

# Random Sampling from a Search Engine's Index

ZIV BAR-YOSSEF

*Technion and Google Haifa, Haifa, Israel*

AND

MAXIM GUREVICH

*Technion, Haifa, Israel*

**Abstract.** We revisit a problem introduced by Bharat and Broder almost a decade ago: How to sample random pages from the corpus of documents indexed by a search engine, using only the search engine's public interface? Such a primitive is particularly useful in creating objective benchmarks for search engines.

The technique of Bharat and Broder suffers from a well-recorded bias: it favors long documents. In this article we introduce two novel sampling algorithms: a lexicon-based algorithm and a random walk algorithm. Our algorithms produce *biased* samples, but each sample is accompanied by a *weight*, which represents its bias. The samples, in conjunction with the weights, are then used to *simulate* near-uniform samples. To this end, we resort to four well-known Monte Carlo simulation methods: *rejection sampling*, *importance sampling*, the *Metropolis–Hastings* algorithm, and the *Maximum Degree* method.

The limited access to search engines force our algorithms to use bias weights that are only “approximate”. We characterize analytically the effect of approximate bias weights on Monte Carlo methods and conclude that our algorithms are *guaranteed* to produce near-uniform samples from the search engine's corpus. Our study of approximate Monte Carlo methods could be of independent interest.

Experiments on a corpus of 2.4 million documents substantiate our analytical findings and show that our algorithms do not have significant bias towards long documents. We use our algorithms to collect comparative statistics about the corpora of the Google, MSN Search, and Yahoo! search engines.

---

A preliminary version of this article appeared in the *Proceedings of the 15th International World-Wide Web Conference (WWW2006)* (Edinburgh, Scotland, May 23–26), ACM, New York, 2006. pp. 367–376.

The research of Z. Bar-Yossef was supported by the European Commission Marie Curie International Re-integration Grant.

Part of this work was done while M. Gurevich was at the IBM Research Lab in Haifa, Israel.

Authors' addresses: Z. Bar-Yossef, Department of Electrical Engineering, Technion, Haifa 32000, Israel, and Google Haifa Engineering Center, Haifa, Israel, e-mail: zivby@ee.technion.ac.il; M. Gurevich, Department of Electrical Engineering, Technion, Haifa 32000, Israel, e-mail: gmax@tx.technion.ac.il.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or direct commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or [permissions@acm.org](mailto:permissions@acm.org).

© 2008 ACM 0004-5411/2008/10-ART24 \$5.00 DOI 10.1145/1411509.1411514 <http://doi.acm.org/10.1145/1411509.1411514>

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

General Terms: Measurement, Algorithms

Additional Key Words and Phrases: Benchmarks, sampling, search engines, size estimation

#### ACM Reference Format:

Bar-Yossef, Z., and Gurevich, M. 2008. Random sampling from a search engine's index. J. ACM 55, 5, Article 24 (October 2008), 74 pages. DOI = 10.1145/1411509.1411514 <http://doi.acm.org/10.1145/1411509.1411514>

## 1. Introduction

Search engine size wars are almost as old as search engines themselves. Self reports of search engine sizes are frequently contradictory, leading to heated public debates. See Price [2005], Mayer [2005], Battelle [2005], and Sullivan [2005] for an account of these incidents. These debates underscore the lack of widely acceptable benchmarks for search engines.

Current evaluation methods for search engines [Gordon and Pathak 1999; Hawking et al. 2001; Can et al. 2004] are labor-intensive and are based on anecdotal sets of queries or on fixed TREC data sets. Such methods do not provide statistical guarantees about their results. Furthermore, when the query test set is known in advance, search engines can manually adapt their results to guarantee success in the benchmark.

In an attempt to come up with reliable automatic benchmarks for search engines, Bharat and Broder [1998] proposed the following problem: can we sample random documents from the corpus of documents indexed by a search engine, using only the engine's public interface? Unlike the manual methods, random sampling offers statistical guarantees about its test results. It is important that the sampling is done only via the public interface, and not by requesting the search engine itself to collect the sample documents, because we would like the tests to be objective and not to rely on the goodwill of search engines. Furthermore, search engines seem to be reluctant to allow random sampling from their corpus, because they do not want third parties to dig into their data.

Random sampling can be used to test the quality of search engines under a multitude of criteria: (1) *Overlap and relative sizes*: we can find out, for example, what fraction of the documents indexed by Yahoo! are also indexed by Google and vice versa. Such size comparisons can indicate which search engines have better recall for narrow-topic queries. (2) *Topical bias*: we can identify themes or topics that are overrepresented or underrepresented in the corpus. (3) *Freshness evaluation*: we can evaluate the freshness of the corpus, by estimating the fraction of "dead links" it contains. (4) *Spam evaluation*: using a spam classifier, we can find the fraction of spam pages in the corpus. (5) *Security evaluation*: using an anti-virus software, we can estimate the fraction of indexed documents that are contaminated by viruses.

Random sampling can be used not only by third parties, but also by search service providers themselves to test the quality of their search engines and to compare them against the competition. The results of the comparison can then be used to improve the quality of the index. Although providers have full access to their own indices, using the public interface has several benefits: (1) It may be technically simpler to

use this interface rather than to implement a procedure that directly samples from the corpus. (2) Applying the same sampling procedure to both the provider's own corpus and the competitors' corpora yields comparable results.

**1.1. THE BHARAT–BRODER APPROACH.** Bharat and Broder proposed the following simple algorithm for uniformly sampling documents from a search engine's corpus. The algorithm successively formulates “random” queries, submits the queries to the search engine, and picks uniformly chosen documents from the result sets returned. In order to construct the random queries, the algorithm requires the availability of a *lexicon* of terms that appear in web documents. Each term in the lexicon should be accompanied by an estimate of its frequency on the web. Random queries are then formulated by randomly selecting a few terms from the lexicon, based on their estimated frequencies, and then taking their conjunction or disjunction. The lexicon is constructed in a pre-processing step by crawling a large corpus of documents (Bharat and Broder crawled the Yahoo! directory).

As Bharat and Broder noted in the original article [Bharat and Broder 1998] and was later confirmed by subsequent studies [Cheney and Perry 2005; Véronis 2005], the method suffers from severe biases. The most significant bias is towards long, “content-rich”, documents, simply because these documents match many more queries than short documents. An extreme example is online dictionaries and word lists (such as the ispell dictionaries), which will be returned as the result of almost any conjunctive query composed of unrelated terms [Cheney and Perry 2005; Véronis 2005].

Another problem is that search engines do not allow access to the full list of results, but rather only to the top  $k$  ones (where  $k$  is usually 1,000). Therefore, the Bharat–Broder algorithm may be biased towards documents with high static rank, if the random queries it generates tend to return more than  $k$  results, which is usually the case for disjunctive queries. To alleviate this problem, Bharat and Broder used the estimated term frequencies to choose queries that are unlikely to return more than  $k$  results. As the number of pages indexed by search engines grew by orders of magnitude over the past decade, while  $k$  has not, this technique for query selection has become ineffective. It is almost impossible to find terms so that disjunctive queries composed of them will return less than  $k$  results. By focusing only on such queries, the algorithm is limited to sampling only from a small subset of the web. It is easier to construct conjunctive queries with less than  $k$  results simply by adding more random terms. This, however, also increases the fraction of dictionaries and word lists among the results.

**1.2. OUR CONTRIBUTIONS.** We propose two novel algorithms for sampling pages from a search engine's corpus. Both algorithms use a lexicon to formulate random queries, but unlike the Bharat–Broder approach, do not need to know term frequencies. The first algorithm, like the Bharat–Broder algorithm, requires a pre-processing step to build the lexicon. The second algorithm is based on a random walk on a virtual graph defined over the documents in the corpus. This algorithm does not need to build the lexicon a priori, but rather creates it “on the fly”, as the random walk proceeds from document to document.

Both algorithms share the same basic framework. The algorithm first produces *biased* samples. That is, some documents are more likely to be sampled than others. Yet, each sample document  $x$  is accompanied by a corresponding “weight”  $w(x)$ , which represents the bias in the sample  $x$ . The weights allow us to apply *stochastic*

*simulation* methods on the samples and consequently obtain *uniform*, unbiased, samples from the search engine’s corpus.<sup>1</sup>

A simulation method accepts samples taken from a *trial distribution*  $p$  and simulates sampling from a *target distribution*  $\pi$ . In order to carry out the simulation, the simulator needs to be able to compute a “bias weight”  $w(x) = \pi(x)/p(x)$ , for any given instance  $x$ . The simulator uses these weights to “promote” certain samples from  $p$ , while “demoting” other samples, thereby eliminating the initial bias in the trial distribution. The simulation has some overhead, which depends on how far  $p$  and  $\pi$  are from each other. In our case  $\pi$  is the uniform distribution over the search engine’s corpus, while  $p$  is some other easy-to-sample-from distribution over the corpus.<sup>2</sup> We employ four Monte Carlo simulation methods: *rejection sampling*, *importance sampling*, the *Metropolis–Hastings* algorithm, and the *Maximum Degree* method.

One technical difficulty in applying simulation methods in our setting is that the weights produced by our samplers are only *approximate*. To the best of our knowledge, stochastic simulation with approximate weights has not been addressed before. We are able to show that all four Monte Carlo methods still work even when provided with approximate weights. The distribution of the samples they generate is no longer identical to the target distribution  $\pi$ , but is rather only close to  $\pi$ . We provide extensive theoretical analysis of the effect of the approximate bias weights on the quality and the performance of the Monte Carlo methods. This study may be of independent interest.

**1.2.1. Pool-Based Sampler.** A *query pool*, or a *lexicon*, is a collection of queries. Our pool-based sampler assumes knowledge of some query pool  $\mathcal{P}$ . The terms constituting queries in the pool can be collected by crawling a large corpus, like the Yahoo!<sup>3</sup> or ODP<sup>4</sup> directories. We stress again that knowledge of the frequencies of these terms is not needed.

The *degree* of a document  $x$  in the corpus is the number of queries from the pool that it matches. The “document degree distribution” is a distribution over the corpus, where each document is selected proportionally to its degree. The inner subroutine of the pool-based sampler generates samples from the document degree distribution. The corresponding bias weights are inverse-proportional to the document degrees, and are thus easy to calculate. The outer subroutine of the pool-based sampler then applies a Monte Carlo method (e.g., rejection sampling) on the samples from the document degree distribution in order to simulate uniform samples.

How does the algorithm generate samples from the document degree distribution? The first solution that comes to mind is as follows. The algorithm picks a random query from the pool, submits the query to the search engine, and then selects a random document from the result set returned. Indeed, documents with high degree match more queries, and are thus more likely to be sampled than documents with low degree. However, the resulting sampling distribution is not exactly the

<sup>1</sup>In fact, our algorithms are more general. They can be used to generate samples from any target distribution over the corpus, not only the uniform one.

<sup>2</sup>For some instances of  $\pi$ , the specific trial distribution  $p$  we propose may be too far from  $\pi$ , and thus inefficient. In such cases, development of alternative trial distributions may be required.

<sup>3</sup>[dir.yahoo.com](http://dir.yahoo.com).

<sup>4</sup>[dmoz.org](http://dmoz.org).

document degree distribution, as the chance of a document to be selected depends also on the number of results returned on queries that this document matches. For example, if documents  $x$  and  $x'$  have both degree 1, but  $x$  matches a query  $q$  with 100 results, while  $x'$  matches a query  $q'$  with 50 results, then the probability of  $x'$  to be sampled is twice as high as the probability of  $x$  to be sampled.

To alleviate this problem, the algorithm does not select the query uniformly at random, but rather proportionally to its *cardinality*. The cardinality of a query  $q$  is the number of results it has. It can be shown that if the algorithm samples queries from the pool proportionally to their cardinalities, then by selecting random documents from their result sets, the algorithm could have obtained samples from the document degree distribution. Sampling queries according to their cardinality is tricky, though, because we do not know a priori the cardinality of queries. What we do instead is sample queries from the pool uniformly, and then *simulate* sampling from the cardinality distribution. To this end, we use Monte Carlo methods again. Hence, Monte Carlo methods are used twice: first to generate the random queries and then to generate the uniform documents.

To demonstrate how the pool-based sampler works, consider a corpus consisting of only 100 documents and a query pool with two queries  $q_1$  and  $q_2$ , whose cardinalities are 99 and 2, respectively. Suppose the result sets of  $q_1$ ,  $q_2$  share a single document  $x$ . The sampler chooses one of  $q_1$ ,  $q_2$  uniformly at random and then applies an acceptance–rejection procedure, in which  $q_1$  is accepted with probability  $99/100$  and  $q_2$  is accepted with probability  $2/100$ . If the query is accepted, it is submitted to the search engine, and a random document is chosen from its result set. A second acceptance–rejection procedure is applied on this document. If it is the document  $x$ , it is accepted with probability  $1/2$ , and otherwise it is accepted with probability 1. If the document is accepted, it is output as a sample. It can be verified that all 100 documents in the corpus are equally likely to be sampled.

We rigorously analyze the pool-based sampler and identify the important properties of the query pool that make this technique accurate and efficient. We find that using a pool of exact phrase queries of a certain length is much more preferable to using conjunctive or disjunctive queries, like the ones used by Bharat and Broder.

**1.2.2. Random Walk Sampler.** We present another sampler, which also uses a query pool, but does not need to construct it a priori. This sampler performs a random walk on a virtual graph defined over the documents in the corpus. The limit distribution of this random walk is the uniform distribution over the documents, and thus if we run the random walk for sufficiently many steps, we are guaranteed to obtain near-uniform samples from the corpus.

The graph is defined as follows: two documents are connected by an edge if and only if they match the same query from the pool. This means that if one submits the shared query to the search engine, both documents are guaranteed to belong to the result set. Running a random walk on this graph is simple: we start from an arbitrary document, at each step choose a random query that the current document matches, submit this query to the search engine, and move to a randomly chosen document from the query's result set.

The random walk as defined does not converge to the uniform distribution. In order to make it uniform, we apply either the Metropolis–Hastings algorithm or the Maximum Degree method. We provide careful analysis of the random walk sampler. Like the pool-based sampler, this sampler too is guaranteed to produce



near-uniform samples. However, theoretical analysis of its performance suggests that it is less efficient than the pool-based sampler.

**1.2.3. *Experimental Results.*** To validate our techniques, we crawled 2.4 million English pages from the ODP hierarchy, and built a search engine over these pages. We used these pages to create the query pool needed for our pool-based sampler and for the Bharat–Broder sampler.

We ran our two samplers as well as the Bharat–Broder sampler on this search engine, and calculated bias towards long documents and towards highly ranked documents. As expected, the Bharat–Broder sampler was found to have significant bias. On the other hand, our pool-based sampler had no bias at all, while the random walk sampler only had a small negative bias towards short documents.

We then ran our pool-based sampler on Google, MSN Search, and Yahoo!. As a query pool, we used 5-term phrases extracted from English pages at the ODP hierarchy. The samples collected enabled us to produce estimates, as of May 2006, of the relative sizes of these search engines as well as interesting statistics about their freshness, their domain name distribution, and their coverage of dynamic URLs.

