientists at the NEC Research Institute in Princeton. Neither focused on management information. Third, we show how some characteristics of marketing/management phrases and of Web pages/URLs affect search engine performance. We also highlight the association between structural characteristics of a search engine (e.g., size of universe covered, depth of search) and that engine's success. Fourth, our empirical model application allows us to do more than just "rate the search engines," enabling us to describe the distinctive patterns of overlap and distinctiveness among them. Finally, for these kinds of management phrases, we are able to estimate the number of URLs *not* found by individual search engines, and indeed not by the collection of engines. We also can calculate the incremental benefit in adding a particular search engine's results to those URLs already found. The next section offers a description of the search process and search outcomes, some summary statistics regarding search engine performance, and a conceptualization of factors thought to affect that performance. The subsequent sections develop our model, validate it empirically, and use it to draw substantive conclusions.

## 2. Searching the Web for Marketing Information

A simple example will help illustrate the research issues of interest. In October 1998, we queried each of six popular Web search engines to find documents containing the phrase "mere exposure effect." Alta Vista found 99 documents. Northern Light located 83; of course, many of these duplicated the ones from Alta

Vista. HotBot found fewer (49), but some had not been discovered by *either* Alta Vista or Northern Light. Finally, engines Infoseek, Excite, and Lycos found fewer documents (22, 21, and 9, respectively) but again some new pages were included. Together, all six engines located 172 documents, so even the "best" search engine (for this phrase) found less than 60% of this total (i.e., Alta Vista's 99 out of 172).

We should highlight that what we refer to simply as "search" (which is of course from the user's perspective) is really the result of a complex process. A search request does not directly cause a real-time search of the Web, but rather a (potentially complicated) look-up in a very large database. This database arises as the result of past webcrawling (i.e., proceeding from URL to URL and indexing the Web page contents) by the search engine and (less often) by specific requests from sites to be included in the engine's database. While any "search" request, then, produces only a search result from a database that is essentially static, a search request *can* affect the database for future searches, e.g., by causing certain URLs to be checked for viability or by influencing the future webcrawling pattern (by changing the engine's inferred popularity/importance for certain words or phrases). Our study simply examines the user's experience upon requesting URLs whose corresponding Web pages contain a particular phrase for these studied search engines.

We also acknowledge at the outset that this study will not attempt to assess the relative "value" of the individual sites found, and indeed one might well be skeptical of any mechanism that claimed to do so. Different searchers will no doubt have different interests or needs. Rather, thinking about this simple example leads directly to the five research issues that we do address:

1. *Search Engine Performance Across Phrases.* Would the search-result pattern above hold up for other marketing phrases? "Mere exposure effect" is relatively new to marketing and is associated more with academic research than with current marketing management practice. Perhaps some engines would do better for longer-established phrases, or those more prevalent among practitioners. Certainly, because Web crawlers proceed from document to document via the links provided, they may end up covering relatively separate,

disparate parts of the space of URLs. Such a propensity can be exacerbated by, for instance, the inclination of academic sites to link to other academic sites (via connection to coauthors, references, etc.).

2. *Factors Affecting Discovery of URLs.* In the example above, several URLs were found by all of the engines, while others were located by only one. For a given phrase, what makes some URLs "easy" to locate? In light of the Web crawler process mentioned, the more sites that link to a URL, the easier finding that URL will be. Of course, this measure is essentially impossible to observe. It is also not directly controllable by a site that *wishes* to be found. Instead, we focus on two factors that are observable and (within limits) controllable: the number of links *on* a URL (to other documents), and the domain type (.com, .edu, .org, etc.). The former should be related to URL discoverability because it is an indicator of sophistication and connectedness and may also stimulate reciprocal linkage (i.e., a linked site electing to provide a link back). The latter factor (domain type) may matter through a propensity for sites to link within (rather than across) these types.

3. *Search Engine Structural Characteristics.* Although search engines' operating details are proprietary, they are known to differ with respect to some basic characteristics. We will summarize the apparent relationship between such structural properties and the engines' search performance.

4. *Overlap and Sequential Search.* We are also interested in the way in which patterns of overlap among the search engines determine their incremental benefit when combined. In our example above imagine that Alta Vista was the search engine used first. Would using a second engine be expected to add substantially to the number of documents found? What about a third? How many engines are needed to find the "lion's share" of relevant documents? Which particular engine would add most to, say, Alta Vista's results? The proposed model will allow us to answer these questions.

5. *How Much Information Did We Miss?* Using all six search engines we found 172 documents mentioning "mere exposure effect." But how many documents did we fail to find? Note that any single URL's search results can be summarized by a binary six-vector, where

the  $i$ th element is a "1" if search engine  $i$  found the URL in question and a "0" if it did not. There are, of course,  $2^6 = 64$  such patterns, and for each phrase searched we can create the full frequency count among these 64 patterns—except for one. The number of URLs associated with the (0,0,0,0,0,0) vector is not available because this represents the number of URLs missed by all six search engines. However, after creating a model that represents well the engines' Web coverage and overlap (by fitting the 63 patterns above), we will forecast the frequency of this 64th pattern—as it indicates the size of the remaining "undiscovered" part of the Web.

To build a model that would address these five issues, we proceeded through four steps to build an appropriate database.

*Step 1: Marketing Phrases for Search.* The marketing phrases searched needed to be diverse enough to represent an interesting universe and also vary on the factors thought to affect search engine performance (i.e., research issue 1, above). Accordingly, phrases were selected via three criteria:

1. They are relatively central to marketing thought, appearing in popular reference works (Bennett et al. 1995, Clemente 1992).
2. They are specific enough so that a Web search need not be refined further (e.g., "marketing management" was found on 44,432 Web pages by Alta Vista—too many to be helpful without more detail).
3. They span the two phrase dimensions discussed earlier: managerial versus academic, and newer versus older. Five phrases were selected in each cell of the resulting  $2 \times 2$  design, leading to 20 phrases overall.

*Step 2: Phrase Search Via Search Engine.* The six search engines examined here (Alta Vista, HotBot, Excite, Infoseek, Northern Light, Lycos, located at <http://www.altavista.com>, <http://www.hotbot.com>, <http://www.excite.com>, <http://infoseek.go.com>, <http://www.northernlight.com>, and <http://www.lycos.com>, respectively) are the most popular based on user awareness, popular press mentions, and inclusion in previous studies and in metasearch programs (*PC Magazine Online* 1998, Beatty 1998, Feldman 1998, Lawrence and Giles 1998 and 1999). Note that while Yahoo! is often mentioned by users as a "search engine," it is actually a directory, and at the time of our

study Yahoo! additionally incorporated Inktomi, the same search engine used by HotBot. Thus, we did not include it. Although, as recently pointed out by Lawrence and Giles (1999), HotBot, Microsoft Snap, and Yahoo! do not return exactly the same information because of filtering and/or different underlying Inktomi databases. The 20 phrases were searched using each of the six engines during October 1998. During the search, two properties of each located URL were recorded: the number of links (0–5, 6–10, or 10+), and the domain type (.com, .edu, .org, or "other") indicating whether the site was commercial, academic, an organization, or other (the latter including non-U.S. sites). This URL information will allow us to address research issue 2 above.

*Step 3: Integrate Search Results.* As noted earlier, the search result for any located URL can be summarized in a binary six-vector. However, meaningfully comparing these results across engines requires substantial care. The same document may be reached by different alphanumeric strings, requiring that the documents themselves be accessed and checked both for similarity across engines and for duplication within an engine. URLs were also checked to verify that they were active and in fact contained the phrase in question. (Both Excite and Infoseek use heuristics that may return URLs similar—but not identical—to the phrase searched. These instances were deleted.)

*Step 4: Search Engine Characteristics.* As in research issue 3, we want to link search engines' performance to their characteristics. Because the number of search engines is small, it would not be useful to formally incorporate these characteristics into the model itself, but we will be able to investigate an association between overall search performance and an engine's structural properties. The key properties of interest are engine size (the total number of pages indexed) and several binary indicators of capability. The latter includes Depth (whether an engine searches an entire site without a preset limit), Frame Support (ability to follow frame links), Image Maps (ability to follow image maps), and Learns Frequency (whether an engine estimates the frequency with which a page's content changes, and uses that information to determine visit frequency). Other search engine characteristics would be interesting to include (such as number of pages

crawled per day) but do not appear to be reliably measured and available (Sullivan 1998). The search engine features above were taken from the Search Engine Watch site (Sullivan 1998) and were measured as of August 4, 1998.

Table 1 shows the 20 marketing phrases, their categorization regarding newness and academic/managerial, and the total number of URLs found by each engine for each phrase. Note that this table is not the complete data, but rather is a summary. For each of the 1588 located URLs, the data used in our model-development are a binary six-vector together with the two URL characteristics (number of links, domain type) and two phrase characteristics (as above).

As a further summary, Table 2 shows how the URLs found are distributed across phrase and URL characteristics. The table entries provide for a given engine, the proportion of all URLs found (by any engine) hav-

ing a particular characteristic. For instance, Alta Vista located 52.1% of all managerial-phrase URLs that were found. It did a little better (53.5%) finding academic-phrase URLs. Relative to the engine's baseline level of performance across all phrases, Infoseek had the greatest skew toward locating academic-phrase URLs (0.163 academic versus 0.125 managerial), and Northern Light had the greatest inclination toward managerial-phrase URLs (0.462 academic versus 0.529 managerial). Overall, Alta Vista had the best performance in finding academic-phrase URLs, while Northern Light had the greatest success finding marketing-managerial ones. Analogous conclusions for other phrase/URL characteristics are available via Table 2. Table 3 provides the structural characteristics of the engines.