The rest of the article is organized as follows. In Section 2, we review some related work. Section 3 overviews some tools from probability theory and statistics used in our analysis. In Section 4, we briefly review the four Monte Carlo simulation methods we use. In Section 5, we analyze the effect of approximate bias weights on Monte Carlo simulation. In Section 6, we describe a formal framework for studying search engine samplers. In Sections 7 and 8, we outline in detail our two samplers. Section 9 includes our experimental results, and Section 10 some concluding remarks.

In order to avoid disturbing the flow of the article, most proofs are postponed to the appendix.

## 2. *Related Work*

Apart from Bharat and Broder, several other studies used queries to search engines to collect random samples from their corpora. Queries were either manually crafted [Bradlow and Schmittlein 2000], collected from user query logs [Lawrence and Giles 1998], or selected randomly using the technique of Bharat and Broder [Gulli and Signorini 2005; Cheney and Perry 2005]. Assuming search engine corpora are independent and uniformly chosen subsets of the web, estimates of the sizes of search engines and of the indexable web have been derived. Due to the bias in the samples, though, these estimates lack any statistical guarantees. Dobra and Fienberg [2004] showed how to avoid the unrealistic independence and uniformity assumptions, but did not address the sampling bias. We believe that their methods could be combined with ours to obtain accurate size estimates.

Several studies [Lawrence and Giles 1999; Henzinger et al. 1999, 2000; Bar-Yossef et al. 2000; Rusmevichientong et al. 2001] developed methods for sampling pages from the indexable web. Such methods can be used to also sample pages from a search engine’s corpus. Yet, since these methods try to solve a harder problem, they also suffer from various biases, which our method does not have. It is interesting to note that the random walk approaches of Henzinger et al. [2000] and Bar-Yossef et al. [2000] implicitly use importance sampling and the Maximum Degree method,

respectively, to make their samples near-uniform. Yet, the bias they suffer towards pages with high in-degree is significant.

Anagnostopoulos et al. [2005] proposed an enhancement to index architecture that could support random sampling from the result sets of broad queries. This is very different from what we do in this article: our techniques do not propose any changes to current search engine architecture and do not rely on internal data of the search engine; moreover, our goal is to sample from the whole corpus and not from the result set of a particular query.

A recent study by Broder et al. [2006] presents an algorithm for accurately estimating the absolute size of a search engine's corpus, using the engine's public interface. The cleverly crafted estimator uses our techniques to generate uniform samples from the search engine's corpus.

The samplers proposed in this article, as well as in the paper of Broder et al. [2006], suffer from a bias caused by inaccurate approximation of document degrees. In this article, we analyzed and quantified the resulting bias. In our subsequent work [Bar-Yossef and Gurevich 2007], we showed a method for overcoming this bias.

### 3. Preliminaries

In this section, we outline our notations and conventions and review some tools from statistics that will be handy in our analysis.

**3.1. CONVENTIONS AND NOTATIONS.** Sets are denoted by uppercase calligraphic letters:  $\mathcal{D}$ ,  $\mathcal{Q}$ ,  $\mathcal{P}$ . Elements in sets are denoted by lowercase Roman letters:  $x$ ,  $y$ ,  $q$ . Random variables are denoted by uppercase Roman letters:  $X$ ,  $Y$ ,  $Q$ . Probability distributions are denoted by small italicized letters, Roman or Greek:  $p$ ,  $q$ ,  $\pi$ .

All probability spaces in this article are discrete and finite. A *probability distribution*  $p$  on a finite probability space  $\mathcal{U}$  is a function  $p : \mathcal{U} \rightarrow [0, 1]$  such that  $\sum_{x \in \mathcal{U}} p(x) = 1$ . The *support* of  $p$  is defined as:

$$\text{supp}(p) = \{x \in \mathcal{U} \mid p(x) > 0\}.$$

A subset  $\mathcal{E}$  of the probability space  $\mathcal{U}$  is called an *event*. For an event  $\mathcal{E} \subseteq \mathcal{U}$ , we define  $p(\mathcal{E})$  to be the probability of this event under  $p$ :

$$p(\mathcal{E}) = \sum_{x \in \mathcal{E}} p(x).$$

A *random variable* with distribution  $p$  and range  $\mathcal{V}$  is a function  $X : \mathcal{U} \rightarrow \mathcal{V}$ . Unless stated otherwise, when we refer to a random variable with distribution  $p$ , we mean the identity random variable:  $\mathcal{V} = \mathcal{U}$  and  $X(x) = x$ , for all  $x \in \mathcal{U}$ .

Given a predicate  $f : \mathcal{V} \rightarrow \{0, 1\}$ ,  $f(X)$  is the following event:

$$f(X) = \{x \in \mathcal{U} \mid f(X(x)) = 1\}.$$

The probability of the event  $f(X)$  under  $p$  is denoted  $\Pr_p(f(X))$ .

When the range of the random variable is  $\mathcal{V} = \mathbb{R}$ , we can define the *expectation* of  $X$  under  $p$ :

$$\mathbb{E}_p(X) = \sum_{x \in \mathcal{U}} p(x) X(x).$$

The *variance* of  $X$  is defined as:

$$\text{var}_p(X) = \mathbb{E}_p((X - \mathbb{E}_p(X))^2) = \mathbb{E}_p(X^2) - (\mathbb{E}_p(X))^2.$$

The *standard deviation* of  $X$  is the square root of its variance:

$$\sigma_p(X) = \sqrt{\text{var}_p(X)}.$$

We define the *mean deviation* of a random variable  $X$  to be its expected absolute deviation from its mean:

$$\text{dev}_p(X) = \mathbb{E}_p(|X - \mathbb{E}_p(X)|).$$

It follows from the Cauchy–Schwartz inequality that the mean deviation is always bounded by the standard deviation:

PROPOSITION 3.1. *For any random variable  $X$ ,  $\text{dev}_p(X) \leq \sigma_p(X)$ .*

The *normalized mean deviation* of  $X$  is the mean deviation, normalized by the mean:

$$\text{ndev}_p(X) = \frac{\text{dev}_p(X)}{\mathbb{E}_p(X)}.$$

The *covariance* of two random variables  $(X, Y)$  with joint distribution  $p$  is defined as:

$$\text{cov}_p(X, Y) = \mathbb{E}_p(XY) - \mathbb{E}_p(X)\mathbb{E}_p(Y).$$

The *average* of a function  $g : \mathcal{D} \rightarrow \mathbb{R}$  is defined as:

$$\text{avg}_{x \in \mathcal{D}} g(x) = \frac{1}{|\mathcal{D}|} \sum_{x \in \mathcal{D}} g(x).$$

Throughout, we assume knowledge of basic probability theory. Everything we use can be found in any standard textbook on the subject.

3.2. TOTAL VARIATION DISTANCE. The *total variation distance* between two distributions  $p, q$  on the same space  $\mathcal{U}$  is defined as:

$$\|p - q\| = \frac{1}{2} \sum_{x \in \mathcal{U}} |p(x) - q(x)|.$$

We use total variation distance, because it has some nice properties described below. Yet, other statistical distance measures, like the Kullback–Leibler divergence could have been used as well.

The following is a standard characterization of the total variation distance:

LEMMA 3.2. *Let  $p, q$  be two distributions on the same probability space  $\mathcal{U}$ . Then,*

$$\|p - q\| = \max_{\mathcal{E} \subseteq \mathcal{U}} |p(\mathcal{E}) - q(\mathcal{E})|.$$

3.3. WALD’S IDENTITY. Suppose we invoke a sampler  $T$  times, where  $T$  is a random variable, and each invocation requires  $n$  search engine queries in expectation. What is then the expected total number of search engine queries made? As  $T$  is a random variable, simple linearity of expectation cannot be used in the



analysis. The following identity from statistical sequential analysis shows that the expectation equals  $n \cdot \mathbb{E}(T)$ .

**THEOREM 3.3 (WALD'S IDENTITY).** *Let  $X_1, X_2, \dots$  be an infinite sequence of independent and identically distributed random variables with mean  $\mu$ . Let  $T$  be a random variable on  $\{0, 1, 2, \dots\}$ , for which the event  $\{T = k\}$  is independent of  $X_{k+1}, X_{k+2}, \dots$  for all  $k$  ( $T$  is called a stopping time random variable). We further assume  $\mathbb{E}(T) < \infty$ . Then,*

$$\mathbb{E} \left( \sum_{i=1}^T X_i \right) = \mu \mathbb{E}(T).$$

A proof of Wald's identity can be found, for example, in the textbook of Siegmund [1985], Section 2.2.

#### 4. Monte Carlo Methods

We briefly review the four Monte Carlo simulation methods we use in this article. For a more elaborate overview, refer to the textbook of Liu [2001].

**4.1. BASIC FRAMEWORK.** The basic question addressed in stochastic simulation is the following. There is a *target distribution*  $\pi$  on a space  $\mathcal{U}$ , which is hard to sample from directly. On the other hand, there is an easy-to-sample-from *trial distribution*  $p$  on the same space  $\mathcal{U}$ . Can we then somehow use the samples from  $p$  to *simulate* sampling from  $\pi$ ? A Monte Carlo simulator is a procedure, which given samples from  $p$  generates samples from  $\pi$ . In order to carry out the simulation, the simulator requires access to three “oracle procedures”, which we describe next.

The first procedure, `getSamplep()`, generates samples from the trial distribution  $p$ . Each invocation returns a single sample  $X$  from  $p$ . Successive invocations generate independent and identically distributed samples  $X_1, X_2, \dots$ .

The two other oracle procedures are used to calculate *unnormalized forms* of the distributions  $\pi$  and  $p$ :

**Definition 4.1 (Unnormalized Form of a Distribution).** Let  $\pi$  be a distribution on a space  $\mathcal{U}$ . An *unnormalized form* of  $\pi$  is a function  $\hat{\pi} : \mathcal{U} \rightarrow [0, \infty)$ , which equals  $\pi$  up to a *normalization constant*  $Z_{\hat{\pi}} > 0$ . That is,

$$\forall x \in \mathcal{U}, \quad \hat{\pi}(x) = \pi(x) \cdot Z_{\hat{\pi}}.$$

$\hat{\pi}(x)$  is a relative weight, which represents the probability of  $x$  to be chosen in the distribution  $\pi$ . For example, if  $\pi$  is the uniform distribution on  $\mathcal{U}$ , then all elements are equally likely to be selected. Hence, the straightforward unnormalized form of  $\pi$  is:  $\hat{\pi}(x) = 1$ , for all  $x \in \mathcal{U}$ . The corresponding normalization constant is  $Z_{\hat{\pi}} = |\mathcal{U}|$ .

A Monte Carlo simulator requires two oracle functions, `getWeight $\hat{\pi}$ ( $x$ )` and `getWeight $\hat{p}$ ( $x$ )`, which given an instance  $x \in \mathcal{U}$  return the unnormalized weights  $\hat{\pi}(x)$  and  $\hat{p}(x)$ , respectively.  $\hat{\pi}$  and  $\hat{p}$  are any unnormalized forms of  $\pi$  and  $p$ . Note that the simulator does *not* need to know the corresponding normalization constants  $Z_{\hat{\pi}}$  and  $Z_{\hat{p}}$ .

**4.2. REJECTION SAMPLING.** *Rejection sampling*, due to von Neumann [1963], is the most classical Monte Carlo method. Rejection sampling makes two

```

1: Function RejectionSampling( $C$ )
2:   while (true) do
3:      $X := \text{getSample}_p()$ 
4:     if (accept( $C, X$ ))
5:       return  $X$ 

1: Function accept( $C, x$ )
2:    $r_{\text{RS}}(x) := \frac{\hat{\pi}(x)}{C \hat{p}(x)}$ 
3:   toss a coin whose heads probability is  $r_{\text{RS}}(x)$ 
4:   return true if and only if coin comes up heads

```

FIG. 1. The rejection sampling procedure.

assumptions: (1)  $\text{supp}(\pi) \subseteq \text{supp}(p)$ ; and (2) there is a known *envelope constant*  $C$  satisfying:

$$C \geq \max_{x \in \text{supp}(p)} \frac{\hat{\pi}(x)}{\hat{p}(x)}.$$

The procedure, described in Figure 1, repeatedly generates samples from the trial distribution  $p$ , until a sample is “accepted”. To decide whether a sample  $X$  is accepted, the procedure applies an acceptance–rejection procedure. The procedure accepts the sample  $X$  with the following acceptance probability:

$$r_{\text{RS}}(X) = \frac{\hat{\pi}(X)}{C \hat{p}(X)}.$$

We call  $r_{\text{RS}}$  the *acceptance function*. Note that  $r_{\text{RS}}(x) \in [0, 1]$ , for all  $x \in \text{supp}(p)$ , due to the property of the envelope constant  $C$ .

Intuitively, rejection sampling uses the acceptance–rejection procedure to bridge the gap between  $p$  and  $\pi$ . For example, when  $\pi$  is the uniform distribution and  $p$  is some nonuniform distribution, then the procedure assigns high acceptance probabilities to instances that have low probability in  $p$  and low acceptance probabilities to instances that have high probabilities in  $p$ . Thus, the acceptance–rejection procedure smoothes the distribution  $p$ . A simple analysis shows that for any  $\pi$  and  $p$ , the distribution of the accepted samples is *exactly* the target distribution  $\pi$ .

The expected number of samples from  $p$  needed in order to generate each sample of  $\pi$  is  $CZ_{\hat{p}}/Z_{\hat{\pi}} \geq \max_{x \in \mathcal{U}} \pi(x)/p(x)$ . Hence, the efficiency of the procedure depends on two factors: (1) the similarity between the target distribution and the trial distribution: the more similar they are the smaller is  $\max_{x \in \mathcal{U}} \pi(x)/p(x)$ ; and (2) the gap between the envelope constant  $C$  and  $\max_{x \in \mathcal{U}} \hat{\pi}(x)/\hat{p}(x)$ .

The main drawback of rejection sampling is the need to know the envelope constant  $C$ . A too high value makes the procedure inefficient, while a too low value violates the envelope condition.

**4.3. IMPORTANCE SAMPLING.** *Importance sampling* [Marshall 1956] does not generate samples from the target distribution  $\pi$ , but rather uses samples from  $p$  to directly estimate statistical parameters relative to the distribution  $\pi$ . For simplicity, we assume the desired statistical parameter is  $\mathbb{E}_{\pi}(f(X))$ , where  $X$  is distributed according to  $\pi$ , and  $f$  is some real-valued function, although the technique is applicable to other statistical parameters as well. Unlike rejection sampling, there is no need to know an “envelope constant”.

```

1: Function ImportanceSampling( $f, n$ )
2:   for  $i = 1$  to  $n$  do
3:      $X_i := \text{getSample}_p()$ 
4:      $w(X_i) := \frac{\hat{\pi}(X_i)}{\hat{p}(X_i)}$ 
5:     compute  $f(X_i)$ 

6:   output  $\frac{\frac{1}{n} \sum_{i=1}^n f(X_i) w(X_i)}{\frac{1}{n} \sum_{i=1}^n w(X_i)}$ 

```

FIG. 2. The importance sampling procedure.