Before developing our model, it was useful to note what would happen if search outcomes for any given phrase were independent—i.e. if each URL had some

**Table 1.** Number of URLs Found By Search Engine and Marketing Phrase

Phrase	Manag.	Newer	Search Engine*						Total: 6 Engines
			AV	HB	EX	IS	NL	LY	
flanker brand	1	1	9	9	1	6	5	0	19
umbrella branding	1	1	38	21	7	4	51	0	76
second mover advantage	1	1	8	9	4	2	20	1	26
professional respondents	1	1	41	19	12	7	31	0	62
audience fragmentation	1	1	106	59	37	36	120	14	215
category development index	1	0	18	11	0	2	19	2	29
modified rebuy	1	0	40	23	5	3	33	1	78
perceived value pricing	1	0	19	14	4	6	21	5	38
simulated test market	1	0	25	15	8	15	35	7	66
unaided recall	1	0	92	45	29	14	67	13	150
low involvement learning	0	1	10	11	5	7	13	4	22
elimination by aspects	0	1	61	35	21	8	57	4	114
mere exposure effect	0	1	99	49	21	22	83	9	172
preference map	0	1	29	21	10	35	41	1	101
decision calculus	0	1	74	54	28	32	55	14	134
multiattribute attitude models	0	0	17	6	2	1	26	2	37
Reilly's law	0	0	27	13	6	3	20	0	40
wheel of retailing	0	0	68	28	12	10	44	2	113
beta binomial model	0	0	39	13	9	10	33	4	64
diffusion of innovation model	0	0	20	13	6	7	11	2	32
Total			840	468	227	230	785	85	1588

\*AV = Alta Vista, HB = HotBot, EX = Excite, IS = Infoseek, NL = Northern Light, LY = Lycos

Note: Manag. = 1 indicates a managerial phrase, 0 an Academic phrase. Newer = 1 a newer phrase, 0 an older phrase.

**Table 2** Search Engine Results by Phrase Age, Phrase Type, URL Number of Links, and Domain Extension

Engine	Age		Type		Links			Domain			
	New	Old	Manag.	Acad.	0–5	6–10	10 +	edu	com	org	other
AV	0.504	0.564	0.521	0.535	0.523	0.545	0.548	0.495	0.557	0.644	0.530
HB	0.304	0.280	0.297	0.293	0.284	0.288	0.328	0.312	0.269	0.328	0.288
Ex	0.155	0.125	0.140	0.144	0.146	0.138	0.137	0.140	0.142	0.164	0.143
IS	0.169	0.109	0.125	0.163	0.135	0.155	0.167	0.153	0.110	0.205	0.147
NL	0.506	0.478	0.529	0.462	0.481	0.551	0.505	0.502	0.526	0.521	0.458
LY	0.050	0.058	0.056	0.050	0.044	0.080	0.066	0.056	0.088	0.041	0.026

**Table 3** Structural Characteristics of Search Engines\*

Characteristics	Search Engine					
	AV	HB	EX	IS	NL	LY
Size (million pages)	140	110	55	30	80	30
Depth of Search	No Limit	No Limit	No Limit	Sample	No Limit	Sample
Frames Support	Yes	No	No	No	Yes	No
Image Maps	Yes	No	No	Yes	Yes	No
Learns Frequency	Yes	Yes	No	Yes	No	No

\*Source Search Engine Watch (Sullivan 1998)

probability of being located (possibly engine-specific) and one engine's finding the URL told us nothing about any other engine's. In such a situation, substantive research issues 1 and 2 (effect of URL and phrase characteristics) could be addressed by a separate simple model (e.g., logistic regression) for each search engine, and research issue 4 (overlap between engines) would have a very simple answer for any set of engines. The independence assumption is also the linchpin of the most careful model published so far for search engine performance (Lawrence and Giles 1998). They consider a model with the top two engines assumed to be independent. Accordingly, we begin by considering the independence assumption in detail.

### 3. Are Search Engine Outcomes Independent?

The simplest, and arguably the most natural, starting point for representing the URLs found by multiple

Web search engines is the independent binomial model. It is based on two assumptions. First, for any given search phrase  $j$ , it imagines that any given search engine  $i$  finds any one of the URLs containing that phrase independently of its finding other such URLs, and with some probability  $p_{ij}$ . Second, the model assumes that the probability  $p_{ij}$  that search engine  $i$  finds any particular URL containing phrase  $j$  does not depend on the set of URLs found by any *other* search engine.

For a single URL containing phrase  $j$ , the data can be written simply as the binary six-vector  $(y_{1jk}, y_{2jk}, y_{3jk}, y_{4jk}, y_{5jk}, y_{6jk})$ , where  $y_{ijk} = 1$  if the  $k$ th URL for phrase  $j$  is found by search engine  $i$ , and is 0 otherwise. For URL  $k$  and phrase  $j$  the likelihood function is

$$L(y_{1jk}, y_{2jk}, y_{3jk}, y_{4jk}, y_{5jk}, y_{6jk}) = \prod_{i=1}^6 p_{ij}^{y_{ijk}} (1 - p_{ij})^{1-y_{ijk}}, \quad (1)$$

where  $p_{ij}$  is the probability that engine  $i$  finds any given URL containing phrase  $j$ . Because the URLs are exchangeable by assumption, the likelihood for the data for phrase  $j$  is the product of (1) across all URLs (in practice, a partial likelihood will be used, since the (0,0,0,0,0,0) vector will be missing).

This independent binomial model has much to recommend it. It is parsimonious: Each search engine  $i$  (for each phrase  $j$ ) can be summarized by a single quantity—its search success probability  $p_{ij}$ . The model can provide an estimate of the number of URLs not found. After any number of search engines have been used, the expected number of *new* URLs from another

search engine  $h$  is simply  $(N_j - m)p_{hj}$ , where  $m$  is the cumulative number of URLs already found and  $N_j$  is the (unknown) number of URLs containing phrase  $j$ .

Lawrence and Giles (1998) expressed concern about the independence assumption, and that concern was well founded. We report in Table 4 the value of  $-\log L$  for this model and the associated BIC statistic. Four particular versions of the independent binomial model were evaluated: (1) constant  $p$ , (2) different  $p$  for each engine but constant across phrases, (3) different  $p$  for each phrase but constant across engines, and (4) different  $p$  for each engine and phrase. A simple chi-square test on the value of  $-2\log L$  rejects each of these four models. Naturally, with over 1,500 observations the power of such a test is very high and may not in itself present a strong case for substantial interdependence. Instead, two other considerations will argue for a model that relaxes the independence assumption. First, we will see later that relevant goodness-of-fit indicators can be improved substantially via a spatial interdependence model. Second, we note that the BIC criterion (which penalizes highly parameterized models for data overfitting) actually prefers, among independence models, the one where location probabilities differ only by search engine (and not by phrase). In other words, search is characterized simply by six  $p_i$  values, one for each search engine (the relative magnitude of the  $p_i$  are given by the total URL count by engine in Table 1).

It is easy to show that an estimate of the number of URLs found by all engines in any three-engine set (denoted 1,2,  $t$  for convenience) under this model is:

$$n_{12t} = \frac{n_{12}^2}{n_1 n_2} n_t. \quad (2)$$

Taking, for instance, Alta Vista and HotBot as engines "1" and "2," the actual three-way overlap  $n_{12t}$  and the overlap predicted by the independence model via (2), are:

Set of Search Engines	Actual Number of URLs	Predicted Number of URLs
Alta Vista, HotBot, Excite	50	22.4
Alta Vista, HotBot, Infoseek	37	22.7
Alta Vista, HotBot, Northern Light	100	77.5
Alta Vista, HotBot, Lycos	19	8.3

In short, looking across our 20 marketing phrases, the independence model substantially underpredicts the actual overlap for these triplets of search engines. These positive residuals suggest that two search engines with high coverage (Alta Vista and HotBot) are inclined to subsume the other four engines. This suggests the use of Rasch-type ability/difficulty models (Rasch 1966, Andersen 1973), whereby the probability that a given URL is located is a function of both a URL "difficulty" parameter and an search engine "ability" parameter. In this kind of model the "easy" URLs will tend to be found by all search engines and the "hard" URLs only by the search engines that find many overall. In other words, Alta Vista and HotBot will overlap somewhat, but the other search engines will overlap even moreso with this pair (and hence produce positive residuals) because the URLs they find will tend to be the "easy" ones. Of course, other search engine triplets could show different discrepancies than those observed above. Our main point is the observation that independence does not appear to be a solidly supported assumption, and a model where spatial location of search engines determines patterns of overlap may have value.

**Table 4** Global Goodness-of-Fit for Independence Models

Model	# Parameters	$-2^*LL$	BIC
Constant $p$	1	11236.72	11245.88
Different $p$ by engine	6	9602.83	9657.80
Different $p$ by phrase	20	11197.58	11380.82
Different $p$ engine by phrase	120	9276.17	10375.60

*Note:* Reported are  $-2^*$  Log-Likelihood, and the BIC criterion.

## 4. A General Proximity Model

We provide initially a heuristic description of our modeling approach for WWW data. This non-formal description is useful to describe our intuition, why we expect this class of models to improve on simpler ones, and the expected limitations and subsequent improvement in fit as our models become more complex. Needed notation and formal models are presented after.



#### 4.1. Heuristic and Graphical Descriptions

We posit a general class of models for the ability of WWW search engines based on the proximity (“distance”) from a specific engine to a given URL and the “reach” of an engine. Our basic model suggests that when an engine and URL are proximate, the engine is likely to find that URL, and unlikely when not. In particular, each engine and URL are hypothesized to “sit” at an unknown location in  $D$ -dimensional space. A URL’s location is modeled to be centered around a mean location determined by both its phrase and covariates specific to the phrase and URL (e.g., type of phrase, URL domain extension, etc.). Then, from an engine’s location, it “throws out a net” and probabilistically captures URLs within its reach. That is, there is a monotonically decreasing relationship between distance from engine to URL and the probability a URL is found. Pushing this analogy farther, inferences of interest under the model are then derived from: (1) the location of each engine (that is, do “weaker” engines find just a subset of those URLs found by the better engines, which would follow if all engines were located at the same place, or do engines “carve” out their own locations), (2) the size of the net for each engine (in our model this is the ability of the engine), (3) the shape of the net (are the underlying dimensions related), (4) the number of underlying dimensions  $D$  adequate to model the data, (5) the effects, if any, of phrase and URL covariates on URL’s locations and hence their probability of being found, and (6) an exponent determining how fast the probability of an engine finding a URL drops off as a function of their proximity. We considered three specific cases of this general proximity model.

As a point of reference for describing the proximity models, consider the graphical representation of the independent binomial model (§ 3) shown in Figure 1, panel A. The horizontal line represents the ( $D = 1$  dimensional) space of URL locations, and the various search engines differ in the degree to which they (probabilistically) cover this space, beginning at the origin. The graph can be interpreted as having each engine stand at the origin and throw out a line, capturing as many URLs (“fish”) as possible. Because engines with longer fishing lines (i.e., more ability) reach out farther

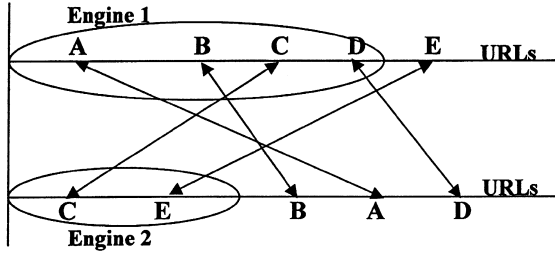
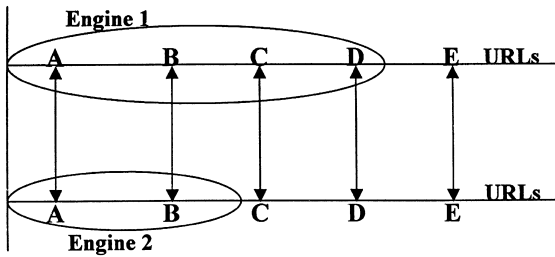
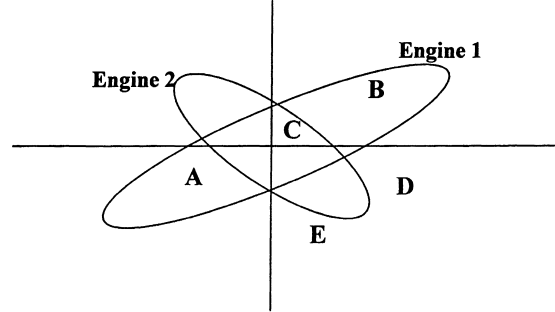
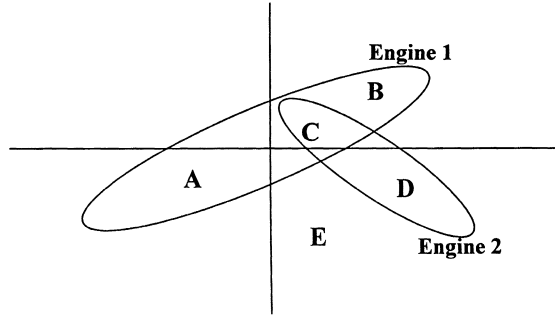
from the origin, they are likely to “catch” more URLs, although *which* URLs the better engine (engine 1) finds is unrelated to the *specific* URLs found by the weaker engine (engine 2). That is, via the independence assumption, it is as if the URLs randomly redistributed their locations in the time elapsed between the search by the two engines.

As an alternative to this independence model, we will examine a  $D = 1$  dimensional proximity model, depicted in Figure 1 panel B and denoted “Model 1” below. Here, each engine is again located at the origin and casts its probabilistic coverage of the line according to its “ability.” But unlike the independence model, here the URL locations remain fixed. Accordingly, some URLs really are more difficult to locate (i.e., those labeled “D” and “E”) than others (e.g., “A” and “B”) as they lie far from the origin. As a result, it is unlikely that the search engines with lesser ability will find URLs not found by the better engines. As suggested in § 2 (and confirmed in § 5.1) this feature of Model 1 does not fit the data particularly well. (Even the weakest engine Lycos finds URLs not found by other engines). This suggested the extension of Model 1 in two ways under our general proximity model structure. First, in Model 2 (Figure 1 panel (C)) we extend to  $D = 2$  dimensions yet leave all of the engine locations at the origin. A more general version considered in Model 3 (Figure 1 panel (D)) also allows the engine locations to vary—i.e., as suggested earlier, a search engine may “stake out” a distinctive part of the URL space. As shown below, the results indicate that Model 3 is necessary to provide an adequate fit to the pattern of Web search results for marketing information.

#### 4.2. Model Notation, Development, and Computational Approach

We consider the case described in § 2 where each of  $i = 1, \dots, I$  search engines is utilized on the WWW to locate URLs for each of  $j = 1, \dots, J$  phrases. Let  $K_j$  denote the total number of distinct URLs found for the  $j$ th phrase (by any of the engines) and  $y_{ijk}$  a binary outcome where  $y_{ijk} = 1$ ,  $k = 1, \dots, K_j$ , if the  $k$ th URL for the  $j$ th phrase is found by engine  $i$ , and 0 otherwise. The collection of all outcomes  $y_{ijk}$  is denoted  $Y$ . In addition, for each URL we obtain covariate vector  $x_{jk} =$

Figure 1 Models for URL Discovery by Web Search Engines

**A. Independent Multinomial Model****B. One-Dimensional Ability/Difficulty Model ("Model 1")****C. Two-Dimensional Ability/Difficulty Model ("Model 2")****D. Two-Dimensional Spatial Coverage Model ("Model 3")**

$(x_{jk1}, \dots, x_{jkP})$  to identify known characteristics of phrases and/or URLs that may make them harder or easier to find. The collection of all covariates is denoted  $X$ .

We posit a proximity model for  $p_{ijk} = \text{Prob}(y_{ijk} = 1)$  defined as a function of the following engine and URL specific parameters. Let  $\theta_i^t = (\theta_{i1}, \dots, \theta_{iD})$  and  $\gamma_{jk}^t = (\gamma_{jk1}, \dots, \gamma_{jkD})$  denote the location of the  $i$ th engine and  $k$ th URL for phrase  $j$  in  $D$ -dimensional space. Additionally, define  $\Sigma_i$ , a  $D \times D$  dimensional scaling matrix for engine  $i$ , and  $d_{ijk} = d(\theta_i, \gamma_{jk}) = (\theta_i - \gamma_{jk})^t \Sigma_i^{-1} (\theta_i - \gamma_{jk})$  a squared Mahalanobis distance between engine  $i$  and the  $k$ -th URL for phrase  $j$ . Thus, the diagonal elements of  $\Sigma_i$  are the abilities ("reach") and the off-diagonal elements indicate the covariation of abilities for engine  $i$  in the  $D$  dimensions.

We assert a model for  $p_{ijk}$  as a function of  $d_{ijk}$  given by

$$p_{ijk} = \frac{1}{1 + d_{ijk}^u}, \quad (3)$$

where  $u$  defines the rate at which the probability an engine finds a given URL drops off. In general, spatial/distance models have been utilized in other marketing contexts, especially brand choice (Elrod 1988, Kamakura and Srivastava 1984). We note that (3) is equivalent to  $\text{logit}(p_{ijk}) = -u \cdot \log(d_{ijk})$ , a logistic link where  $u$  is the slope of regressor  $\log(d_{ijk})$ . Assuming conditional independence of engines, phrases, and URLs within phrase this yields a product Bernoulli likelihood for parameters  $\Omega_1 = (\theta_1, \dots, \theta_I, \gamma_{11}, \dots, \gamma_{JK_I}, \Sigma_1, \dots, \Sigma_I, u)$  equal to

$$p(Y|\Omega_1) = \prod_i \prod_j \prod_k \left( \frac{1}{1 + d_{ijk}^u} \right)^{y_{ijk}} \left( \frac{d_{ijk}^u}{1 + d_{ijk}^u} \right)^{1 - y_{ijk}}. \quad (4)$$

Because commonalities are likely to exist among the engines, the phrases, and the URLs, we extend the

model for  $Y$  given in (4) to include a set of prior distributions for  $\Omega_1$ , allowing for the sharing of information across units. The choice of priors for the components of  $\Omega_1$  were made in the following manner. Because the six engines that we consider represent the engines of interest, we treat the engine specific parameters as fixed effects and put non-informative priors on  $\theta_i$ ,  $\Sigma_i$ ,  $i = 1, \dots, I$ . A non-informative prior is also adopted for  $u$  reflecting our lack of knowledge regarding this parameter. In contrast, it is of interest to summarize the location of phrase  $j$  for which we may regard  $\gamma_{jk}$ ,  $k = 1, \dots, K_j$  as a random sample of URLs drawn from a population distribution. By convention and for computational convenience, we put a hierarchical multivariate normal-Inverse Wishart prior structure on the URL locations:

$$\begin{aligned}\gamma_{jk} &\sim \text{MVN}_D(\alpha_j + \beta x_{jk}, A_j) \\ \alpha_j &\sim \text{MVN}_D(\bar{\alpha}, \Sigma_\alpha) \\ A_j &\sim W_\nu^{-1}(S),\end{aligned}\quad (5)$$

where  $\text{MVN}_D(x, y)$  denotes a  $D$ -dimensional multivariate normal distribution with mean vector  $x$  and covariance matrix  $y$ ,  $\alpha_j = (\alpha_{j1}, \dots, \alpha_{jD})$  the mean location of phrase  $j$ ,  $\beta$  a  $D \times P$  dimensional coefficient matrix where  $\beta_{dp}$  is the slope for the  $p$ th covariate in dimension  $d$ ,  $\bar{\alpha} = (\bar{\alpha}_1, \dots, \bar{\alpha}_D)$  the population mean of the phrase locations,  $A_j$  and  $\Sigma_\alpha$  are  $D \times D$ -dimensional covariance matrices for phrase  $j$  and the population of phrase means, and  $W_\kappa^{-1}(Q)$  denotes an Inverse-Wishart distribution with  $\kappa$  degrees of freedom and scale matrix  $Q$ . The values of  $\nu$  and  $S$  were chosen as uninformative, allowing the data to fully specify the values of  $A_j$ . As well, a noninformative prior distribution was utilized for  $\beta$ . We denote the prior level parameters by  $\Omega_2 = (\alpha_1, \dots, \alpha_I, \beta, A_1, \dots, A_I, \bar{\alpha}, \Sigma_\alpha)$  and the prior distribution by  $p(\Omega_1 | \Omega_2)$ .