In order to estimate  $\mathbb{E}_\pi(f(X))$ , the importance sampling procedure (see Figure 2) generates  $n$  independent samples  $X_1, \dots, X_n$  from the trial distribution  $p$ . If  $p = \pi$ , then clearly  $\frac{1}{n} \sum_{i=1}^n f(X_i)$  is an unbiased estimator of  $\mathbb{E}_\pi(f(X))$ . However, when  $p \neq \pi$ , the samples  $X_1, \dots, X_n$  are “weighted” and the weights have to be accounted for in the estimation. The *importance ratio* at  $x$ , which is defined as

$$w(x) = \frac{\hat{\pi}(x)}{\hat{p}(x)},$$

is exactly the desired weight. Hence,  $\frac{1}{n} \sum_{i=1}^n f(X_i) w(X_i)$  is an unbiased estimator of  $\mathbb{E}_\pi(f(X))$ , modulo normalization. In order to get a correct estimator, we need to estimate also the ratio between the normalization constants of  $\hat{\pi}$  and  $\hat{p}$ . Hence, the final estimator is:

$$\hat{\mu} = \frac{\frac{1}{n} \sum_{i=1}^n f(X_i) w(X_i)}{\frac{1}{n} \sum_{i=1}^n w(X_i)}.$$

*Remark 4.2.* This is one of several possible importance sampling estimators. The estimator is biased (i.e., its expectation is not necessarily  $\mathbb{E}_\pi(f(X))$ ), however it is guaranteed to be close to the true value with high probability, as long as  $n$  is sufficiently large. See more details in Hesterberg [1988].

The efficiency of importance sampling depends on how close is the “shape” of  $\hat{p}(x)$  to the “shape” of  $f(x)\hat{\pi}(x)$ . An appropriately chosen  $p$  can lead to less samples than even sampling from  $\pi$  directly. See more details in Liu [1996]. In this article, we rely on the Liu’s “rule of thumb” [Liu 1996]. It states that the number of samples from  $p$ , required to estimate  $\mathbb{E}_\pi(f(X))$  with the same confidence level (variance) as if using  $n$  independent samples from  $\pi$ , is at most  $n(1 + \text{var}_p(\pi(X)/p(X)))$ .

*Remark 4.3.*  $\text{var}_p(\pi(X)/p(X))$  can be either calculated exactly on local data, or estimated from samples on real search engines (e.g., using again importance sampling).

**4.4. MARKOV CHAIN MONTE CARLO METHODS.** In some situations even generating i.i.d. samples from the trial distribution  $p$  is infeasible. Instead, we are given a *random walk* that converges in the limit to  $p$ . Can we then transform this random walk into a new random walk that converges to the target distribution  $\pi$ ? This is the question addressed by Markov Chain Monte Carlo (MCMC) methods. In this article, we focus on two of these methods: the Metropolis–Hastings (MH) algorithm [Metropolis et al. 1953; Hastings 1970] and the Maximum Degree (MD) method (cf. Bar-Yossef et al. [2000] and Boyd et al. [2004]).

```

1: Function RandomWalk( $P, B, x_0$ )
2:  $X := x_0$ 
3: for  $t = 1$  to  $B$  do
4:    $Y :=$  sample generated according to  $P_X$ 
5:    $X := Y$ 
6: return  $X$ 

```

FIG. 3. A random walk sampler.

A *Markov Chain* (a.k.a. *random walk*) on a finite state space  $\mathcal{U}$  is a stochastic process in which states of  $\mathcal{U}$  are visited successively. The Markov Chain is specified by a  $|\mathcal{U}| \times |\mathcal{U}|$  *probability transition matrix*  $P$ .  $P$  is a *stochastic matrix*, meaning that every row  $x$  of  $P$  specifies a probability distribution  $P_x$  on  $\mathcal{U}$ .  $P$  induces a directed graph  $G_P$  on  $\mathcal{U}$  with nonnegative edge weights. There is an edge  $x \rightarrow y$  in the graph if  $P(x, y) > 0$  and the corresponding weight is  $P(x, y)$ .

The Markov Chain is called *ergodic*, if it satisfies two conditions: (1) it is *irreducible*, meaning that the graph  $G_P$  is strongly connected; and (2) it is *aperiodic*, meaning that the greatest common denominator of the lengths of directed paths connecting any two nodes in  $G_P$  is 1.

A random walk process starts at some initial state  $x_0 \in \mathcal{U}$ . The initial state is chosen in accordance with an initial state distribution  $p_0$  on  $\mathcal{U}$  (typically, the mass of the initial distribution is concentrated on a single fixed state of  $\mathcal{U}$ ). The random walk then successively moves between states of  $\mathcal{U}$ . After visiting state  $x$ , the next state is chosen randomly in accordance with the probability distribution  $P_x$ . This process goes on indefinitely.

Each step  $t$  of the random walk induces a probability distribution  $p_t$  on the state space  $\mathcal{U}$ . The initial distribution is  $p_0$ . Successive distributions are given by the recursive formula:  $p_t = p_{t-1}P$ . Therefore,  $p_t = p_0P^t$ . A fundamental theorem of the theory of Markov Chains states that if a Markov Chain is ergodic, then *regardless* of the initial distribution  $p_0$ , the sequence of distributions  $p_0, p_1, p_2, \dots$  is guaranteed to converge to a unique limit distribution  $p$ . That is,

$$\|p_t - p\| \xrightarrow{t \rightarrow \infty} 0.$$

Furthermore, the unique limit distribution  $p$  is also the unique *stationary distribution* of  $P$ :

$$pP = p.$$

A random walk sampler (see Figure 3) uses a Markov Chain to generate samples from a distribution which is close to  $p$ . The algorithm starts the random walk from any given initial state  $x_0$  and runs it for  $B$  steps ( $B$  is called the “burn-in period”). The reached state  $x_B$  is then returned as the sample. By the above, the distribution of this sample is  $p_B$ , and thus if  $B$  is sufficiently large,  $\|p_B - p\|$  is small. (We discuss below how to choose a sufficiently large burn-in period.) To generate more samples, the algorithm runs the random walk again and again. (There are more efficient sampling procedures, which we mention below.)

MCMC methods allow us to transform a given ergodic Markov Chain  $P$  whose limit distribution is  $p$  into a new Markov Chain  $P_{\text{MCMC}}$  whose limit distribution is  $\pi$ . The two MCMC methods we consider in this article, MH and MD, use the same framework to perform the transformation. The MCMC sampler (see Figure 4)

```

1: Function MCMC( $P, B, x_0$ )
2:  $X := x_0$ 
3: for  $t = 1$  to  $B$  do
4:    $Y :=$  sample generated according to  $P_X$ 
5:   if (accept( $P, X, Y$ ))
6:      $X := Y$ 
7: return  $X$ 

1: Function accept( $P, x, y$ )
2: compute  $r_{\text{MCMC}}(x, y)$  from  $\hat{\pi}(x), \hat{\pi}(y), \hat{p}(x), \hat{p}(y), P(x, y)$ , and  $P(y, x)$ .
3: toss a coin whose probability is  $r_{\text{MCMC}}(x, y)$ 
4: return true if and only if coin comes up heads

```

FIG. 4. An MCMC sampler.

runs a random walk similarly to the random walk sampler, except that it applies an acceptance–rejection procedure at each step. The procedure is used to determine whether the “proposed state”  $Y$  is “accepted” and thus will become the next step of the random walk or not. The acceptance function  $r_{\text{MCMC}}(x, y)$  depends on both the current state and the proposed state. MH and MD differ in the choice of the acceptance function.

The acceptance–rejection procedure effectively modifies the transition probabilities of the random walk and defines a new transition matrix  $P_{\text{MCMC}}$ . A careful choice of the acceptance function  $r_{\text{MCMC}}$  guarantees that the limit distribution of  $P_{\text{MCMC}}$  is the target distribution  $\pi$ .

*Remark 4.4.* Throughout this article, when dealing with MCMC samplers, we will assume that the target distribution  $\pi$  satisfies  $\text{supp}(\pi) = \mathcal{U}$ , that is, every node in the Markov Chain’s state space has positive probability under the target distribution. This assumption is made only to make the presentation simpler.

**4.4.1. The Metropolis–Hastings Algorithm.** The MH algorithm requires that the initial Markov Chain  $P$  is not only ergodic but also satisfies the following condition: for all states  $x, y \in \mathcal{U}$ ,  $P(x, y) > 0 \Leftrightarrow P(y, x) > 0$ . Also the supports of the limit distribution  $p$  and of the target distribution  $\pi$  must be equal (i.e.,  $\text{supp}(p) = \text{supp}(\pi)$ ).

The acceptance function used by the MH algorithm is the following:

$$r_{\text{MH}}(x, y) = \min \left\{ \frac{\pi(y) P(y, x)}{\pi(x) P(x, y)}, 1 \right\}.$$

Note that since  $\frac{\hat{\pi}(y)}{\hat{\pi}(x)} = \frac{\pi(y)}{\pi(x)}$ , oracle access to an unnormalized form of  $\pi$  suffices for computing this acceptance function.

The probability transition matrix of the resulting Markov Chain is:

$$P_{\text{MH}}(x, y) = \begin{cases} P(x, y) r_{\text{MH}}(x, y), & \text{if } x \neq y, \\ P(x, x) r_{\text{MH}}(x, x) + 1 - \sum_{z \in \mathcal{U}} P(x, z) r_{\text{MH}}(x, z), & \text{if } x = y. \end{cases}$$

It can be shown (see, e.g., Diaconis and Saloff-Coste [1998]) that the limit distribution of this Markov Chain is exactly  $\pi$ .

**4.4.2. The Maximum Degree Method.** The MD method applies to arbitrary ergodic Markov Chains. The supports of the limit distribution  $p$  and of the target distribution  $\pi$  should satisfy  $\text{supp}(p) = \text{supp}(\pi)$ . Similarly to rejection sampling, application of this method requires the availability of an “envelope constant”  $C$ .  $C$  must satisfy the following condition:

$$C \geq \max_{x \in \mathcal{U}} \frac{\hat{p}(x)}{\hat{\pi}(x)}.$$

The need for an envelope constant is the major disadvantage of the Maximum Degree method relative to the Metropolis–Hastings algorithm. However, if the envelope constant is chosen sufficiently close to its lower bound, then the MD method can become significantly more efficient than the MH algorithm.

The acceptance function used by the MD algorithm is the following:

$$r_{\text{MD}}(x) = \frac{\hat{p}(x)}{C \hat{\pi}(x)}.$$

*Remark 4.5.* Notice the reversed roles of  $\hat{p}(x)$  and  $\hat{\pi}(x)$ , comparing to the rejection sampling acceptance function. This is not accidental. Rejection sampling is similar (but not identical) to the application of the MD method on a degenerate random walk, which converges in one step to the limit distribution  $p$  (that is, all the rows of the transition matrix  $P$  equal  $p$ ). The reversed roles are due to the reversed semantics of acceptance in rejection sampling versus MCMC methods. In rejection sampling, when the current state is accepted, the process stops and outputs the current state. In MCMC methods, when a *proposed* state is accepted, then implicitly the current state is *rejected*, and the process continues.

The probability transition matrix of the MD Markov Chain is:

$$P_{\text{MD}}(x, y) = \begin{cases} P(x, y) r_{\text{MD}}(x), & \text{if } x \neq y, \\ P(x, x) r_{\text{MD}}(x) + 1 - r_{\text{MD}}(x), & \text{if } x = y. \end{cases}$$

**THEOREM 4.6.** *The limit distribution of the Markov Chain defined by  $P_{\text{MD}}$  is  $\pi$ .*

*Remark 4.7.* To the best of our knowledge, the formulation above is the first application of the Maximum Degree method to arbitrary ergodic Markov Chains and to arbitrary target distributions. Previous studies [Bar-Yossef et al. 2000; Boyd et al. 2004] applied the MD method only in the special case  $P$  is a simple random walk on an undirected graph  $G$  and  $\pi$  is the uniform distribution over the vertices of  $G$ . The limit distribution of the simple random walk is the degree distribution (i.e., each node is chosen proportionally to its degree). In this case,  $C$  must be set as an upper bound on the maximum degree of the graph, and that is why the method is called “Maximum Degree”.

The acceptance function of the MD method has the peculiar property that it depends only on the current state  $x$  and not on the proposed state  $y$ . This fact allows a more efficient implementation of the MD sampler (see Figure 5). This sampler postpones the selection of the proposed state  $Y$  to until after acceptance is achieved. Hence, rather than selecting a proposed state every time the acceptance–rejection procedure is called, the proposed state is selected only once. Since in many situations selection of a proposed state is costly, this amounts to significant savings in running time.



```

1: Function MD( $P, B, x_0, C$ )
2:  $X := x_0$ 
3: for  $t = 1$  to  $B$  do
4:   if ( $\text{accept}(P, C, X)$ )
5:      $Y :=$  sample generated according to  $P_X$ 
6:      $X := Y$ 
7: return  $X$ 

1: Function accept( $P, C, x$ )
2:  $r_{MD}(x) := \frac{\hat{p}(x)}{C \cdot \hat{\pi}(x)}$ 
3: toss a coin whose heads probability is  $r_{MD}(x)$ 
4: return true if and only if coin comes up heads

```

FIG. 5. The Maximum Degree sampler.

A further improvement is possible. Note that the number of steps the random walk spends at each state  $x$  is a geometric random variable with a known success probability. Therefore, when the sampler moves to a new state  $x$ , it can *simulate* the number of steps that it is going spend at the state by generating an appropriate geometric random variable. It can then immediately select the next state. This saves the need to perform the iterative coin tosses.

**4.4.3. The Burn-in Period.** How should we set the burn-in period  $B$  of a random walk in order to guarantee that the selected sample has distribution which is close to the limit distribution? To this end, a rich pool of techniques, ranging from algebraic to geometric, is available from the theory of Markov Chains (see a survey by Diaconis and Saloff-Coste [1998] for a detailed review). In this article we focus on a popular technique, which is based on estimating the *spectral gap* of the Markov Chain's transition matrix.

Let  $P$  be the transition matrix of an ergodic Markov Chain, whose limit distribution is  $p$ . For each  $\varepsilon > 0$ , we define the  $\varepsilon$ -burn-in period (a.k.a.  $\varepsilon$ -mixing time) of the Markov Chain as:

$$T_\varepsilon(P) = \min\{t \mid \text{for all initial distributions } p_0 \text{ and } \forall t' \geq t, \|p_{t'} - p\| < \varepsilon\}.$$

That is,  $T_\varepsilon(P)$  is the first step  $t$ , for which  $p_t$  is guaranteed to be at most  $\varepsilon$  away from the limit distribution  $p$ .

The spectral gap technique for bounding the burn-in period is applicable only to *reversible* Markov Chains. A Markov Chain is called *reversible*, if for all states  $x, y \in \mathcal{U}$ ,  $p(x)P(x, y) = p(y)P(y, x)$ . It can be shown that a reversible Markov Chain is equivalent to a random walk on a weighted and undirected graph.

Let  $n = |\mathcal{U}|$  and let  $\lambda_1, \dots, \lambda_n$  be the eigenvalues of  $P$ , ordered from largest to smallest by absolute value (i.e.,  $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|$ ). The *spectral gap* of  $P$  is defined as the difference between its first and the second eigenvalues:

$$\alpha(P) = |\lambda_1| - |\lambda_2| = 1 - |\lambda_2|.$$

(For a transition matrix  $P$ ,  $\lambda_1 = 1$  always.) Using the spectral gap, one can bound the pointwise distance between  $p_t$  and  $p$ :

**THEOREM 4.8** ([SINCLAIR 1993], PROPOSITION 2.1, PAGE 47). *Let  $P$  be the transition matrix of a reversible Markov Chain whose limit distribution is  $p$ . Then,*

for any initial distribution, for every  $t \geq 0$ , and for every  $x \in \mathcal{U}$ ,

$$|p_t(x) - p(x)| \leq \sqrt{\frac{p(x)}{p_{\min}}} \cdot (1 - \alpha(P))^t,$$

where  $p_{\min} = \min_{x \in \mathcal{U}} p(x)$ .

That is, the larger the spectral gap, the faster the convergence to the limit distribution.

An immediate corollary of the above theorem is the following bound on the burn-in period of a reversible Markov Chain:

**COROLLARY 4.9.** *Let  $P$  be the transition matrix of a reversible Markov Chain whose limit distribution is  $p$ . Then, for every  $\varepsilon > 0$ ,*

$$T_\varepsilon(P) \leq \frac{1}{\alpha(P)} \left( \ln \frac{1}{p_{\min}} + \ln \frac{1}{\varepsilon} \right).$$

**PROOF.** We first bound the total variation distance between  $p$  and  $p_t$ :

$$\begin{aligned} \|p - p_t\| &= \frac{1}{2} \sum_{x \in \mathcal{U}} |p(x) - p_t(x)| \\ &\leq \frac{1}{2} \sum_{x \in \mathcal{U}} \sqrt{\frac{p(x)}{p_{\min}}} \cdot (1 - \alpha(P))^t \quad (\text{By Theorem 4.8}) \\ &\leq \frac{(1 - \alpha(P))^t}{2 \cdot \sqrt{p_{\min}}} \cdot \sqrt{|\mathcal{U}|} \quad (\text{Cauchy-Schwartz}) \\ &\leq \frac{(1 - \alpha(P))^t}{2 \cdot p_{\min}} \quad (\text{Using the fact } p_{\min} \leq 1/|\mathcal{U}|). \end{aligned}$$

Hence, in order for  $\|p - p_t\| \leq \varepsilon$ , we need

$$t \geq \frac{\ln \frac{1}{\varepsilon} + \ln \frac{1}{p_{\min}} - \ln 2}{\ln \frac{1}{1 - \alpha(P)}}.$$

The corollary follows using the fact  $1 - \alpha < \exp(-\alpha)$  for  $0 < \alpha < 1$ .  $\square$

It can be shown that if  $P$  is reversible, then also  $P_{\text{MH}}$  and  $P_{\text{MD}}$  are reversible, and thus we can use the spectral gap technique to estimate the required burn-in periods of the MH and the MD samplers.

**4.4.4. Efficient Random Walk Sampling.** The naive method for generating multiple independent samples from (a distribution that is close to) the limit distribution  $p$  of a Markov Chain is to run a new random walk for each desired sample. This incurs high overhead, if the required number of samples is large. It turns out that in some situations it is possible to use a *single* random walk to generate multiple samples.

Aldous [1987] (with further improvements by Gillman [1998] and Kahale [1997]) proposed an efficient procedure for generating dependent samples that can be used to perform accurate *density estimations*.

Suppose  $\mathcal{A} \subseteq \mathcal{U}$  is a subset of the state space. The *density* of  $\mathcal{A}$  under a probability measure  $p$  is the quantity  $p(\mathcal{A})$ . For example, when  $p$  is the uniform distribution on  $\mathcal{U}$ , the density of  $\mathcal{A}$  is the ratio  $|\mathcal{A}|/|\mathcal{U}|$ . Suppose that we are given an “oracle” procedure, which on input  $x \in \mathcal{U}$  can determine whether  $x \in \mathcal{A}$  or not, and that we would like to use this procedure to estimate the density of  $\mathcal{A}$ . This type of estimation is very popular. One example from our domain is the estimation of the relative overlap between two search engines.

Given a Markov Chain  $P$  whose limit distribution is  $p$ , the most obvious method for estimating  $p(\mathcal{A})$  would be to run  $n$  random walks, each for  $T_\varepsilon(P)$  steps, and thereby obtain  $n$  i.i.d. samples from (a distribution that is close to)  $p$ . The estimator for  $p(\mathcal{A})$  would then be the fraction of the  $n$  samples that fall into  $\mathcal{A}$ .

Aldous's procedure is more efficient. Instead of running  $n$  walks, Aldous suggests running only a single walk of length  $T_\varepsilon(P) + n \frac{1}{\alpha(P)}$ , and use the last  $n \frac{1}{\alpha(P)}$  states visited as the samples, disregarding the first  $T_\varepsilon(P)$  steps as *sample delay*. As shown in Aldous [1987], Gillman [1998], and Kahale [1997], these  $n \frac{1}{\alpha(P)}$  dependent samples can be used to produce an estimate for  $p(\mathcal{A})$ , which is as good as the estimate obtained from the  $n$  independent samples. Overall, the Aldous procedure saves a factor of  $\log(1/p_{\min})$  in the number of random walk steps over the naive procedure.

## 5. Approximate Monte Carlo Methods

All Monte Carlo methods assume that the trial distribution  $p$  is known, up to normalization. This assumption turns out to be very problematic in our setting, since the trial distributions we construct depend on unknown internal data of the search engine. An *approximate Monte Carlo method* employs an “approximate trial distribution”  $q$  rather than the true trial distribution  $p$  in the acceptance function calculations. The mismatch between the trial samples (that are generated from  $p$ ) and the acceptance function (that is based on  $q$ ) implies that the sampling distributions of the resulting procedures are no longer guaranteed to equal the target distribution  $\pi$ . To the best of our knowledge, no previous study addressed this scenario before.

We show that the sampling distribution of approximate rejection sampling and the limit distributions of approximate Metropolis–Hastings and of approximate Maximum Degree are all identical to some distribution  $\pi'$ , for which we give a closed form formula. We then prove that  $\pi'$  is close to the target distribution  $\pi$ , as long as the trial distribution  $p$  and the approximate trial distribution  $q$  are similar. We also prove that the estimations generated by approximate importance sampling are close to the true values. All proofs are provided in Appendix B.

**5.1. APPROXIMATE REJECTION SAMPLING.** Consider the following modified (“approximate”) form of rejection sampling. The procedure is given the following three oracle procedures: (1)  $\text{getSample}_p()$ , which generates samples from  $p$ , (2)  $\text{getWeight}_{\hat{\pi}}(x)$ , which calculates an unnormalized form of the target distribution  $\pi$ ; and (3)  $\text{getWeight}_{\hat{q}}(x)$ , which calculates an unnormalized form of an “approximate trial distribution”  $q$ .  $\pi$ ,  $p$ ,  $q$  are assumed to satisfy:

$$\text{supp}(\pi) \subseteq \text{supp}(p) \subseteq \text{supp}(q).$$

The approximate rejection sampling procedure works exactly like the standard procedure, except that the acceptance function it uses is the following:

$$r'_{\text{RS}}(x) = \frac{\hat{\pi}(x)}{C \hat{q}(x)},$$

where  $C \geq \max_{x \in \text{supp}(p)} \frac{\hat{\pi}(x)}{\hat{q}(x)}$  is an envelope constant.

The following theorem characterizes the sampling distribution of the approximate rejection sampling procedure and analyzes its sample complexity:

**THEOREM 5.1.** *The sampling distribution of the approximate rejection sampling procedure is:*

$$\pi'(x) = \pi(x) \frac{\frac{p(x)}{q(x)}}{\mathbb{E}_{\pi} \left( \frac{p(X)}{q(X)} \right)}.$$

The expected number of samples from  $p$  needed to generate each sample from  $\pi'$  is:

$$\frac{C Z_{\hat{q}}}{Z_{\hat{\pi}} \mathbb{E}_{\pi} \left( \frac{p(X)}{q(X)} \right)}.$$

The following proposition shows that as long as the trial distribution  $p$  and the approximate trial distribution  $q$  are “similar” relative to the target distribution  $\pi$  (in the sense that the ratio  $p(X)/q(X)$  has low variance when  $X$  is chosen according to  $\pi$ ), then the sampling distribution of approximate rejection sampling is close to the target distribution:

**PROPOSITION 5.2**

$$\|\pi' - \pi\| = \frac{1}{2} \text{ndev}_{\pi} \left( \frac{p(X)}{q(X)} \right).$$

**5.2. APPROXIMATE IMPORTANCE SAMPLING.** The approximate importance sampling procedure assumes oracle access to the approximate trial distribution  $q$ .  $\pi$ ,  $p$ ,  $q$  are assumed to satisfy:

$$\text{supp}(\pi) \subseteq \text{supp}(p) \subseteq \text{supp}(q).$$

The approximate importance sampling procedure works exactly like the standard importance sampling procedure, except that the following approximate importance ratios are used:

$$w'(x) = \frac{\hat{\pi}(x)}{\hat{q}(x)}.$$

The following theorem shows that as long as the ratio  $p(X)/q(X)$  is uncorrelated with the function  $f(X)$  whose expectation we need to estimate, then the estimate produced by approximate importance sampling is close to the desired parameter:

**THEOREM 5.3.** *Let*

$$\hat{\mu}' = \frac{A}{B} = \frac{\frac{1}{n} \sum_{i=1}^n f(X_i) w'(X_i)}{\frac{1}{n} \sum_{i=1}^n w'(X_i)}$$

be the estimator produced by the approximate importance sampling procedure for the parameter  $\mathbb{E}_\pi(f(X))$ . Then,

$$\frac{\mathbb{E}_p(A)}{\mathbb{E}_p(B)} = \mathbb{E}_\pi(f(X)) + \frac{\text{cov}_\pi\left(f(X), \frac{p(X)}{q(X)}\right)}{\mathbb{E}_\pi\left(\frac{p(X)}{q(X)}\right)}.$$

The next theorem shows that as  $n$  grows to infinity, the difference between  $\mathbb{E}_p(\hat{\mu}')$  and  $\frac{\mathbb{E}_p(A)}{\mathbb{E}_p(B)}$  diminishes to 0.

**THEOREM 5.4.** *Suppose  $A$  and  $B$  are two estimators such that  $\frac{\mathbb{E}(A)}{\mathbb{E}(B)} = I$ . Let  $A_1, \dots, A_n$  and  $B_1, \dots, B_n$  be  $n$  independent instances of  $A$  and  $B$ , respectively. Then,*

$$\left| \mathbb{E} \left( \frac{\frac{1}{n} \sum_{i=1}^n A_i}{\frac{1}{n} \sum_{i=1}^n B_i} \right) - I \right| \leq \frac{1}{n} \cdot C,$$

where  $C = I \cdot \frac{\text{var}(B)}{\mathbb{E}^2(B)} + \frac{|\text{cov}(A, B)|}{\mathbb{E}^2(B)} + o(1)$ .

The above theorem can be proved using the Delta method in statistics (see, e.g., Stuart and Ord [1994, p. 350]).

**5.3. APPROXIMATE METROPOLIS–HASTINGS.** Next, we discuss an approximate variant of the Metropolis–Hastings algorithm. Like in approximate rejection sampling, we assume oracle access to an approximate trial distribution  $q$ . We assume that  $\text{supp}(q) = \text{supp}(p) = \text{supp}(\pi) = \mathcal{U}$ .

We also need to assume that the base Markov Chain  $P$  is reversible (i.e.,  $p(x)P(x, y) = p(y)P(y, x)$ , for all  $x, y \in \mathcal{U}$ ). When  $P$  is reversible, the acceptance function of the standard Metropolis–Hastings algorithm can be rewritten as follows:

$$r_{\text{MH}}(x, y) = \min \left\{ \frac{\pi(y) P(y, x)}{\pi(x) P(x, y)}, 1 \right\} = \min \left\{ \frac{\pi(y) p(x)}{\pi(x) p(y)}, 1 \right\}.$$

The approximate Metropolis–Hastings procedure is identical to the standard procedure, except that it uses the following acceptance function:

$$r'_{\text{MH}}(x, y) = \min \left\{ \frac{\pi(y) q(x)}{\pi(x) q(y)}, 1 \right\}.$$

Note that since  $\frac{\hat{\pi}(y)}{\hat{\pi}(x)} = \frac{\pi(y)}{\pi(x)}$  and  $\frac{\hat{q}(x)}{\hat{q}(y)} = \frac{q(x)}{q(y)}$ , this acceptance function can be computed using oracle access to unnormalized forms of  $\pi$  and  $q$ .

The following theorem shows that the resulting Markov Chain is ergodic and that its unique limit distribution is  $\pi'$ , where  $\pi'$  is defined as in Theorem 5.1:  $\pi'(x) = \pi(x) \cdot \frac{p(x)}{q(x)} / \mathbb{E}_\pi\left(\frac{p(X)}{q(X)}\right)$ .

**THEOREM 5.5.** *Let  $P'_{\text{MH}}$  be the transition matrix of the approximate Metropolis–Hastings algorithm. Then,  $P'_{\text{MH}}$  forms an ergodic Markov Chain and its unique limit distribution is  $\pi'$ .*

It follows from Proposition 5.2 that the limit distribution of the approximate MH random walk is close to the target distribution  $\pi$  as long as  $p$  and  $q$  are similar relative to  $\pi$ .

**5.4. APPROXIMATE MAXIMUM DEGREE.** As before, we assume oracle access to an approximate trial distribution  $q$  and that  $\text{supp}(q) = \text{supp}(p) = \text{supp}(\pi) = \mathcal{U}$ . This time, we do not need to assume that the base Markov Chain  $P$  is reversible.

Approximate MD is identical to the standard MD, except that the following modified acceptance function is used:

$$r'_{\text{MD}}(x) = \frac{\hat{q}(x)}{C \hat{\pi}(x)}, \quad \text{where } C \geq \max_{x \in \mathcal{U}} \frac{\hat{q}(x)}{\hat{\pi}(x)}.$$

The following theorem shows that the resulting Markov Chain is ergodic and that its unique limit distribution equals  $\pi'$ , where  $\pi'$  is defined as in Theorem 5.1:  $\pi'(x) = \pi(x) \cdot \frac{p(x)}{q(x)} / \mathbb{E}_{\pi}(\frac{p(X)}{q(X)})$ .

**THEOREM 5.6.** *Let  $P'_{\text{MD}}$  be the transition matrix of the approximate Maximum Degree procedure. Then,  $P'_{\text{MD}}$  forms an ergodic Markov Chain. The unique limit distribution of  $P'_{\text{MD}}$  is  $\pi'$ .*

It follows from Proposition 5.2 that the limit distribution of the approximate MD random walk is close to the target distribution  $\pi$  as long as  $p$  and  $q$  are close relative to  $\pi$ .

## 6. Formal Framework

In this section we lay out the formal framework for the design and analysis of search engine samplers.