Inferences for the model parameters  $\Omega_1$  and  $\Omega_2$  were derived by obtaining samples from the marginal posterior distributions  $p(\Omega_1 | Y, X)$  and  $p(\Omega_2 | Y, X)$  using a Markov chain Monte Carlo (MCMC) sampler (Gelfand et al. 1990, Rossi et al. 1996). For each of Model 1, Model 2, and Model 3, we report results obtained by running three independent chains for 3000 draws from overdispersed starting positions, discarding the initial 500 draws of each chain after determining convergence

(German and Rubin 1992) and estimating the quantities of interest using the remaining 7500 draws. Further details are provided in the appendix.

## 5. Results

### 5.1. Model 1: One-Dimensional Ability/Difficulty Model

We first considered a simple special case of the general proximity model defined by (3), (4), and (5), which consisted of a  $D = 1$  dimensional model with all engines located at the origin  $\theta_1 = \dots = \theta_I = 0$ . To identify the model, we set as a reference point  $\Sigma_1 = 1$ , the ability of Alta Vista indexed as  $i = 1$ , and set the rate factor  $u = 0.5$ . This model, in which each engine ("examinee") has a unidimensional ability  $\Sigma_i$  and each URL has a unidimensional location  $\gamma_{jk}$  ("test item difficulty"), is similar in spirit to the Rasch (1960) model commonly used in educational testing.

Model 1 was applied to the set of 20 phrases and 1588 URLs described in § 2. A summary of results for engine abilities, presented as  $\Sigma_i$  is given in column 2 of Table 5. The ordering of engine abilities suggested (Alta Vista, Northern Light, HotBot, Excite = Infoseek, Lycos) is unambiguous in all comparisons (true for all 7500 draws) except for the comparison (a) Alta Vista  $>$  Northern Light,  $p = 0.78$  and (b) Excite  $\geq$  Infoseek,  $p = 0.48$ . The results from Models 2 and 3, better fitting models described later, will further refine these relations.

Inferences under Model 1 for the phrase and URL covariates (1) domain extension: .edu, .com, .org, other, (2) # of Links on the URL: 0–5, 6–10, 10+, (3) Type of Phrase: Managerial/Academic, (4) Age of Phrase: Newer/Older, and (5) the interaction between (3) and (4), on the mean of URL locations  $\gamma_{jk}$  and hence  $p_{ijk}$ , are given in column 2 of Tables 6 and 7. In Table 6 we report the posterior median, standard error, and probability of the effect being greater than 0 for each covariate. Table 7 gives the adjusted phrase mean for URLs with a given covariate level. To interpret these findings, recall that all engines for Model 1 are located at the origin, thus any positive coefficient suggests that the covariate level makes URLs of that type harder to find, and vice-versa. We observe strong evidence that

**Table 5** Posterior Medium Engine Abilities on Dimensions 1 and 2 ( $\Sigma_{11}$ ,  $\Sigma_{22}$ ), Correlation ( $\rho_{12}$ ), and Distance Factor  $u$  for Models 1, 2, and 3

Engine	Model 1	Model 2			Model 3		
	Dim 1	Dim 1	Dim 2	$\rho$	Dim 1	Dim 2	$\rho$
AV	1.000 (—)	1.000 (—)	1.760 (0.41)	0.005 (0.01)	1.000 (—)	1.960 (0.40)	0.020 (0.01)
HB	0.128 (0.03)	0.157 (0.01)	0.199 (0.00)	0.006 (0.01)	0.074 (0.01)	0.761 (0.01)	0.020 (0.00)
Ex	0.017 (0.00)	0.040 (0.01)	0.024 (0.01)	−0.586 (0.01)	0.020 (0.01)	0.060 (0.01)	−0.630 (0.01)
IS	0.017 (0.00)	0.010 (0.01)	0.040 (0.01)	0.593 (0.01)	0.055 (0.01)	0.052 (0.01)	0.640 (0.01)
NL	0.945 (0.08)	2.670 (0.48)	1.020 (0.40)	−0.007 (0.09)	3.720 (0.52)	1.870 (0.45)	−0.003 (0.00)
LY	0.001 (0.00)	.0004 (0.00)	0.001 (0.00)	0.000 (0.00)	0.0003 (0.00)	.0008 (0.00)	0.000 (0.00)
$u$	0.500 (—)		0.369 (0.02)			0.354 (0.015)	

Note: Posterior standard errors are in parenthesis.

**Table 6** Phrase and URL Covariate Slopes ( $\beta$ ) for Models 1, 2, and 3.

Cov.	Model 1		Model 2			Model 3				
	Dim 1		Dim 2		Dim 1		Dim 2		Dim 2	
.edu	−0.052 (0.06)	0.177	−0.104 (0.05)	0.050	−0.145 (0.07)	0.008	−0.208 (0.06)	0.000	0.017 (0.04)	0.653
.com	0.018 (0.07)	0.610	−0.223 (0.09)	0.003	0.008 (0.08)	0.423	−0.213 (0.09)	0.003	−0.182 (0.04)	0.000
.org	−0.090 (0.12)	0.117	−0.134 (0.11)	0.005	0.020 (0.10)	0.633	−0.222 (0.10)	0.018	0.057 (0.05)	0.998
OL-5L	0.099 (0.05)	0.957	−0.280 (0.06)	0.005	−0.030 (0.08)	0.470	0.138 (0.06)	0.990	−0.146 (0.03)	0.000
6L-10L	0.136 (0.09)	0.947	0.090 (0.10)	0.733	−0.017 (0.09)	0.360	0.042 (0.09)	0.673	−0.132 (0.05)	0.005
Man.	0.067 (0.11)	0.733	0.198 (0.11)	0.990	0.370 (0.12)	1.000	−0.114 (0.10)	0.148	0.148 (0.10)	0.913
Newer	0.106 (0.11)	0.807	0.137 (0.08)	0.930	0.082 (0.15)	0.635	−0.230 (0.08)	0.000	0.280 (0.04)	1.000
Int.	−0.112 (0.15)	0.237	−0.457 (0.15)	0.000	−0.388 (0.22)	0.010	0.026 (0.13)	0.560	−0.374 (0.10)	0.000

Note: Reported are the posterior medians (standard deviations) and posterior probability of the effect being greater than 0. Int. is the interaction between managerial and newer.

URLs with fewer links are harder to find than those with the most number of links (10+) and modest evidence that URLs having domain extensions .edu or .org are slightly easier to find. These results are also confirmed by Lawrence and Giles (1999). Other inferences were: (1) there was no significant difference in the phrase locations (posterior median of  $\Sigma_\alpha = 0.001$ ), which is consistent with the stable hit rates for each engine by phrase reported in Table 1 and the log-likelihoods reported in Table 4, and (2) URL variances  $\Sigma_j$  were inversely related to the number of URLs found ( $r = -0.85$ ).

A more detailed and informative look at the performance of Model 1 is presented in columns 3 through

6 of Table 8. Here we consider the number of URLs, showing each of the  $2^6 = 64$  possible engine-hit patterns. The table provides the observed number  $n_{obs}$  for each pattern (excluding (0,0,0,0,0,0)), as well as the 2.5%, 50%, and 97.5% percentiles for the predicted frequency. Some interesting residuals are evident. First, we note that Model 1 tends to underpredict the number of unique URLs found by each engine as seen in the unique engine-hit patterns 32, 48, 56, 60, and 62 (pattern 63 is slightly overpredicted). Second, and related to the underprediction in the number of uniques, Model 1 also tends to overpredict the number of URLs found by exactly two engines, as seen in patterns 16, 24, 28, 30, 44, 46, 47, 52, 54, 59, and 61 (patterns 31, 55,

**Table 7** Effect of Phrase and URL Covariates  $x_{jk}$  on the Mean Phrase Location

$x_{jk}$	Model 1	Model 2	Model 3
	$\mu_1$	$(\mu_1, \mu_2)$	$(\mu_1, \mu_2)$
.edu	0.968	(0.026, -0.225)	(0.017, 0.085)
.com	1.038	(-0.093, -0.072)	(0.012, -0.114)
.org	1.011	(0.004, -0.060)	(0.003, 0.125)
OL-5L	1.119	(-0.105, -0.110)	(0.363, -0.078)
6L-10L	1.156	(0.220, -0.097)	(0.268, -0.064)
Man.	1.087	(0.328, 0.290)	(0.111, 0.216)
Newer	1.126	(0.267, 0.000)	(-0.005, 0.348)
New + Man.	1.081	(-0.008, -0.016)	(-0.113, 0.112)

Note:  $(\bar{\alpha}_1, \bar{\alpha}_2)$  is the mean phrase location with all covariates at baseline levels.  $(\mu_1, \mu_2) = (\alpha_1 + \beta_1 x_{jk}, \alpha_2 + \beta_2 x_{jk})$  are the new coordinates including the covariate effects. Model 1:  $\bar{\alpha} = 1.020$ , Model 2:  $(\bar{\alpha}_1, \bar{\alpha}_2) = (0.130, -0.080)$ , Model 3:  $(\bar{\alpha}_1, \bar{\alpha}_2) = (0.225, 0.068)$

and 58 are adequately fit, and pattern 40 is underpredicted). These results were not surprising because in Model 1 each engine is located at the origin and is casting its "fishing line" in the same direction.