### 6.1. SEARCH ENGINES

**Definition 6.1.** (*Search engine*). A *Search Engine* is a 4-tuple  $\langle \mathcal{D}, \mathcal{Q}, \text{results}(\cdot), k \rangle$ , where:

- (1)  $\mathcal{D}$  is the document corpus indexed. Documents are assumed to have been pre-processed (e.g., they may be truncated to some maximum size limit).
- (2)  $\mathcal{Q}$  is the space of queries supported by the search engine.
- (3)  $\text{results}(\cdot)$  is a mapping that maps every query  $q \in \mathcal{Q}$  to an ordered sequence of documents, called *results*. The *cardinality* of  $q$  is the number of results:  $\text{card}(q) = |\text{results}(q)|$ .
- (4)  $k$  is the *result limit*. Only the top  $k$  results are actually returned to the user via the search engine's public interface.

A query  $q$  is said to be *overflowing*, if  $\text{card}(q) > k$ . It is said to be *underflowing*, if  $\text{card}(q) = 0$ . If  $q$  neither overflows nor underflows, it is called *valid*.

A document  $x$  *matches* a query  $q$ , if  $x \in \text{results}(q)$ . The set of queries that a document  $x$  matches is denoted  $\text{queries}(x)$ .

### 6.2. SEARCH ENGINE SAMPLERS

**Definition 6.2** (*Search Engine Sampler*). Let  $\pi$  be a target distribution on a document corpus  $\mathcal{D}$  indexed by a search engine. A *search engine sampler* is a



randomized procedure, which is given access to three “oracle” procedures:

- (1) `getWeight $\hat{\pi}$ (x)`: Returns the unnormalized weight of an instance  $x \in \mathcal{D}$  under the target distribution  $\pi$ .
- (2) `getResults(q)`: Returns the top  $k$  results from the search engine on the query  $q$ .
- (3) `getDocument(x)`: Returns the HTTP header and the content of the document  $x$ .

Each invocation of the sampler outputs a random document  $X$  from the corpus  $\mathcal{D}$ . The distribution of the sample  $X$  is called the *sampling distribution* and is denoted by  $\eta$ . Successive invocations of the sampler produce independent samples from  $\eta$ .

If the unnormalized form of  $\pi$  is independent of the corpus  $\mathcal{D}$  (e.g., when  $\pi$  is the uniform distribution and  $\hat{\pi}(x) = 1$  for all  $x$ ), then the same sampler can be used to sample from different search engines. All that needs to be changed is the implementation of the procedure `getResults(q)`.

When the sampling distribution  $\eta$  of the sampler equals exactly the target distribution  $\pi$ , then the sampler is said to be *perfect*. Otherwise, it is *biased*. We discuss below the two main metrics for measuring the quality of a sampler with respect to a given target distribution: the *sampling recall* and the *sampling bias*.

Not all documents in  $\mathcal{D}$  are practically reachable via the public interface of the search engine. Some pages have no text content and others have very low static rank, and thus formulating a query that returns them as one of the top  $k$  results may be impossible. Thus, search engine samplers usually generate samples only from large subsets of  $\mathcal{D}$  and not from the whole corpus  $\mathcal{D}$ . The *sampling recall* of a sampler with target  $\pi$  and sampling distribution  $\eta$  is defined as  $\pi(\text{supp}(\eta))$ . For instance, when  $\pi$  is the uniform distribution, the sampling recall is  $|\text{supp}(\eta)|/|\mathcal{D}|$ , that is, the fraction of documents that the sampler can actually return as samples. Ideally, we would like the recall to be as close to 1 as possible. Note that even if the recall is lower than 1, but  $\text{supp}(\eta)$  is sufficiently representative of  $\mathcal{D}$ , then estimators that use samples from  $\text{supp}(\eta)$  can produce accurate estimates of parameters of the whole corpus  $\mathcal{D}$ .

Since samplers sample only from large subsets of  $\mathcal{D}$  and not from  $\mathcal{D}$  in its entirety, it is unfair to measure the bias of a sampler directly with respect to the target distribution  $\pi$ . Rather, we measure the bias with respect to the distribution  $\pi$  restricted to  $\text{supp}(\eta)$ . Formally, let  $\pi_{\text{supp}(\eta)}$  be the following distribution on  $\text{supp}(\eta)$ :

$$\pi_{\text{supp}(\eta)}(x) = \frac{\pi(x)}{\pi(\text{supp}(\eta))}, \quad \forall x \in \text{supp}(\eta).$$

The *sampling bias* of the sampler is defined as the total variation distance between  $\eta$  and  $\pi_{\text{supp}(\eta)}$ :

$$\|\eta - \pi_{\text{supp}(\eta)}\| = \frac{1}{2} \sum_{x \in \text{supp}(\eta)} \left| \eta(x) - \frac{\pi(x)}{\pi(\text{supp}(\eta))} \right|.$$

For example, if a sampler generates truly uniform samples from a subset  $\mathcal{D}'$  of  $\mathcal{D}$  that constitutes 80% of  $\mathcal{D}$ , then its sampling recall is 0.8 and its sampling bias is 0.

The two most expensive resources of a search engine sampler are: (1) the amount of queries submitted to the search engine; and (2) the amount of additional web pages fetched. Search engine queries and web page fetches consume significant amount of time and require network bandwidth. In addition, the rate at which a

sampler can submit queries to the search engine is usually very restricted, since search engines impose hard daily limits on the number of queries they accept from any single user. We thus measure the complexity of search engine samplers in terms of their *query cost* (expected number of calls to the subroutine `getResults()` per sample generated) and their *fetch cost* (expected number of calls to the subroutine `getDocument()` per sample generated).

**6.3. QUERY POOLS.** Consider a search engine  $\mathcal{S}$ , whose corpus is  $\mathcal{D}$  and whose query space is  $\mathcal{Q}$ .

*Definition 6.3 (Query Pool).* A query pool is a fragment  $\mathcal{P} \subseteq \mathcal{Q}$  of the query space.

A query pool may be specified either *explicitly* as a set of queries, for example,  
 $\mathcal{P} = \{[\text{Java software}], [\text{“Michael Jordan” -basketball}], [\text{Car OR Automobile}]\},$   
 or *implicitly*, for example,

$$\mathcal{P} = \text{All single-term queries.}$$

Note that in the latter case in order to transform  $\mathcal{P}$  into an explicit form, we need a list of all the terms that occur in the corpus  $\mathcal{D}$ . All the samplers we consider in this article fix some query pool  $\mathcal{P}$  and use only queries from  $\mathcal{P}$  in order to generate the sample documents from  $\mathcal{D}$ .

**6.3.1. Queries-Documents Graph.** Every query pool  $\mathcal{P}$  naturally induces a bipartite graph  $B_{\mathcal{P}}$  on  $\mathcal{P} \times \mathcal{D}$ . Its left side consists of all queries in  $\mathcal{P}$  and its right side consists of all documents in  $\mathcal{D}$ . A query  $q \in \mathcal{P}$  is connected to a document  $x \in \mathcal{D}$  if and only if  $x \in \text{results}(q)$ .<sup>5</sup>

The *cardinality* of a query  $q$ , denoted  $\text{card}(q)$ , is its degree in  $B_{\mathcal{P}}$ . This is exactly the number of documents in the result set of the query. The cardinality of a set of queries  $\mathcal{P}' \subseteq \mathcal{P}$  is defined as:

$$\text{card}(\mathcal{P}') = \sum_{q \in \mathcal{P}'} \text{card}(q).$$

For a document  $x$ , we denote by  $\text{queries}_{\mathcal{P}}(x)$  the set of its neighbors in  $B_{\mathcal{P}}$ . These are exactly all the queries in  $\mathcal{P}$  that  $x$  matches. For example, if  $\mathcal{P}$  is the pool of all single term queries, then  $\text{queries}_{\mathcal{P}}(x)$  is the set of all distinct terms that occur in the text of  $x$ . The *degree* of  $x$  is:  $\text{deg}_{\mathcal{P}}(x) = |\text{queries}_{\mathcal{P}}(x)|$ . The degree of a set of documents  $\mathcal{D}' \subseteq \mathcal{D}$  is defined as:

$$\text{deg}_{\mathcal{P}}(\mathcal{D}') = \sum_{x \in \mathcal{D}'} \text{deg}_{\mathcal{P}}(x).$$

Note that  $\text{card}(\mathcal{P})$  is the sum of the degrees of all nodes on the left side of  $B_{\mathcal{P}}$ , while  $\text{deg}_{\mathcal{P}}(\mathcal{D})$  is the sum of the degrees of all nodes on the right side of  $B_{\mathcal{P}}$ . Both sums equal to the number of edges in  $B_{\mathcal{P}}$ , and we thus obtain the following result:

**PROPOSITION 6.4.** *Let  $\mathcal{P}$  be any query pool. Then,  $\text{card}(\mathcal{P}) = \text{deg}_{\mathcal{P}}(\mathcal{D})$ .*

<sup>5</sup> A similar graph was suggested by Davison [2004] in a different context.

**6.3.2. Recall.** We say that a query pool  $\mathcal{P}$  *covers* a document  $x$ , if  $\deg_{\mathcal{P}}(x) > 0$ . That is, at least one query in  $\mathcal{P}$  returns  $x$  as a result. Let  $\mathcal{D}_{\mathcal{P}}$  be the collection of documents covered by  $\mathcal{P}$ . Note that a sampler that uses only queries from  $\mathcal{P}$  can never reach documents outside  $\mathcal{D}_{\mathcal{P}}$ .

For a distribution  $\pi$  on  $\mathcal{D}$ , the *recall of  $\mathcal{P}$  with respect to  $\pi$* , denoted  $\text{recall}_{\pi}(\mathcal{P})$ , is the probability that a random document selected from  $\pi$  is covered by  $\mathcal{P}$ . That is,

$$\text{recall}_{\pi}(\mathcal{P}) = \pi(\mathcal{D}_{\mathcal{P}}).$$

In the case  $\pi$  is the uniform distribution on  $\mathcal{D}$ ,  $\text{recall}_{\pi}(\mathcal{P}) = |\mathcal{D}_{\mathcal{P}}|/|\mathcal{D}|$ .

**6.3.3. Overflow Probability.** Recall that a query  $q$  is *valid* if it neither overflows nor underflows. The set of valid queries  $q \in \mathcal{P}$  is denoted  $\mathcal{P}_{+}$  and the set of invalid queries is denoted  $\mathcal{P}_{-}$ . Given a distribution  $\phi$  on  $\mathcal{P}$ , we define the *overflow probability* of  $\phi$ , denoted  $\text{ovprob}(\phi)$ , to be the probability that a random query  $Q$  chosen from  $\phi$  overflows:

$$\text{ovprob}(\phi) = \Pr_{\phi}(\text{card}(Q) > k).$$

**6.4. INCIDENCE COMPUTATION.** The samplers we use in this article require “local accessibility” to the queries-documents graph  $B_{\mathcal{P}}$ . By that we mean that the sampler needs efficient implementations of the following procedures that compute incidences in the graph:

- (1) `getIncidentDocs $_{\mathcal{P}}$ (q)`: Given a query  $q \in \mathcal{P}$ , returns all documents that are incident to  $q$  in  $B_{\mathcal{P}}$ , that is,  $\text{results}(q)$ .
- (2) `getIncidentQueries $_{\mathcal{P}}$ (x)`: Given a document  $x \in \mathcal{D}$ , returns all queries that are incident to  $x$  in  $B_{\mathcal{P}}$ , that is,  $\text{queries}_{\mathcal{P}}(x)$ .

Implementing the above procedures efficiently is crucial since our algorithms call them over and over again. Moreover, the query and the fetch costs of our samplers are “consumed” only through these procedures. We first describe a simple, yet impractical, implementations of the above procedures, and then propose more realistic implementations.

**6.4.1. Naive Implementations.** The implementation of `getIncidentDocs $_{\mathcal{P}}$ (q)` is trivial: just submit  $q$  to the search engine and output all the results returned. The cost of this implementation is a single search engine query. It has one caveat, though: it is applicable only to non-overflowing queries. If the given query overflows, the procedure returns only the top  $k$  documents in the result set. The second procedure (`getIncidentQueries $_{\mathcal{P}}$ (x)`) is implemented by submitting each query from  $\mathcal{P}$  to the search engine and returning those that have  $x$  as their result. This implementation is usually impractical, due to a large number of queries in a typical query pool.

**6.4.2. Efficient Implementations.** The implementation of `getIncidentDocs $_{\mathcal{P}}$ (q)` is as before: submit  $q$  to the search engine and output all the results returned. To implement `getIncidentQueries $_{\mathcal{P}}$ (x)`, we first fetch the content of  $x$ . Then, for a pool consisting of term/phrase queries, we can compute  $\text{queries}_{\mathcal{P}}(x)$  by extracting all the terms/phrases *directly* from the content of  $x$ , and without submitting queries to the search engine. We call the set of queries extracted from the content of  $x$  the *predicted queries*, and denote them by  $\text{pqueries}_{\mathcal{P}}(x)$ . The cost of computing  $\text{pqueries}_{\mathcal{P}}(x)$  by the above procedure is a single page fetch. We note that not

all the pools allow such an implementation. Pools that consist solely of standard term/phrase queries (e.g., [java], ["Michael Jordan"]) do. Pools that contain other types of queries, like complex Boolean queries or link queries (which ask for documents that contain a link to a given URL), may not allow such an implementation.

**6.4.3. Inaccuracies in Incidence Calculation.** Using  $pqueries_{\mathcal{P}}(x)$  as a substitute for  $queries_{\mathcal{P}}(x)$  introduces some inaccuracies in the incidence calculations. We distinguish between two types of inaccuracies: (1) *incidence recall deficiency*:  $pqueries_{\mathcal{P}}(x)$  may miss some queries that belong to  $queries_{\mathcal{P}}(x)$ ; (2) *incidence precision deficiency*:  $pqueries_{\mathcal{P}}(x)$  may contain queries that do not belong to  $queries_{\mathcal{P}}(x)$ .

In the following we list several factors that cause these deficiencies:

(1) *Overflowing Queries*. Some of the queries that  $x$  matches may be overflowing (have more than  $k$  matches) and consequently may not return  $x$  as one of their top  $k$  results.

(2) *Duplicates and Near-Duplicates*. If the search engine filters duplicate or near-duplicate documents, a query that  $x$  matches may not return  $x$  as a result, if one of the documents that are similar to  $x$  is returned as a result. Note that although we requested search engines not to eliminate duplicates from search results, duplicate elimination may be done at crawl/index time already.

(3) *Host Collapsing*. If the search engine collapses documents belonging to the same host, a query that matches  $x$  may not return  $x$  as a result, if another document from the same host is returned as a result. Note that we requested search engines not to collapse results from the same host, so this can be a problem only if the search engine limits the number of documents it indexes per host.

(4) *Indexing Depth*. We assumed the first  $d$  (where  $d$  is some constant,  $d$  was set to 10,000 in our experiments) terms (and the corresponding phrases) in each document are indexed. If the search engine's indexing depth is smaller than  $d$ , queries that we find  $x$  to match may not return  $x$  as a result. While the exact indexing depth is not disclosed by search engines, and may even vary from document to document, our experiments, as well as those of [Bondar 2006], showed that the majority of the documents are indexed by the first  $d = 10,000$  terms at the least.