One further inference that can be derived from the model is an estimate of the number of URLs not found by any of the engines. This question has managerial relevance from two perspectives: (1) A manager searching for URLs on a specific topic may wish to know the fraction of those related URLs he or she is likely to find by using these six engines; and (2) consider the owner of a URL wanting his or her Webpage to be found. Under the model, we can compute the posterior distribution of the number of URLs not found,  $K$ , by noting that

P(all engines miss a URL)

$$= \prod_i (1 - p_{ijk}) \Rightarrow$$

P(at least one finds it)

$$= 1 - \prod_i (1 - p_{ijk}) \Rightarrow$$

$$n_{obs} = \left(1 - \prod_i (1 - p_{ijk})\right) * K \Rightarrow$$

$$K = \frac{n_{obs}}{1 - \prod_i (1 - p_{ijk})}. \quad (6)$$

These results are shown in pattern 64 and suggest that the 95% posterior interval for the number of missing URLs for the 20 phrases is (253.94, 330.30) with posterior median 283.43. This indicates that Model 1 predicts  $283.43 / (1588 + 283.43) \approx 15\%$  of the URLs are missed by using all six engines.

## 5.2. Model 2 and Model 3 Results

We considered two additional special cases of the general proximity model to improve on Model 1. Model 2 consisted of a  $D = 2$  dimensional version where each engine was located at the origin ( $\theta_{11} = \theta_{12} = \dots \theta_{1n} = \theta_{12} = 0$ ). As a scale identifiability constraint we set  $\Sigma_{11}$ , the ability of Alta Vista on dimension 1, equal to 1. By definition, the addition of a second dimension would improve the fit; however, we suspected that locating each engine at the origin, as per a pure ability/difficulty model, would still provide an inadequate fit. In Model 3, we generalize Model 2 to allow individual search engines to carve out a distinctive portion of (two-dimensional) URL space, i.e., the engine locations ( $\theta_{i1}, \theta_{i2}$ ) were allowed to vary. In fitting Model 3, we set  $\theta_{11} = \theta_{12} = 0$ ,  $\theta_{21} = 0$ , restricted  $\theta_{31} > 0$ , and put  $\Sigma_{11} = 1$  as shift, y-axis rotation, x-axis rotation, and scale identifiability constraints respectively.

Models 2 and 3 results for engine abilities  $\Sigma_{i11}$ ,  $\Sigma_{i22}$ , and the correlation between dimensions  $\rho_{i12}$  is given in Table 5 (columns 3–8). A graphical representation of the engine performances for Model 3 is given in Figure 2, panels A and B. The results suggest that there are indeed two unique dimensions in which engines operate. Model 2 findings give the ordering in dimension 1 of Northern Light, Alta Vista, HotBot, Excite, Infoseek, and Lycos, whereas dimension 2 results give the ordering Alta Vista, Northern Light, HotBot, Infoseek, Excite, and Lycos. This is consistent with Model 1 findings of an ambiguous ordering of Alta Vista versus Northern Light and Excite versus Infoseek. However, we note that the total "area" covered by Northern Light is superior to that of Alta Vista because its posterior median abilities (2.670, 1.020) suggest greater coverage than Alta Vista's (1.000, 1.760). These findings are replicated in Model 3, in which Northern Light is far superior to Alta Vista on dimension 1 (3.720 versus 1.000) and almost equal on dimension 2 (1.870 versus 1.960). This is suggested by Northern Light's high

**Table 8. Table of Web Engine Patterns and 95% Confidence Intervals for Models 1 through 3**

Number	Pattern	$n_{obs}$	Model 1			Model 2			Model 3		
			2.5%	50%	97.5%	2.5%	50%	97.5%	2.5%	50%	97.5%
1	111111	2	0.020	0.032	0.052	0.019	0.113	0.825	0.017	0.104	0.623
2	111110	3	0.573	0.901	1.262	0.453	2.009	9.105	0.385	1.796	11.195
3	111101	1	0.024	0.041	0.060	0.016	0.096	1.009	0.021	0.115	0.843
4	111100	2	0.809	1.120	1.415	0.372	2.006	11.473	0.467	2.067	9.218
5	111011	3	0.175	0.284	0.450	0.153	0.805	3.599	0.155	0.616	3.553
6	111010	21	5.641	7.762	10.218	3.647	13.798	38.908	3.794	10.159	46.521
7	111001	0	0.225	0.353	0.509	0.148	0.740	4.358	0.162	0.706	4.677
8	111000	18	7.946	9.418	11.372	3.569	13.459	55.313	4.115	10.819	43.796
9	110111	3	0.179	0.289	0.425	0.190	0.625	6.061	0.213	0.788	3.075
10	110110	16	5.911	7.831	9.669	5.174	11.833	44.923	4.465	13.090	38.602
11	110101	1	0.246	0.357	0.509	0.121	0.619	4.927	0.263	0.794	3.828
12	110100	9	8.188	9.518	11.848	3.874	11.461	45.196	4.783	13.902	42.832
13	110011	7	1.721	2.556	3.427	1.243	4.902	20.107	1.478	4.284	16.198
14	110010	45	56.218	68.057	76.816	39.428	89.492	227.409	33.747	80.826	149.755
15	110001	2	2.220	3.086	4.035	1.198	4.844	18.233	1.495	5.245	19.813
16	110000	64	77.743	82.908	94.274	36.453	83.855	203.724	38.321	89.130	174.093
17	101111	3	0.071	0.111	0.176	0.045	0.220	1.527	0.053	0.252	1.727
18	101110	5	2.090	2.957	3.983	0.807	4.139	19.556	0.944	4.334	25.004
19	101101	0	0.089	0.139	0.212	0.040	0.195	1.360	0.057	0.264	1.962
20	101100	4	2.947	3.648	4.577	0.751	3.936	20.847	1.097	4.850	20.581
21	101011	4	0.631	0.973	1.431	0.253	1.671	7.866	0.400	1.569	8.107
22	101010	25	20.300	25.408	32.607	6.517	27.677	91.151	9.312	26.298	123.853
23	101001	1	0.825	1.210	1.782	0.371	1.481	7.023	0.423	1.669	9.798
24	101000	25	28.132	31.949	36.674	5.151	27.985	109.350	11.617	28.970	106.308
25	100111	2	0.681	0.971	1.415	0.412	1.388	11.372	0.539	1.893	8.642
26	100110	20	21.908	26.016	31.929	8.630	24.931	92.002	13.482	32.703	98.609
27	100101	4	0.883	1.176	1.710	0.343	1.290	8.354	0.545	2.091	9.668
28	100100	19	28.950	32.011	37.390	6.627	24.896	115.699	12.236	35.924	129.168
29	100011	9	6.364	8.431	11.262	3.053	10.262	41.608	3.704	12.416	38.137
30	100010	174	206.952	225.715	244.697	89.316	194.995	394.634	102.703	209.021	379.458
31	100001	8	7.843	10.398	14.566	2.678	9.751	30.043	4.200	13.789	42.763
32	100000	340	259.575	279.628	303.729	47.639	186.828	347.876	197.822	290.697	378.992
33	011111	2	0.024	0.040	0.059	0.019	0.115	0.716	0.017	0.091	0.435
34	011110	1	0.747	1.097	1.440	0.688	2.015	7.967	0.386	1.451	8.018
35	011101	0	0.031	0.050	0.075	0.014	0.107	0.703	0.019	0.096	0.587
36	011100	3	1.025	1.331	1.661	0.479	1.970	11.391	0.414	1.675	7.726
37	011011	0	0.217	0.348	0.512	0.146	0.822	3.502	0.101	0.566	3.300
38	011010	9	7.311	9.215	11.637	4.358	15.006	51.108	3.175	9.224	33.472
39	011001	0	0.289	0.424	0.631	0.176	0.718	2.900	0.098	0.600	4.482
40	011000	24	9.909	11.462	13.131	3.415	13.234	60.630	2.877	10.249	36.014
41	010111	3	0.235	0.351	0.495	0.190	0.752	4.655	0.229	0.723	2.877
42	010110	11	7.567	9.521	11.487	4.499	12.636	45.393	3.630	11.113	40.041
43	010101	2	0.303	0.432	0.624	0.147	0.666	3.678	0.140	0.799	3.564
44	010100	8	9.834	11.614	14.324	3.495	11.751	58.152	2.893	13.000	40.329
45	010011	3	2.199	3.041	3.941	1.410	5.082	22.546	1.281	4.111	14.191
46	010010	54	73.643	81.158	91.576	31.988	99.461	192.606	26.370	76.380	144.021
47	010001	2	2.703	3.684	5.172	3.182	5.192	17.191	1.211	4.632	19.091

**Table 8.** (Continued) Table of Web Engine Patterns and 95% Confidence Intervals for Models 1 through 3

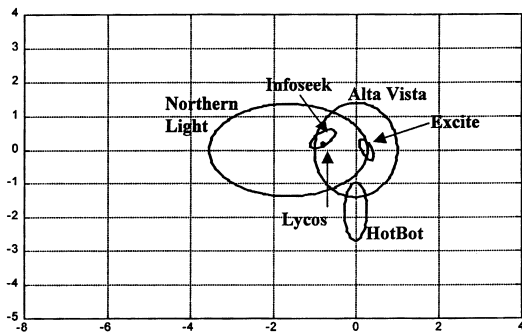
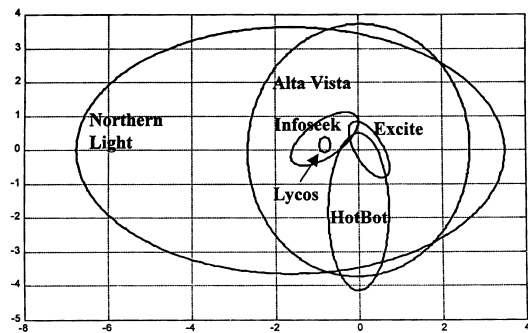
Number	Pattern	$n_{obs}$	Model 1			Model 2			Model 3		
			2.5%	50%	97.5%	2.5%	50%	97.5%	2.5%	50%	97.5%
48	010000	149	88.824	100.612	113.904	28.300	94.545	244.392	68.209	124.176	164.069
49	001111	1	0.087	0.135	0.201	0.056	0.245	1.740	0.046	0.236	1.440
50	001110	4	2.732	3.597	4.634	1.248	4.310	21.954	0.908	4.136	22.126
51	001101	0	0.111	0.166	0.261	0.043	0.204	1.232	0.054	0.245	1.673
52	001100	2	3.636	4.416	5.447	0.787	3.936	19.224	1.057	4.550	21.405
53	001011	0	0.798	1.180	1.700	0.294	1.680	9.401	0.260	1.500	8.731
54	001010	16	26.524	31.352	37.340	8.416	31.783	127.606	6.588	26.053	85.965
55	001001	1	0.986	1.417	2.314	0.304	1.759	6.933	0.320	1.580	10.571
56	001000	47	31.994	37.984	46.043	5.575	31.091	116.979	6.935	47.154	90.273
57	000111	5	0.873	1.159	1.651	0.426	1.675	9.338	0.496	1.734	7.448
58	000110	36	27.684	31.448	36.390	7.405	28.992	86.199	9.828	29.417	98.250
59	000101	0	1.058	1.401	2.210	0.370	1.423	6.813	0.405	1.862	9.071
60	000100	58	33.064	39.282	45.388	8.765	25.982	87.026	8.498	52.983	114.146
61	000011	7	7.793	10.255	13.397	8.992	51.819	36.697	3.479	11.115	38.810
62	000010	291	260.756	273.615	289.516	80.482	213.347	416.143	111.628	289.659	340.592
63	000001	9	9.430	12.453	18.961	2.735	10.615	27.031	2.946	11.953	40.908
64	000000	NA	253.938	283.431	330.301	44.592	174.756	327.190	82.149	192.704	354.000

number of unique finds (pattern 62), indicating its location far from the other engines, but still high hit rate 785/1588 (i.e., high ability to “compensate” for a distant location). The remaining ordering of engines for Model 3 are similar to those described for Model 2.

The engine locations for Model 3 are given in Table 9 (also seen in Figure 2) and suggest that the engines do carve out different locations. Northern Light, and HotBot are located the farthest distance from Alta Vista, indicating their abilities to have unique finds. Infoseek and Lycos are located “half-way” between Northern Light and Alta Vista and in a sense are “maximizing” their ability to find URLs that happen not to be found by either of the two best-performing engines. Excite’s location near Alta Vista suggests, as described more fully in § 6.3, that the additional benefit of using Excite if Alta Vista has already been used is less than that for Infoseek, despite the fact that they are “equally able” engines.

The effects of the phrase and URL covariates on dimensions 1 and 2 for Models 2 and 3 are given in Tables 6 and 7. The posterior probabilities of the effects being greater than 0 (Table 6 columns 4, 6, 8, 10) indicate that in fact domain extension, number of links,

and type and age of phrase do have a significant impact on the mean phrase location. To interpret their effects on the probability that a given URL is found, consider Table 7, which gives the coordinates of the mean phrase location for a URL with each of the given covariate attribute levels, and that of a URL with each covariate level at the baseline condition. Because under Model 2 all engines are located at the origin, and the mean phrase under the baseline condition is at (0.130, −0.080), any covariate level that brings the phrase mean closer to the origin will increase the probability a URL is found, and vice-versa. The results indicate that fewer than 10+ links and managerial phrases move the mean farther from the origin and hence lower the find probabilities. The domain extension .com, .org, and the interaction of new and managerial phrase condition move the mean phrase locations closer to the origin. The remaining covariate levels have results that depend on the ability of a given engine in each dimension. The covariate-effect results for Model 3 generally need to be examined separately for each search engine because the locations of the engines vary. This examination is straightforward, using the phrase/URL locations from Table 7 and the search en-

**Figure 2** URL Coverage by Six Popular Web Search Engines**A. Iso-Probability Ellipse: URL Discovery Probability = .5****B. Iso-Probability Ellipse: URL Discovery Probability = .33****Table 9.** 2.5%, 50%, and 97.5% Posterior Percentiles for Model 3 Engine Locations ( $\theta_1, \theta_2$ )

Engine	Dimension 1			Dimension 2		
	2.5%	50%	97.5%	2.5%	50%	97.5%
AV	0.000	0.000	0.000	0.000	0.000	0.000
HB	0.000	0.000	0.000	-2.100	-1.820	-1.340
EX	0.150	0.265	0.626	-0.321	0.010	0.187
IS	-1.080	-0.800	-0.350	-0.110	0.330	0.600
NL	-2.140	-1.640	-1.377	-0.316	-0.004	0.184
LY	-0.949	-0.799	-0.645	-0.091	0.162	0.418

gine locations from Table 9. The results for Model 3 do, however, indicate one consistent finding across search engines: The 0–5 and 6–10 link conditions move the mean phrase locations further away from the locations of the engines, decreasing the predicted probability that they are found. For the remaining cases, the results depend on the covariate and the specific engine.

A more detailed analysis for Models 2 (columns 7–9) and 3 (columns 10–12) of the  $2^6$  engine-hit patterns with observed counts  $n_{obs}$  and 2.5%, 50%, 97.5% quantiles is provided in Table 8. We observe a significant improvement in Model 3 fit for the uniques (patterns 32, 48, 56, 60, 62, 63) relative to Models 1 and 2. We also note that for 14 of the 15 engine pairs (excluding pattern 59), the 95% interval for Model 3 contains the observed value compared to 3 out of 15 for Model 1 and 12 out of 15 for Model 2. An estimate of the fraction of URLs not found is also obtained in pattern 64. The estimates under Model 2 ( $174.756/(1588 + 174.756) \approx 10\%$ ) and that for Model 3 ( $192.704/(1588 + 192.704) \approx 11\%$ ) are consistent with each other and suggest that these six engines as a whole, for these 20 phrases, cover a significant proportion of the Web. A global comparison of model fit is presented next.

### 5.3. Model Comparison and Cross-Validation

A global goodness-of-fit comparison was performed for each of Models 1 through 3 against the simple “strawman” independence models described earlier: (1) constant  $p$ , (2) different  $p$  for each engine but constant across phrases, (3) different  $p$  for each phrase but constant across engines, and (4) different  $p$  for each engine and phrase. Table 10 presents the number of parameters and the natural log of the Bayes Factor  $\log(p(M_i | Y, X)/p(M_1 | Y, X))$ , as described in Newton and Raftery (1994), comparing each models marginal likelihood  $p(M_i | Y, X)$  to the constant  $p$  model,  $p(M_1 | Y, X)$  in turn. For the independence models,  $p(M_i | Y, X)$  is evaluated at the MLE, for Models 1 through 3  $p(M_i | Y, X)$  is computed as the harmonic mean of the log-likelihood evaluated at the 7500 MCMC draws. Larger values of the Bayes factor indicate model superiority. In the end, Model 3 is selected as superior. Interestingly, we note that Model 1 does not defeat the simple model of constant  $p$  for each engine and phrase or a different  $p$  by engine.

To assess the predictive ability of our model, we employed a version of Bayesian cross-validation (Rust and Schmittlein 1985) where we dropped out in turn each of the 1588 URLs, re-estimated the model, and predicted the engine find pattern for the left out URL. To make this approach computationally feasible under a MCMC simulation structure, we employed the



**Table 10. Global Goodness-of-Fit for Various Models**

Model	# Parameters	log (Bayes Factors)
Constant $p$	1	0
Different $p$ by engine	6	816.90
Different $p$ by phrase	20	19.50
Different $p$ engine by phrase	120	980.28
Model 1	63	253.46
Model 2	140	1501.04
Model 3	149	1564.96

Note: Reported are log(Bayes Factor) comparing each model in turn to the Constant  $p$  model.

method of Bradlow and Zaslavsky (1997), in which case deletion of URLs is implemented by importance reweighting the parameter draws from the full data posterior distribution. As a result of the conditional independence structure of the likelihood given in (4), the importance reweighting scheme is trivial and computationally cheap in that each parameter draw is reweighted for URL  $jk$  by the inverse of its contribution to the likelihood, i.e.,  $p(y_{jk} | \Omega_1, \Omega_2)^{-1}$ . The total number of predictions made under this approach 9528 (1588 URLs by 6 engines) provides an adequate basis for validation. The results of the validation experiment indicated that Models 1 through 3, were able to predict 58%, 72%, and 81%, respectively, of the URL correctly (all results significant at the 0.05 level), suggesting an adequate predictive ability of the modeling approach and a substantial preference for Model 3.