(5) *Parsing and tokenization*. Different search engines may have slightly different algorithms for parsing and tokenizing documents. For example, two words separated by a comma may or may not be indexed as a phrase. If our parser determines a sequence of terms in  $x$  to be a phrase, while the search engine's parser does not, the corresponding phrase will not return  $x$  as a result. Conversely, search engine's parser may detect a phrase that our parser does not. We designed our parser to mimic the search engines' parsers as closely as we could.

(6) *Indexing by terms not appearing in document content*. Search engines index documents under terms that do not occur at their text at all (e.g., anchor text terms). Our parser, obviously, will not find a document to match such terms (unless they appear in the document content too).

**6.4.4. Overcoming the Inaccuracies.** Incidence recall deficiency can be easily alleviated by discarding from  $B_{\mathcal{P}}$  edges corresponding to queries that belong to  $queries_{\mathcal{P}}(x)$  but not to  $pqueries_{\mathcal{P}}(x)$ . To this end, we modify  $getIncidentDocs_{\mathcal{P}}(q)$  so that after computing  $results(q)$ , it fetches all the documents in  $results(q)$  and returns only those for which  $q \in pqueries_{\mathcal{P}}(x)$  (note that  $|results(q)|$  is at most  $k$ ).

Incidence precision deficiency is not as easy to handle. Theoretically, we could have submitted all the queries in  $pqueries_P(x)$  to the search engine and discard the ones that do not return  $x$  as a result, but that would have required sending many (sometimes, thousands of) queries per sample document.

To tackle the incidence precision problem more realistically, we note that our algorithms do not really need to know the whole set of incident queries of a given document. Rather, they need: (1) to sample random queries from  $queries_P(x)$ ; and (2) calculate the degree of  $x$ :  $deg_P(x) = |queries_P(x)|$ . Sampling random queries from  $queries_P(x)$  can be done by repeatedly selecting queries from  $pqueries_P(x)$  and submitting them to the search engine until hitting a query  $q$  so that  $x \in results(q)$ . As for the degree calculation, we do not have an accurate method of estimating document degrees. Rather, we use  $|pqueries_P(x)|$  as an (over-)estimate of  $deg_P(x)$ . In our subsequent article [Bar-Yossef and Gurevich 2007], we presented an efficient technique for avoiding the degree overestimation.

In Section 5, we provide a theoretical analysis of the effect of the inaccuracies in degree calculations on the bias of our samplers. From now on, we focus on overflowing queries as the only source of deficiency, as we believe these to be the most significant factor causing degree overestimation.

## 6.5. OTHER ASSUMPTIONS.

**6.5.1. Dynamic Corpora.** Our algorithms assume that search engine corpora do not change during the sampling process. Obviously, this assumption does not hold in practice as search engine indices are constantly being updated. In our experiments we noticed only slight differences in the results returned for the same queries at different steps of the experiment. We note that the duration of our experiments was determined by the limited resources we used. Having more resources could have shortened this duration and drastically diminished the effect of corpus changes.

**6.5.2. Versioned Indices.** Search engines may maintain multiple nonidentical versions of the index simultaneously, and serve users from different versions of the index (based on the user's profile or based on load-balancing criteria). Our algorithms assume all queries are served from a single coherent index. If all the queries are indeed served from the same version of the index, then the results produced by our algorithms reflect properties of the specific index version used. Some anomalies may occur, if the samplers work with multiple index versions simultaneously, assuming the differences among the versions are significant (which we do not believe to be the case in most search engines).

## 7. Pool-Based Sampler

In this section, we describe our pool-based (PB) sampler. The sampler assumes knowledge of an *explicit* and *admissible* query pool  $\mathcal{P}$ . Such a pool can be constructed by crawling a large corpus of web documents, such as the ODP directory, and collecting terms or phrases that occur in its pages. We can run the PB sampler with any such pool, yet the choice of the pool may affect the bias, the recall, and the query and fetch costs of the sampler. In the end of the section, we provide analysis of the impact of the choice of the query pool on the performance of the PB sampler.

Let  $\pi$  be a target distribution on the corpus  $\mathcal{D}$ . We assume our sampler is given a black box procedure  $getWeight_{\hat{\pi}}(x)$  that computes an unnormalized form  $\hat{\pi}$

```

1: Function PBSampler( $SE, C$ )
2:   while (true) do
3:      $X :=$  random document generated by DDSampler( $SE$ )
4:     toss a coin whose heads probability is  $\frac{\hat{\pi}(X)}{C \deg_P(X)}$ 
5:     if (coin comes up heads)
6:       break
7:   return  $X$ 

```

FIG. 6. The outer procedure of the PB sampler.

of  $\pi$ . We denote by  $qcost(\hat{\pi})$  (query cost) and by  $fcost(\hat{\pi})$  (fetch cost) the worst-case number of search engine queries and page fetches, respectively, the procedure uses to compute the weight  $\hat{\pi}(x)$ . As a running example, we think of  $\pi$  as the uniform distribution on  $\mathcal{D}$ . The unnormalized form we use is  $\forall x, \hat{\pi}(x) = 1$ , and thus  $qcost(\hat{\pi}) = fcost(\hat{\pi}) = 0$ .

*Remark 7.1.* There are natural examples of target distributions  $\pi$ , for which  $qcost(\hat{\pi}) > 0$  and/or  $fcost(\hat{\pi}) > 0$ . For example, if  $\pi$  is the distribution of documents by in-degree, then  $qcost(\hat{\pi}) > 0$ , as we need to query a search engine in order to find the number of in-links of a given document. If  $\pi$  is the distribution of documents by out-degree, then  $fcost(\hat{\pi}) > 0$ , because we need to fetch the content of a document in order to find the number of out-links it has.

The PB sampler does not directly generate samples from the target distribution  $\pi$ . Instead, it uses another sampler—the *degree distribution sampler*—that generates samples from the “document degree distribution” (see definition below). An unnormalized form of the document degree distribution can be efficiently computed. The PB sampler therefore applies a Monte Carlo method (e.g., rejection sampling) on the samples from the degree distribution in order to generate samples from  $\pi$ .

**7.1. THE OUTER PROCEDURE.** Recall that the degree of a document  $x \in \mathcal{D}$  is the number of queries it matches:

$$\deg_P(x) = |\text{queries}_P(x)| = |\{q \in \mathcal{P} \mid x \in \text{results}(q)\}|.$$

The distribution of documents by degree is called the *document degree distribution*:

$$d_P(x) = \frac{\deg_P(x)}{\deg_P(\mathcal{D})}.$$

Note that the support of  $d_P$  is exactly  $\mathcal{D}_P$ —the set of documents that are covered by  $\mathcal{P}$ . The document degree distribution has an unnormalized form, which is easy to compute:

$$\hat{d}_P(x) = \deg_P(x).$$

Recall that  $\mathcal{P}$  is an admissible query pool, and thus computing  $\deg_P(x) = |\text{queries}_P(x)|$  can be done by fetching only  $x$  and without submitting queries to the search engine.

Assume for the moment that we already have a sampler (DDSampler) that generates random documents sampled from the degree distribution  $d_P$  (we show in the next subsection how to construct such a sampler). Figure 6 shows the outer function of the PB sampler, which uses DDSampler as a subroutine.



The PB sampler applies rejection sampling with trial distribution  $d_{\mathcal{P}}$  and target distribution  $\pi$ . The unnormalized weights used for document  $x$  are  $\hat{\pi}(x)$  (which is computed by calling the `getWeight $\hat{\pi}$ ( $x$ )` procedure) and  $\hat{d}_{\mathcal{P}}(x) = \deg_{\mathcal{P}}(x)$ . Recall that the latter can be computed by a single page fetch. An envelope constant  $C$  satisfying

$$C \geq \max_{x \in \mathcal{D}_{\mathcal{P}}} \frac{\hat{\pi}(x)}{\deg_{\mathcal{P}}(x)}$$

must be given to the PB sampler as input. In the case  $\pi$  is the uniform distribution on  $\mathcal{D}$ ,  $\hat{\pi}(x) = 1$  for all  $x \in \mathcal{D}$  while  $\deg_{\mathcal{P}}(x) \geq 1$  for all  $x \in \text{supp}(d_{\mathcal{P}}) = \mathcal{D}_{\mathcal{P}}$ . Therefore, in this case, an envelope constant of  $C = 1$  will do. The resulting acceptance probability (Line 4) is simply  $1 / \deg_{\mathcal{P}}(x)$ .

We next analyze the recall and the bias of the PB sampler, under the assumption that DDSampler generates samples from  $d_{\mathcal{P}}$ :

**PROPOSITION 7.2.** *Suppose the sampling distribution of DDSampler is exactly the degree distribution  $d_{\mathcal{P}}$ . Then, the sampling recall of the PB sampler is:*

$$\text{recall}_{\pi}(\mathcal{P}) = \pi(\mathcal{D}_{\mathcal{P}})$$

*and it is a perfect sampler; that is, it has a sampling bias of 0.*

**PROOF.** Let  $\eta$  be the sampling distribution of the PB sampler. We first show that  $\text{supp}(\eta) = \mathcal{D}_{\mathcal{P}} \cap \text{supp}(\pi)$ . Clearly,  $\text{supp}(\eta) \subseteq \mathcal{D}_{\mathcal{P}}$ , because the PB sampler cannot output documents that are not covered by  $\mathcal{P}$ . Also,  $\text{supp}(\eta) \subseteq \text{supp}(\pi)$ , because only documents that have non-zero probability under  $\pi$  can be accepted by the acceptance–rejection procedure. Therefore,  $\text{supp}(\eta) \subseteq \mathcal{D}_{\mathcal{P}} \cap \text{supp}(\pi)$ .

To show containment in the other direction, consider any document  $x \in \mathcal{D}_{\mathcal{P}} \cap \text{supp}(\pi)$ . This means that: (1)  $\deg_{\mathcal{P}}(x) > 0$ ; and (2)  $\hat{\pi}(x) > 0$ . The first condition implies that  $x$  has a positive probability to be selected by DDSampler. The second condition implies that  $x$  has a positive probability to be accepted by the acceptance–rejection procedure. Therefore,  $x$  has an overall positive probability to be returned by the PB sampler and thus  $\mathcal{D}_{\mathcal{P}} \cap \text{supp}(\pi) \subseteq \text{supp}(\eta)$ .

We can now calculate the recall of the PB sampler:

$$\text{recall}_{\pi}(PB) = \pi(\text{supp}(\eta)) = \pi(\mathcal{D}_{\mathcal{P}} \cap \text{supp}(\pi)) = \pi(\mathcal{D}_{\mathcal{P}}) = \text{recall}_{\pi}(\mathcal{P}).$$

Next, we analyze the sampling bias of the PB sampler. To this end, we need to calculate the distance between  $\eta$  and the distribution obtained by restricting  $\pi$  to  $\text{supp}(\eta)$ , that is,  $\pi_{\text{supp}(\eta)}$ . The main thing to observe is that the unnormalized form  $\hat{\pi}$  of  $\pi$  also gives an unnormalized form of  $\pi_{\text{supp}(\eta)}$ . Let  $Z_{\hat{\pi}}$  be the normalization constant of  $\hat{\pi}$ . Define

$$Z_{\hat{\pi}_{\text{supp}(\eta)}} = Z_{\hat{\pi}} \cdot \pi(\mathcal{D}_{\mathcal{P}}).$$

Hence, for every  $x \in \text{supp}(\eta)$ , we have:

$$\pi_{\text{supp}(\eta)}(x) = \frac{\pi(x)}{\pi(\text{supp}(\eta))} = \frac{\hat{\pi}(x)}{Z_{\hat{\pi}} \cdot \pi(\mathcal{D}_{\mathcal{P}} \cap \text{supp}(\pi))} = \frac{\hat{\pi}(x)}{Z_{\hat{\pi}} \cdot \pi(\mathcal{D}_{\mathcal{P}})} = \frac{\hat{\pi}(x)}{Z_{\hat{\pi}_{\text{supp}(\eta)}}}.$$

Therefore,  $\hat{\pi}$  is indeed an unnormalized form of  $\pi_{\text{supp}(\eta)}$ . So the right way to view the PB sampler is as a rejection sampling procedure with  $\pi_{\text{supp}(\eta)}$  (and not  $\pi$ ) as the target distribution and with  $d_{\mathcal{P}}$  as the trial distribution. Note that  $\text{supp}(\pi_{\text{supp}(\eta)}) \subseteq$

```

1: Function DDSampler(SE)
2:  Q := random query generated by QCSampler(SE)
3:  submit Q to the search engine SE
4:  results(Q) := results returned from SE
5:  X := document chosen uniformly at random from results(Q)
6:  return X

```

FIG. 7. The degree distribution sampler.

$\mathcal{D}_{\mathcal{P}} = \text{supp}(d_{\mathcal{P}})$  and hence the necessary precondition of rejection sampling is met. It now follows from the analysis of rejection sampling that  $\eta = \pi_{\text{supp}(\eta)}$ .  $\square$

We note that one could implement the PB sampler with other Monte Carlo methods as well. We chose to present here rejection sampling, due to its simplicity.

**7.2. DEGREE DISTRIBUTION SAMPLER.** Next, we describe DDSampler—the sampler that samples documents from the degree distribution. To this end, we need to sample queries from the query pool according to the *query cardinality distribution*. Recall that the cardinality of a query is the number of documents in its result set. The distribution of queries by cardinality is defined as:

$$c_{\mathcal{P}}(q) = \frac{\text{card}(q)}{\text{card}(\mathcal{P})}.$$

In Figure 7, we describe the degree distribution sampler. For the time being, we make two unrealistic assumptions: (1) There is a sampler QCSampler that samples queries from the cardinality distribution  $c_{\mathcal{P}}$ . (This seems unrealistic, because we do not know a priori the cardinalities of all the queries in  $\mathcal{P}$ , and so it is not clear how to sample queries proportionally to their cardinalities.) (2) No query in  $\mathcal{P}$  overflows. (This is unrealistic, because it is not clear how to create a large explicit pool of queries that does not have even a single overflowing query.) We later show how to remove these assumptions.