## 6. Discussion and Conclusions

We set out to better understand the performance of popular Web search engines in finding marketing phrases. This required development of a model (Model 3) able to capture distinctive patterns of overlap and coverage among the engines. Furthermore, we wanted to understand how some characteristics of the phrase being searched, and of the URL being sought, would affect search outcomes. As discussed in § 5.2, two phrase characteristics (newer/older and managerial/academic) and two URL characteristics (number of links, domain type) significantly affected search engine outcomes. The effect of number-of-links happens to be consistent across engines: The more links, the more

likely the document will be located. Given the disparity in Web engine coverage patterns (as in Figure 2), the other substantive effects differed by engine. For instance, a search for an academic phrase (as opposed to managerial) aided Infoseek's prospect for locating URLs, but hindered that of Northern Light.

To elaborate on our empirical and model-based results, we conclude by addressing four simple questions:

- What search engine “works best”?
- Why do certain search engines find more URLs than other engines?
- What are the benefits to sequential search?
- How much information is still unaccounted for?

### 6.1. What Search Engine “Works Best?”

We again acknowledge that “best” here means simply locating more URLs containing the desired marketing phrase. Overall, based on the Model 3 estimates in Table 8 (and consistent with Table 1), we can make five simple statements concerning the “best engine question”:

1. Overall, for a randomly chosen marketing phrase and URL, the search engine most likely to find it is Alta Vista.
2. BUT, Northern Light is a very close second and, in fact, does slightly better than Alta Vista in finding managerial phrases.
3. HotBot is a very respectable third, locating a little over 50%–60% as many URLs as Alta Vista or Northern Light.
4. Excite and Infoseek trail more substantially, locating 20%–30% as many documents as the two leading engines.
5. Lycos found 10%–15% as many documents as the two leaders.

Of course, these findings pertain specifically to the time period of search (October 1998), the information domain of interest to us (marketing phrases), and the particular 20 phrases selected. With respect to this last restriction, however, we note that the variation in mean locations across phrases (after accounting for our covariates) was very small. (The variance across phrases in the baseline mean phrase location ( $\bar{\alpha}_1$ ,  $\bar{\alpha}_2$ ) from Table 7 is only 0.0027 for  $\bar{\alpha}_1$  and 0.002 for  $\bar{\alpha}_2$ .) That is, another set of 20 phrases drawn at random from our marketing-phrase universe would have essentially no

chance to change our findings. We next consider possible explanations for the engines' differential performance.

## 6.2. Why Do Certain Engines Find More URLs?

Research issue 3 in § 2 addressed how structural characteristics of search engines would affect the search results. Recall that some fundamental measures of this sort were provided in Table 3. Because the number of popular search engines (here, six) is small relative to the information, it was not desirable to embed these features formally in our URL-location model. Armed, however, with overall performance statistics engine-by-engine we can conduct an exploratory analysis linking search engine properties to overall search effectiveness.

Of course, the factor that looms largest in such an analysis is search engine size—i.e., the total number of Web pages indexed. Not only would it be extraordinary if “size did not matter,” but it could be well argued that “size is everything,” i.e., that the number of URLs found by search engine *A* relative to engine *B* is entirely predicted by their relative sizes. This last hypothesis was essentially tested with the independence model of search outcomes, and rejected, in § 3. In other words, our Model 3, with search engines that are somewhat distinct in the space that they cover, argues that structural characteristics beyond size may have an impact on search outcomes and motivated us to examine the full set of engine characteristics in Table 3.

Accordingly, our profiling search outcomes based on engine characteristics was done in two sequential steps. The first examined the relationship between size and overall URLs found. The second looked at any deviations from a “size/total-URLs” connection to see if those deviations are associated with other engine properties from Table 3. Essentially, the factor size represents a very simple “par” model for engine performance, and we examine in step 2 engines that overperform (and underperform) relative to size.

Table 11 reports the results of these analyses. Columns (a) and (b) show clearly that our marketing phrase search outcomes are correlated substantially with engine size ( $\rho = 0.833$ ). They also show that size is far from the only factor. Column (c) reports the ratio of URLs found to engine size. The variation in these

values shows that much more is going on than simply engines indexing more pages. Based on column (c), three engines did substantially better in locating URLs than their size would indicate: Northern Light, Alta Vista, and Infoseek. At the other extreme, not only was Lycos tied for smallest size, but it also found fewer URLs relative to its size than any of the other engines. Taking the overperformance of Northern Light and Alta Vista alone, one might suggest a convex relationship between size and URLs found (increasing returns to size) as opposed to a linear one posited in column (c), but this explanation is inconsistent with HotBot's underperformance and Infoseek's overperformance.

Instead, we sought to understand the variation in column (c) via the other search engine characteristics. Specifically, we created a simple index of search sophistication from the characteristics Depth of Search, Frames Support, Image Maps, and Learns Frequency. For each engine, we summed the binary indicators for each of the four variables (“1” = more sophisticated search, “0” = less sophisticated) and report the resulting index in Table 11 column (d).

Our measure of sophistication does a good job of explaining which engines overperform relative to their size. The three overperforming engines in column (c) are also leaders with respect to the sophistication index, although Infoseek and HotBot were admittedly tied. Overall, the correlation between overperformance in column (c) and the sophistication index in (d) is  $\rho = 0.658$ , which shows that these structural properties of search engines are substantially related to engine

**Table 11. The Relation Between Search Engine Performance and Search Engine Structural Characteristics**

Engine	(a) Total URLs Found	(b) Size (millions)	(c) URLs/Size	(d) Sophistication Index*
AV	840	140	6.0	4
NL	785	80	9.8	3
HB	468	110	4.3	2
IS	230	30	7.7	2
EX	227	55	4.1	0
LY	85	30	2.8	0

\*Sum of indicators for high performance in Depth of Search, Frames Support, Image Maps, and Learns Frequency from Table 3

performance, and in a way not reflected in the engine's size.

### 6.3. Sequential Search

One practical question of managerial interest is "which search engine should I use?" We believe that the previous two subsections summarize what our data and modeling say about that. Another practical question is "Now that I have used search engine **W** should I do an additional search, and if so what engine **Z** should I use?" Let's examine the first part of this question. Based on the results for Model 3 (Table 8), Alta Vista would be one's best single search engine choice, expected to find 48% of the marketing/management phrase URLs that exist. This is pretty good, but there is still plenty to find. More to the point, there is still plenty that can readily be found. Now turning to the second part above, if one added a second search engine after using Alta Vista, which should it be? Figure 2 by itself does not provide a clear answer. Instead, this figure shows that a putative case could be made for four of the other engines. HotBot's coverage does not overlap much with Alta Vista's, but Northern Light also does not overlap completely and covers a great deal of the URL space. Alternatively, Alta Vista will not actually find all URLs in its Figure 2 coverage area as indicated by the probability values 0.5 and 0.33 for the iso-probability curves, and many URLs exist to be found close to the origin. Excite and Infoseek are centered near the origin and accordingly are well-positioned to locate those residual URLs.

As it turns out, Northern Light is easily the best choice here for finding additional URLs. This can be established both by Table 8 using the actual search pattern finds (column 3) or Model 3's predicted search pattern outcomes. For our purposes it will suffice to simply tally the *incremental* URLs (not found by Alta Vista) for each of the remaining five engines. These are, in order, Northern Light (actual incremental = 443, predicted incremental using Model 3 = 468), HotBot (actual = 271, predicted = 259), Infoseek (actual = 136, predicted = 124), Excite (actual = 110, predicted = 109), and Lycos (actual = 35, predicted = 42). Thus we conclude that in general it is important to consider both overall coverage ability and overlap in selecting combinations of search engines.

### 6.4. How Much Information Is Still Unaccounted For?

We have seen that combined search outcomes from multiple engines improves greatly on any one engine's performance. Yet, how much marketing information remains unlocated, even after using all six engines? For our 20 marketing phrases, the results in Table 8 provide an answer to that question. Based on the estimate from Model 3, the fraction of total relevant URLs *missed* by all six search engines is just 10.8% (192,704/1,786,967). Given the small variation in phrase location for our 20 marketing phrases searched, the reader should feel confident that the search engines cover about 90% of what exists to be found for these kind of phrases.

This is quite different—and much better—than the Web coverage estimated by Lawrence and Giles (1998) for their scientific-phrase searches. There, the six search engines were estimated to cover about 60% of the indexable URLs. In their updated 1999 article, this figure is even lower and, as they suggest, states that "engines aren't keeping up." Two explanations for the discrepancy across studies suggest themselves readily. First, the estimated number of URLs *not* found could be highly sensitive to the particular model specification selected. As we have seen, our marketing data reject the independent binomial model used by Lawrence and Giles because that model does not effectively capture the patterns of overlap for sets of engines. So if we had to select one model to estimate the size of the Web, we would propose our Model 3 as a more appealing approach. Nonetheless, if the estimated Web size is so sensitive to model specification, one might well question the ability of *any* of these models to provide a reliable estimate—at least without exhaustive checking of individual assumptions. Fortunately, this situation has not arisen. While we do not recommend using the independent binomial model, its estimate of cumulative URL coverage by our six search engines (across all 20 phrases) is 89.6%—very close to the value found using our Model 3. In short, while the independent binomial model methodology is suspect, it too indicates high coverage of marketing information. Accordingly, the differences between our results

do not stem from hypersensitivity to model assumptions.

This brings us to the second explanation: namely, that these kinds of marketing/management documents are relatively easy to locate. While we cannot prove this, it is a reasonable hypothesis. Parts of the Web are of course much more "active" than others, with respect to both availability of hyperlinks from one document to another, and the degree of use of those links. This interconnectedness is the key to a search engine's performance. Documents containing our marketing research and marketing management phrases may well be relatively active in this respect. That is, other Web documents may be particularly likely to link to the commercial sites, educational sites, or organizations' sites that contain the information. While our results do not say that Web-based marketing information providers can simply count on search engines bringing multitudes to their location, they do indicate that much of the marketing information currently on the Web can be located readily—if one uses multiple search engines.