The sampler is very similar to the Bharat–Broder sampler, except that it samples random queries proportionally to their cardinalities. Since no query overflows, all documents that match a query are included in its result set. It follows that the probability of a document to be sampled is proportional to the number of queries in  $\mathcal{P}$  that it matches:

**PROPOSITION 7.3.** *Suppose the sampling distribution of QCSampler is the query cardinality distribution and that  $\mathcal{P}$  has no overflowing queries. Then, the sampling distribution of DDSampler is  $d_{\mathcal{P}}$ .*

**PROOF.** Let  $p$  be the sampling distribution of DDSampler, and let  $X$  denote a random document selected by DDSampler. Let  $Q$  denote a random query chosen from the cardinality distribution  $c_{\mathcal{P}}$ . To calculate  $p(x)$ , we expand over all choices for  $Q$ :

$$p(x) = \Pr_p(X = x) = \sum_{q \in \mathcal{P}} \Pr_{p, c_{\mathcal{P}}}(X = x | Q = q) \cdot \Pr_{c_{\mathcal{P}}}(Q = q).$$

Note that given  $Q = q$ , the probability that  $X = x$  is 0 if  $x \notin \text{results}(q)$  and is  $1/\text{card}(q)$  otherwise. Hence, the only terms left in the sum are ones that belong to

$\text{queries}_{\mathcal{P}}(x)$ :

$$\begin{aligned} \sum_{q \in \mathcal{P}} \Pr_{\mathcal{P}, c_{\mathcal{P}}} (X = x | Q = q) \cdot \Pr_{c_{\mathcal{P}}} (Q = q) &= \sum_{q \in \text{queries}_{\mathcal{P}}(x)} \frac{1}{\text{card}(q)} \cdot \frac{\text{card}(q)}{\text{card}(\mathcal{P})} \\ &= \frac{|\text{queries}_{\mathcal{P}}(x)|}{\text{card}(\mathcal{P})} = \frac{\deg_{\mathcal{P}}(x)}{\text{card}(\mathcal{P})}. \end{aligned}$$

Since  $\text{card}(\mathcal{P}) = \deg_{\mathcal{P}}(\mathcal{D})$  (Proposition 6.4), then the left-hand side of the last expression equals  $\frac{\deg_{\mathcal{P}}(x)}{\deg_{\mathcal{P}}(\mathcal{D})} = d_{\mathcal{P}}(x)$ .  $\square$

We next address the unrealistic assumption that none of the queries in  $\mathcal{P}$  overflows. Rather than using  $\mathcal{P}$ , which is likely to have overflowing queries, we use the query pool  $\mathcal{P}_+$  (recall that  $\mathcal{P}_+$  is the set of valid queries in  $\mathcal{P}$ ).  $\mathcal{P}_+$  does not have any overflowing queries by definition.

In the next subsection, we show an efficient implementation of QCSampler that generates samples from  $c_{\mathcal{P}_+}$  (the cardinality distribution of  $\mathcal{P}_+$ ) rather than from  $c_{\mathcal{P}}$ . Since  $\mathcal{P}_+$  has no overflowing queries, then by Proposition 7.3, the sampling distribution of DDSampler in this case equals the degree distribution  $d_{\mathcal{P}_+}$  induced by  $\mathcal{P}_+$ .

**7.2.1. Recall and Bias Analysis.** Let us now return to the outer function of the PB sampler. That function assumed DDSampler generates samples from  $d_{\mathcal{P}}$ . What happens if instead it generates samples from  $d_{\mathcal{P}_+}$ ? Note that now there is a mismatch between the trial distribution used by the PB sampler (i.e.,  $d_{\mathcal{P}_+}$ ) and the unnormalized weights it uses (i.e.,  $\deg_{\mathcal{P}}(x)$ ).

One possible solution could be to try to compute the unnormalized weights of  $d_{\mathcal{P}_+}$ , that is,  $\hat{d}_{\mathcal{P}_+}(x) = \deg_{\mathcal{P}_+}(x)$ . However, this is impossible to do efficiently, because  $\mathcal{P}_+$  is no longer an admissible query pool. Instead, we opt for a different solution: we leave the outer function of the PB sampler as is; that is, the trial distribution will be  $d_{\mathcal{P}_+}$  but the unnormalized weights will remain those of  $d_{\mathcal{P}}$  (i.e.,  $\deg_{\mathcal{P}}(x)$ ). This means that the PB sampler is in fact an approximate rejection sampling procedure, and we can thus use Theorem 5.1 to bound its sampling bias.

Theorem 7.5 below bounds the recall and the bias of the PB sampler. The upper bound on the bias is given in terms of a property of documents, which we call the *validity density*:

**Definition 7.4 (Validity Density).** Let  $\mathcal{P}$  be a query pool. The *validity density* of a document  $x \in \mathcal{D}_{\mathcal{P}}$  relative to  $\mathcal{P}$  is:

$$\text{vdensity}_{\mathcal{P}}(x) = \frac{\deg_{\mathcal{P}_+}(x)}{\deg_{\mathcal{P}}(x)}.$$

That is, the validity density of  $x$  is the fraction of valid queries among the queries from  $\mathcal{P}$  that  $x$  matches. The bias of the PB sampler is bounded by half the normalized mean deviation of the validity density of documents, where documents are weighted by the target distribution. Hence, if all documents have more-or-less the same validity density, then we can expect the PB sampler to be accurate.

**THEOREM 7.5.** *The sampling recall of the PB sampler is equal to:*

$$\text{recall}_{\pi}(\mathcal{P}_+) = \pi(\mathcal{D}_{\mathcal{P}_+}).$$

*The sampling bias of the PB sampler is at most:*

$$\frac{1}{2} \text{ndev}_{\pi_{\mathcal{P}_+}}(\text{vdensity}_{\mathcal{P}}(\mathbf{X})),$$

where  $\pi_{\mathcal{P}_+}$  is the restriction of  $\pi$  to  $\mathcal{D}_{\mathcal{P}_+} \cap \text{supp}(\pi)$ .

For the proof, see Appendix C.

As we later show in Table I, the normalized mean deviation of the validity density of our test pool is relatively small (0.34), implying the bias is at most 0.17.

Another factor that affects the variance of the validity density is the fraction of invalid queries among the queries in the pool. If invalid queries are rare, then the validity density of most documents will be close to 1, implying the variance of the validity density is small. This is formalized by Theorem 7.6 below.

To state the theorem, we first need to define a distribution over queries, induced by the distribution  $\pi$  on documents, with respect to which we measure the overflow probability. For every document  $x \in \mathcal{D}$ , we define the “weight” of  $x$  to be its probability under the target distribution  $\pi_{\mathcal{P}_+}$ , that is,  $\pi_{\mathcal{P}_+}(x)$ . The *weight* of a query is the sum of all the weights it “absorbs” from the documents it is connected to:

$$w_{\mathcal{P}}(q) = \sum_{x \in \text{results}(q)} \frac{\pi_{\mathcal{P}_+}(x)}{\deg_{\mathcal{P}}(x)}.$$

Note that each document  $x$  distributes its weight  $\pi_{\mathcal{P}_+}(x)$  among all the queries it is connected to. Thus, its contribution to one query  $q$  is only  $\pi_{\mathcal{P}_+}(x)/\deg_{\mathcal{P}}(x)$ . It follows that the sum of all query weights equals the sum of all document weights:

$$w_{\mathcal{P}}(\mathcal{P}) = \sum_{q \in \mathcal{P}} w_{\mathcal{P}}(q) = \sum_{x \in \mathcal{D}} \pi_{\mathcal{P}_+}(x) = 1.$$

We can thus view  $w_{\mathcal{P}}$  as a probability distribution over  $\mathcal{P}$ . We call this distribution the *query weight distribution*.

**THEOREM 7.6.** *The sampling bias of the PB sampler is at most:*

$$\frac{\text{ovprob}(w_{\mathcal{P}})}{1 - \text{ovprob}(w_{\mathcal{P}})}.$$

That is, if the overflowing queries in the pool have a relatively low mass under the query weight distribution (i.e., they are few in number and they do not overflow by “much”), then the bias of the sampler is low. The proof of the theorem appears in Appendix C.

**7.3. CARDINALITY DISTRIBUTION SAMPLER.** We are left to show how to efficiently sample queries from  $\mathcal{P}$  in accordance with the cardinality distribution  $c_{\mathcal{P}_+}$ . Sampling queries uniformly from  $\mathcal{P}$  is easy, since we have  $\mathcal{P}$  in explicit form. But how do we sample queries from  $\mathcal{P}_+$  proportionally to their cardinalities? This seems impossible to do, because we do not know a priori which queries belong to  $\mathcal{P}_+$  and what are the cardinalities of these queries.

Our most crucial observation is that an unnormalized form of  $c_{\mathcal{P}_+}$  can be computed efficiently. Given a query  $q \in \mathcal{P}$ , a corresponding unnormalized weight is

```

1: Function QCSampler(SE)
2:   k := SE.result_limit
3:   while (true) do
4:     Q := uniformly chosen query from  $\mathcal{P}$ 
5:     submit Q to the search engine SE
6:     card(Q) := number of results returned from SE
7:     if (card(Q) > k)
8:       W := 0
9:     else
10:      W := card(Q)
11:    toss a coin whose heads probability is  $\frac{W}{k}$ 
12:    if (coin comes up heads)
13:      break
14:  return Q

```

FIG. 8. The cardinality distribution sampler.

the following:

$$\hat{c}_{\mathcal{P}_+}(q) = \begin{cases} \text{card}(q), & \text{if } \text{card}(q) \leq k, \\ 0, & \text{if } \text{card}(q) > k. \end{cases}$$

$\hat{c}_{\mathcal{P}_+}(q)$  can be computed by submitting  $q$  to the search engine and counting the number of matches it has.

*Remark 7.7.* Since we need to know  $\text{card}(q)$  exactly only when  $q$  does not overflow, then we can compute  $\text{card}(q)$  by physically counting the number of results returned by the search engine on the query  $q$ . We do not need to rely on the number of results reported by the search engine, which is notoriously inaccurate.

Now, since we know  $c_{\mathcal{P}_+}$  in unnormalized form, we can apply rejection sampling with the uniform distribution on  $\mathcal{P}$  as the trial distribution and with  $c_{\mathcal{P}_+}$  as the target distribution. This will give us samples from  $c_{\mathcal{P}_+}$ .

The query cardinality sampler (QCSampler) is depicted in Figure 8. The sampler applies rejection sampling with the cardinality distribution  $c_{\mathcal{P}_+}$  as the target distribution and with the uniform distribution on  $\mathcal{P}$  (which we denote by  $u_{\mathcal{P}}$ ) as the trial distribution. The unnormalized form used for the target distribution is  $\hat{c}_{\mathcal{P}_+}$ , as described above. The unnormalized form used for the trial distribution is:

$$\forall q \in \mathcal{P}, \hat{u}_{\mathcal{P}}(q) = 1.$$

Since for every  $q \in \mathcal{P}$ ,  $\hat{c}_{\mathcal{P}_+}(q) \leq k$ , then

$$\max_{q \in \mathcal{P}} \frac{\hat{c}_{\mathcal{P}_+}(q)}{\hat{u}_{\mathcal{P}}(q)} \leq k,$$

and thus the sampler uses the envelope constant  $C = k$ .

The following now follows directly from the correctness of rejection sampling:

**PROPOSITION 7.8.** *The sampling distribution of QCSampler is  $c_{\mathcal{P}_+}$ .*

**7.4. COST ANALYSIS.** We now analyze the query cost and the fetch cost of the PB sampler.

THEOREM 7.9. *The query cost of the PB sampler is at most:*

$$\frac{C \cdot \deg_{\mathcal{P}_+}(\mathcal{D}_{\mathcal{P}_+})}{Z_{\hat{\pi}} \cdot \pi(\mathcal{D}_{\mathcal{P}_+}) \cdot (1 - \text{ovprob}(w_{\mathcal{P}}))} \cdot \left( \frac{|\mathcal{P}|}{|\mathcal{P}_+|} \cdot \frac{k}{\text{avg}_{q \in \mathcal{P}_+} \text{card}(q)} + \text{qcost}(\hat{\pi}) \right).$$

*The fetch cost of the PB sampler is at most:*

$$\frac{C \cdot \deg_{\mathcal{P}_+}(\mathcal{D}_{\mathcal{P}_+})}{Z_{\hat{\pi}} \cdot \pi(\mathcal{D}_{\mathcal{P}_+}) \cdot (1 - \text{ovprob}(w_{\mathcal{P}}))} \cdot (1 + \text{fcost}(\hat{\pi})).$$

The proof can be found in Appendix C.

The above expressions may seem hard to parse, so we would like to simplify and interpret them, at least for the case  $\pi$  is the uniform distribution on  $\mathcal{D}$ .

When  $\pi$  is uniform,  $Z_{\hat{\pi}} = |\mathcal{D}|$  and  $\pi(\mathcal{D}_{\mathcal{P}_+}) = |\mathcal{D}_{\mathcal{P}_+}|/|\mathcal{D}|$ . Therefore, the term  $Z_{\hat{\pi}} \cdot \pi(\mathcal{D}_{\mathcal{P}_+})$  is  $|\mathcal{D}_{\mathcal{P}_+}|$ . Also, an envelope constant  $C = 1$  can be used in this case. It follows that

$$\frac{C \cdot \deg_{\mathcal{P}_+}(\mathcal{D}_{\mathcal{P}_+})}{Z_{\hat{\pi}} \cdot \pi(\mathcal{D}_{\mathcal{P}_+})} = \frac{\deg_{\mathcal{P}_+}(\mathcal{D}_{\mathcal{P}_+})}{|\mathcal{D}_{\mathcal{P}_+}|} = \text{avg}_{x \in \mathcal{D}_{\mathcal{P}_+}} \deg_{\mathcal{P}_+}(x).$$

Also,  $\text{qcost}(\hat{\pi}) = \text{fcost}(\hat{\pi}) = 0$  in this case. Therefore, the query cost is:

$$\text{avg}_{x \in \mathcal{D}_{\mathcal{P}_+}} \deg_{\mathcal{P}_+}(x) \cdot \frac{1}{1 - \text{ovprob}(w_{\mathcal{P}})} \cdot \frac{|\mathcal{P}|}{|\mathcal{P}_+|} \cdot \frac{k}{\text{avg}_{q \in \mathcal{P}_+} \text{card}(q)}.$$

Thus, the query cost is the product of four components: (1) The average degree of documents that are covered by the valid queries in  $\mathcal{P}$ ; (2) The inverse of the “validity probability” of queries, when selected proportionally to their weights; (3) The ratio between the total number of queries in  $\mathcal{P}$  and the valid queries in  $\mathcal{P}$ ; and (3) The ratio between the maximum cardinality of valid queries (i.e.,  $k$ ) and the average cardinality of valid queries.

Similarly, the fetch cost in this case is:

$$\text{avg}_{x \in \mathcal{D}_{\mathcal{P}_+}} \deg_{\mathcal{P}_+}(x) \cdot \frac{1}{1 - \text{ovprob}(w_{\mathcal{P}})}.$$

7.5. CHOOSING THE QUERY POOL. We next review the parameters of the query pool that impact the PB sampler.

*Pool’s Recall.* The sampler’s recall equals the recall of  $\mathcal{P}_+$ —the pool of valid queries among the queries in  $\mathcal{P}$ . Therefore, we would like pools whose valid queries cover most of the documents in the corpus  $\mathcal{D}$ . In order to guarantee such high recall, the pool must consist of enough terms/phrases that are not too popular (and thus would not overflow), but yet almost every document in  $\mathcal{D}$  contains at least one of them. We can obtain such a collection of terms/phrases by crawling a large corpus of web documents, such as the ODP directory.

*Validity Density Deviation.* The bias of the PB sampler is bounded by the normalized mean deviation of the validity density of documents, where documents are weighted by the target distribution. That is, in order to keep the bias low, we need to make sure that all documents have roughly similar validity densities. If the query pool has few overflowing queries and the queries that do overflow do not overflow by much, then we should expect most documents to have a validity



density that is close 1, implying that the variance of the validity density is small. Obtaining such a pool whose recall is still high may be tricky. A pool consisting of disjunctions of terms, for example, may be problematic, because such queries are likely to overflow. We thus opted for exact phrase queries. Our experiments indicate that phrases of length at least 5 are unlikely to overflow. If the phrases are collected from a sufficiently large and representative corpus, then the corresponding recall is still reasonable.