### 6.5. Limitations and Future Research

This study is limited in that it used six specific search engines (the ones discussed most often in the popular press and examined in other systematic studies) during one specific time period (October 1998) to search for Web pages containing each of 20 specific marketing/management phrases (obtained by surveying common marketing reference sources). In addition, our analysis treats each Web page containing the search phrase as fully and equally valued, i.e., we do not judgmentally assess how "good" a page is (for an unspecified search purpose). To be sure, we are skeptical of attempts to do this assessment. In this area, we essentially assume that the searcher is able to articulate what is in fact being sought. Accordingly, we also do not evaluate the heuristics used by search engines to rank URLs reported in a search.

Changing any of these study design elements may materially affect the empirical results. We note in particular that the relative performance of search engines has been observed to vary over time (Lawrence and Giles 1999). We are less concerned about selection of the search phrases because search phrase locations did

not vary substantially across the 20 examined here. Our investigation of the role played by the search phrase characteristics and search engine characteristics is limited by judgmental coding of the former and the need to rely on nonproprietary factors for the latter. The study found significant effects for each despite these limitations.

We hope that this paper has provided some useful data, and some insight, concerning use of Web search engines to find managerial information. Our proposed (and validated) spatial coverage model provides both a "snapshot summary" of the search engines vis-a-vis each other (as in Figure 2), and also yields predictions regarding cumulative performance of engine combinations. We have shown that certain characteristics of search engines, search phrases, and URL locations affect the probability that a given engine will locate a given URL. Of course, the search engines themselves will evolve, and patterns of coverage and overlap can change accordingly. This evolution (and its causes) will be interesting to explore in future research. We are hopeful that our model framework will continue to provide a basis for summarizing these patterns. The marketing information base on the Web is evolving—and expanding—very rapidly. For many purposes it has (and will continue to) outstrip the ability of managed directories, lists, and the like to provide focused useful direction, or even to keep up with change. The Web search engines are well positioned to meet this challenge in the future, and currently they collectively—if not individually—can do so for the kind of marketing information examined here.<sup>1</sup>

### Appendix

Inferences for parameters  $\Omega_1$  and  $\Omega_2$  are obtained from the marginal posterior distributions

$$p(\Omega_1 | Y) \propto \int p(Y | \Omega_1) p(\Omega_1 | \Omega_2) p(\Omega_2) d\Omega_2, \text{ and} \quad (7)$$

$$p(\Omega_2 | Y) \propto \int p(Y | \Omega_1) p(\Omega_1 | \Omega_2) p(\Omega_2) d\Omega_1, \quad (8)$$

defined by the likelihood and priors given in (4) and (5). The non-conjugate likelihood and prior structure prevent closed-form inte-

<sup>1</sup>The authors thank the Special Issue editor, area editor, and three anonymous reviewers for useful suggestions.

gration of (7) and (8). The approach taken here to solve these intractable integrals is iterative simulation via a Markov chain Monte Carlo (MCMC) sampler. This approach states that under certain regularity conditions, samples from (7) and (8) may be obtained by repeatedly sampling values  $\Omega_1^{(t+1)}$  from the conditional distribution  $p(\Omega_1 | Y, \Omega_2^{(t)})$  and  $\Omega_2^{(t+1)}$  from  $p(\Omega_2 | Y, \Omega_1^{(t+1)})$  until convergence, and treating draws thereafter as draws from the desired marginal posterior distributions.

Unfortunately, for our model the conditional distributions  $p(\Omega_1 | Y, \Omega_2^{(t)})$  necessary to straightforwardly implement an MCMC sampler cannot be sampled from directly. We note that the conditional distribution of  $p(\Omega_2 | Y, \Omega_1^{(t+1)})$  can be sampled directly because of the conjugate multivariate normal – Inverse Wishart prior structure chosen for  $\Omega_1$ . To sample  $\Omega_1^{(t+1)}$  from  $p(\Omega_1 | Y, \Omega_2^{(t)})$  we implemented a Metropolis-Hastings jumping algorithm (Hastings 1970), where for each parameter that was unconstrained, we utilized a symmetric Gaussian jumping distribution with mean at the previously drawn value  $\Omega_1^{(t)}$ , and variance set to provide a high acceptance rate. For those parameters constrained to the positive real line (variances,  $u$ , and  $\theta_{31}$  in Model 3), we utilized a Gamma distribution kernel with shape parameter  $k(\Omega_1^{(t)})^2$  and scale parameter  $k\Omega_1^{(t)}$ , which has mean equal to the previous draw  $\Omega_1^{(t)}$  and variance  $1/k$ . The value of  $k$  was set differently for each parameter to obtain an adequate acceptance rate.

Three independent streams for each of the three models were run using overdispersed starting values obtained from an initial run. Computing times for Models 1 through 3 were 3, 12, and 14 seconds, respectively, per iteration on an HP7000 workstation using Fortran 77 code.

## References

- Alba, Joseph, John Lynch, Barton Weitz, Chris Janiszewski, Richard Lutz, Alan Sawyer, Stacy Wood. 1997. Interactive home shopping: incentives for consumers, retailers, and manufacturers to participate in electronic marketplaces. *J. Marketing* **61** (July) 38–53.
- Andersen, Erling B. 1973. Conditional inference for multiple-choice questionnaires. *British J. Math. Statist. Psycho.* **26** 31–44.
- Bakos, Yannis. 1997. Reducing buyer search costs: implications for electronic marketplaces. *Management Sci.* **43**(12) 1676–1692.
- Beatty, Sally. 1998. NBC puts its firepower behind snap! *The Wall Street Journal*. September 15.
- Bennett, Peter D. (ed.). 1995. *Dictionary of Marketing Terms* (2nd edition). NTC Business Books, Lincolnwood, IL.
- Bradlow, E. T., A. M. Zaslavsky. 1997. Case influence analysis in Bayesian inference. *J. Comput. Graphical Statist.* **6** (September) 314–331.
- Burke, Ray. 1996. Virtual shopping: breakthrough in marketing research. *Harvard Bus. Review* **74**(2) 120–131.
- Clemente, Mark N. 1992. *The Marketing Glossary*. AMACOM (American Management Association), New York.
- Elrod, Terry. 1988. Choice map: inferring a product-market map from panel data. *Marketing Sci.* **7** 21–40.
- Feldman, Susan. 1998. Web search services in 1998: trends and challenges. *Searcher* **6** Web document address <http://www.infotoday.com/searcher/jun/story2.htm#chart>.
- Feller, William. 1968. *An Introduction to Probability Theory and Its Applications*. 3rd edition. Wiley, New York.
- Gelfand, Alan E., Susan E. Hills, Amy Racine-Poon, Adrian F. M. Smith. 1990. Illustration of Bayesian inference in normal data models using Gibbs sampling. *J. American Statist. Assoc.* **85** 972–985.
- Gelman, Andrew, Donald B. Rubin. 1992. Inference from iterative simulation using multiple sequences. *Statist. Sci.* **7** 457–511.
- Hastings, W. K. 1970. Monte Carlo sampling methods using Markov chains and their applications. *Biometrika* **57** 97–109.
- Hoffman, Donna L., William D. Kalsbeek, Thomas P. Novak. 1996. Internet and Web use in the United States: baselines for commercial development. *Comm. ACM* **39** (December) 36–46.
- , Thomas P. Novak. 1996. Marketing in hypermedia computer-mediated environments: conceptual foundations. *J. Marketing* **60**(July) 50–68.
- Kamakura, Wagner A., Rajendra K. Srivastava. 1984. Predicting choice shares under conditions of brand interdependence. *J. Marketing Res.* **21** 420–434.
- Lawrence, Steve, C. Lee Giles. 1998. Searching the World Wide Web. *Science* **280** (3) 98–100.
- , —. 1999. Accessibility of information on the Web. *Nature* **400**(July 8) 107–109.
- Newton, Michael A., Adrian E. Raftery. 1994. Approximate Bayesian inference with the weighted likelihood bootstrap. *J. Royal Statist. Soc. Series B.* **56**(1) 3–48.
- PC Magazine. 1998. Web search sites: metasearch sites (December 1).
- Rasch, G. 1960. *Probabilistic Models for Some Intelligence and Attainment Tests*. Nielson and Lydiche (for Danmarks Paedagogiske Institut), Copenhagen, Denmark.
- . 1966. An item analysis which takes individual differences into account. *British J. Math. Statist. Psych.* **19** (Part 1) 49–57.
- Ritter, Christian, Martin A. Tanner. 1992. Facilitating the Gibbs sampler: the Gibbs stopper and the griddy-Gibbs sampler. *J. Am. Statist. Assoc.* **87** 861–868.
- Rossi, Peter E., Robert E. McCulloch, Greg M. Allenby. 1996. The value of purchase history data in target marketing. *Marketing Sci.* **15**(4) 321–340.
- Rust, Roland T., David C. Schmittlein. 1985. A Bayesian cross-validated likelihood method for comparing alternative specifications of quantitative models. *Marketing Sci.* **4**(1) 20–40.
- Selberg, Erik, Oren Etzioni. 1996. Multi-engine search and comparison using the MetaCrawler. In *Proceedings of the Fourth International World Wide Web Conference*. Boston, MA, 195.
- Sullivan, Danny. 1998. Search engine watch: search engines feature chart. Web document address <http://searchenginewatch.interest.com/webmasters/features.html>.
- The Wall Street Journal*. 1998. Web's vastness foils even best search engines. (April 3).

This paper was received June 18, 1998, and has been with the authors 7 months for 4 revisions; processed by Greg Allenby.