*Average Degree.* The query and fetch costs depend on the average degree of documents that are covered by the valid queries in the pool. Hence, we would like to find pools for which the degree of documents grows moderately with the document length. Exact phrase queries are a good example, because then the degree of documents grows linearly with the document length. Conjunctions or disjunctions of  $m$  terms are poor choices, because there the growth rate is exponential in  $m$ .

*Overflow and Underflow Probabilities.* High density of overflowing queries in the pool has two negative effects: (1) it potentially increases the sampling bias of the sampler; and (2) it increases the query and fetch costs. High density of underflowing queries does not impact the sampling bias, but may increase the query cost. We therefore would like to keep the overflow and underflow probabilities as small as possible.

*Average Cardinality.* The query cost depends also on the ratio between the maximum cardinality of valid queries ( $k$ ) and the average cardinality of valid queries. We would like thus the average cardinality to be as close as possible to  $k$ . Of course, this may interfere with the overflow probability: if the average cardinality is too high, many queries will simply overflow.

## 8. Random Walk Based Sampler

We propose two variants of a random walk sampler. One is based on the Metropolis–Hastings algorithm and another on the Maximum Degree method. Both samplers perform a random walk on a virtual graph whose nodes are the documents indexed by the search engine.

Like the pool-based sampler, this sampler too selects its queries from a fixed admissible query pool  $\mathcal{P}$ . However, here the pool may be *implicit* rather than explicit, and thus does not require a pre-processing step for constructing the query pool.

Let  $\pi$  be a target distribution on the corpus  $\mathcal{D}$ , so that  $\text{supp}(\pi) = \mathcal{D}$  (recall our simplifying assumption from Section 4.4). As with the pool-based sampler, we assume access to an oracle procedure  $\text{getWeight}_{\hat{\pi}}(x)$  that computes an unnormalized form  $\hat{\pi}$  of  $\pi$ .  $\text{qcost}(\hat{\pi})$  and  $\text{fcost}(\hat{\pi})$  denote, respectively, the worst-case query and fetch costs of this oracle procedure.

**8.1. DOCUMENT GRAPH.** The random walk sampler performs a random walk on a virtual graph  $G_{\mathcal{P}}$ , which we call the *document graph*.<sup>6</sup>  $G_{\mathcal{P}}$  is obtained from the queries-documents graph  $B_{\mathcal{P}}$  defined in Section 6.3.  $G_{\mathcal{P}}$  is an undirected weighted graph whose vertex set is  $\mathcal{D}$ —the documents indexed by the search engine. Two documents  $x, y$  are connected by an edge in  $G_{\mathcal{P}}$  if and only if they share a neighboring

<sup>6</sup>The random walk is actually performed on  $G_{\mathcal{P}_+}$ , which is derived from the pool of valid queries in  $\mathcal{P}$ . See details below.

query  $q$  in the graph  $B_{\mathcal{P}}$ . The edge weight is the number of such shared neighbor queries, where each query is normalized by its cardinality. In other words,  $(x, y)$  is an edge if and only if  $\text{queries}_{\mathcal{P}}(x) \cap \text{queries}_{\mathcal{P}}(y) \neq \emptyset$ . The weight of the edge  $(x, y)$  is:

$$\text{weight}_{\mathcal{P}}(x, y) = \sum_{q \in \text{queries}_{\mathcal{P}}(x) \cap \text{queries}_{\mathcal{P}}(y)} \frac{1}{\text{card}(q)}.$$

The *degree* of a document  $x$  in  $G_{\mathcal{P}}$  is defined as:

$$\deg_{G_{\mathcal{P}}}(x) = \sum_{y \in \mathcal{D}} \text{weight}_{\mathcal{P}}(x, y).$$

We next observe that the degree of a document  $x$  in the graph  $G_{\mathcal{P}}$  coincides with its degree in the graph  $B_{\mathcal{P}}$ :

PROPOSITION 8.1. *For every document  $x \in \mathcal{D}$ ,*

$$\deg_{G_{\mathcal{P}}}(x) = |\text{queries}_{\mathcal{P}}(x)| = \deg_{B_{\mathcal{P}}}(x).$$

PROOF. For every triple  $(x, y, q)$ , where  $x, y \in \mathcal{D}$  are documents and  $q \in \mathcal{P}$  is a query, define the predicate  $A(x, y, q)$  to be 1 if and only if both  $x$  and  $y$  belong to  $\text{results}(q)$ . Now, we use the predicate  $A$  to rewrite the degree of a document  $x$ :

$$\begin{aligned} \deg_{G_{\mathcal{P}}}(x) &= \sum_{y \in \mathcal{D}} \text{weight}_{\mathcal{P}}(x, y) \\ &= \sum_{y \in \mathcal{D}} \sum_{q \in \text{queries}_{\mathcal{P}}(x) \cap \text{queries}_{\mathcal{P}}(y)} \frac{1}{\text{card}(q)} \\ &= \sum_{y \in \mathcal{D}} \sum_{q \in \mathcal{P}} \frac{A(x, y, q)}{\text{card}(q)} \\ &= \sum_{q \in \mathcal{P}} \frac{1}{\text{card}(q)} \sum_{y \in \mathcal{D}} A(x, y, q). \end{aligned}$$

If  $q \in \text{queries}_{\mathcal{P}}(x)$ , then  $\sum_{y \in \mathcal{D}} A(x, y, q) = |\text{results}(q)| = \text{card}(q)$ . However, if  $q \notin \text{queries}_{\mathcal{P}}(x)$ , then  $\sum_{y \in \mathcal{D}} A(x, y, q) = 0$ . Hence, we have:

$$\begin{aligned} \sum_{q \in \mathcal{P}} \frac{1}{\text{card}(q)} \sum_{y \in \mathcal{D}} A(x, y, q) &= \sum_{q \in \text{queries}_{\mathcal{P}}(x)} \frac{1}{\text{card}(q)} \cdot \text{card}(q) \\ &= |\text{queries}_{\mathcal{P}}(x)| = \deg_{B_{\mathcal{P}}}(x). \quad \square \end{aligned}$$

8.2. SKELETON OF THE RANDOM WALK SAMPLER. The random walk sampler runs a random walk on the document graph  $G_{\mathcal{P}_+}$  (like the pool-based sampler, it ignores invalid queries). The transition matrix  $P$  of a simple random walk on this graph is the following:

$$P(x, y) = \frac{\text{weight}_{\mathcal{P}_+}(x, y)}{\deg_{\mathcal{P}_+}(x)}.$$

That is, having visited a node  $x$  in the graph, a neighbor  $y$  is chosen proportionally to the weight of the edge connecting  $x$  and  $y$ . It can be shown that this random walk is a reversible Markov Chain whose limit distribution is the document degree

```

1: Function RWSampler( $SE, B, x_0$ )
2:  $X := x_0$ 
3: for  $t = 1$  to  $B$  do
4:    $Y := \text{sampleNeighbor}(SE, X)$ 
5:   if ( $\text{accept}(X, Y)$ )
6:      $X := Y$ 
7: return  $X$ 

1: Function  $\text{accept}(x, y)$ 
2: compute  $r_{\text{MCMC}}(x, y)$ 
3: toss a coin whose heads probability is  $r_{\text{MCMC}}(x, y)$ 
4: return true if and only if coin comes up heads

```

FIG. 9. Skeleton of the random walk sampler.

distribution  $d_{\mathcal{P}_+}$  (recall definition from Section 7.1). In order to transform this simple random walk into a Markov Chain that converges to the target distribution  $\pi$ , an MCMC algorithm is applied. In this article, we focus on the Metropolis–Hastings and the Maximum Degree methods.

The skeleton of the random walk sampler is described in Figure 9. The sampler runs a simple random walk on the graph  $G_{\mathcal{P}_+}$  augmented with an acceptance–rejection procedure. Having visited a node  $X$  in the graph, the function  $\text{sampleNeighbor}$  selects a random neighbor  $Y$  with probability  $P(X, Y) = \frac{\text{weight}_{\mathcal{P}_+}(X, Y)}{\deg_{\mathcal{P}_+}(X)}$  (the implementation of this function is described below). An acceptance–rejection procedure is then applied on  $X$  and  $Y$  in order to determine whether  $Y$  should be accepted as the next step of the random walk. The choice of the acceptance function depends on the particular MCMC method used. The two acceptance functions we employ are:

$$r_{\text{MH}}(x, y) = \min \left\{ \frac{\pi(y) P(y, x)}{\pi(x) P(x, y)}, 1 \right\} = \min \left\{ \frac{\pi(y) \deg_{\mathcal{P}_+}(x)}{\pi(x) \deg_{\mathcal{P}_+}(y)}, 1 \right\}; \text{ and}$$

$$r_{\text{MD}}(x) = \frac{\hat{d}_{\mathcal{P}_+}(x)}{C \hat{\pi}(x)} = \frac{\deg_{\mathcal{P}_+}(x)}{C \hat{\pi}(x)},$$

where

$$C \geq \max_{x \in \mathcal{D}_{\mathcal{P}_+}} \frac{\deg_{\mathcal{P}_+}(x)}{\hat{\pi}(x)}$$

is an envelope constant.

In order to implement the random walk sampler, we need to address four issues: (1) how to select the start node  $x_0$ ? (2) how to set the length of the burn-in period  $B$ ? (3) how to implement the neighbor sampling procedure? and (4) how to calculate the acceptance functions  $r_{\text{MH}}$  and  $r_{\text{MD}}$ ?

**8.3. SELECTING THE START NODE.** The graph  $G_{\mathcal{P}_+}$  is not necessarily connected. When we choose the start node  $x_0$  from some connected component  $F$  of  $G_{\mathcal{P}_+}$ , then the random walk will never reach nodes outside  $F$ . This implies that the Markov Chain we produce will not necessarily converge to the target distribution  $\pi$ , but

rather to the distribution  $\pi_F$  obtained by restricting  $\pi$  to  $F$ :

$$\pi_F(x) = \frac{\pi(x)}{\pi(F)}, \quad \text{for all } x \in F.$$

Therefore, the recall of the sampler in this case will be at most  $\pi(F)$ . In order to maximize recall, we would like the start node to belong to the component that has the highest mass under  $\pi$ . We do not have any rigorous technique for making such a selection. On the other hand, we speculate that for sufficiently rich query pools,  $G_{\mathcal{P}_+}$  is expected to have a giant connected component, and thus almost any document chosen as a start node will do. Our experimental results (see Section 9.3) support this speculation, as the largest connected component of  $G_{\mathcal{P}_+}$  in a small search engine that we built constituted close to 99% of the nodes.

**8.4. SETTING THE BURN-IN PERIOD.** As shown in Section 4.4.3, the main factor that determines the length of the burn-in period is the *spectral gap* of the underlying transition matrix. Thus, in order to set  $B$ , we need an estimate of the spectral gaps of the transition matrices  $P_{\text{MH}}$  and  $P_{\text{MD}}$  obtained by applying the MH and the MD methods, respectively, on the simple random walk on  $G_{\mathcal{P}_+}$ .

We experimentally estimated these gaps for a small search engine that we built. The results, discussed more thoroughly in Section 9.3, provide the following estimates:

$$\alpha(P_{\text{MD}}) \geq \frac{1}{20,000}, \quad \alpha(P_{\text{MH}}) \geq \frac{1}{20,000}.$$

As the connectivity of  $G_{\mathcal{P}_+}$  in larger search engines is expected to be similar to the connectivity of  $G_{\mathcal{P}_+}$  in the small search engine, we expect similar bounds to be applicable also to random walks on real search engines. See more details in Section 9.3.

**8.5. SAMPLING NEIGHBORS.** We now address the problem of sampling neighbors in accordance with the transition matrix  $P(x, y)$ . The naive method to select such a random neighbor would be the following: given a document  $x$ , find all valid queries  $q \in \text{queries}_{\mathcal{P}_+}(x)$  that match  $x$ , choose one of them at random, submit the query to the search engine, and then pick one of its results at random. The only problem with this algorithm is that we do not know how to compute the set  $\text{queries}_{\mathcal{P}_+}(x)$  efficiently, since  $\mathcal{P}_+$  is not an admissible query pool.

The solution to the above problem emanates from the following observation:  $\text{queries}_{\mathcal{P}_+}(x)$  is a subset of  $\text{queries}_{\mathcal{P}}(x)$ , which we can compute efficiently, since  $\mathcal{P}$  is an admissible query pool. So we can simply select random queries from  $\text{queries}_{\mathcal{P}}(x)$  until hitting a valid query. This random valid query will be uniform in  $\text{queries}_{\mathcal{P}_+}(x)$ . The neighbor sampling procedure, described in Figure 10, implements this idea.

The following proposition proves the correctness of the neighbor sampling procedure:

**PROPOSITION 8.2.** *Let  $x$  be any document in  $\mathcal{D}_{\mathcal{P}_+}$  and let  $Y$  be the random neighbor selected by the procedure `sampleNeighbor`, when given  $x$  as input. Then, for every neighbor  $y$  of  $x$  in the graph  $G_{\mathcal{P}_+}$ ,*

$$\Pr(Y = y) = \frac{\text{weight}_{\mathcal{P}_+}(x, y)}{\deg_{\mathcal{P}_+}(x)}.$$

```

1: Function sampleNeighbor( $SE, x$ )
2:  $queries_{\mathcal{P}}(x) := \text{getIncidentQueries}_{\mathcal{P}}(x)$ 
3: while (true) do
4:    $Q :=$  query chosen uniformly from  $queries_{\mathcal{P}}(x)$ 
5:   submit  $Q$  to the search engine  $SE$ 
6:   if ( $Q$  neither overflows nor underflows)
7:     break
8:    $results(Q) :=$  results returned from  $SE$ 
9:    $Y :=$  document chosen uniformly at random from  $results(Q)$ 
10: return  $Y$ 

```

FIG. 10. The neighbor sampling procedure.

PROOF. To calculate  $\Pr(Y = y)$ , we expand over all possibilities for the random query  $Q$  chosen from  $queries_{\mathcal{P}_+}(x)$ . Obviously,  $\Pr(Q = q) = 0$  for all  $q \notin queries_{\mathcal{P}_+}(x)$ .

$$\Pr(Y = y) = \sum_{q \in queries_{\mathcal{P}_+}(x)} \Pr(Y = y | Q = q) \cdot \Pr(Q = q).$$

For fixed  $q$  and  $y$ ,  $\Pr(Y = y | Q = q) = \frac{1}{\text{card}(q)}$ , if  $q \in queries_{\mathcal{P}_+}(y)$ , and  $\Pr(Y = y | Q = q) = 0$ , otherwise.  $\Pr(Q = q) = \frac{1}{|queries_{\mathcal{P}_+}(x)|} = \frac{1}{\deg_{\mathcal{P}_+}(x)}$ . Therefore,

$$\begin{aligned}
& \sum_{q \in queries_{\mathcal{P}_+}(x)} \Pr(Y = y | Q = q) \cdot \Pr(Q = q) \\
&= \sum_{q \in queries_{\mathcal{P}_+}(x) \cap queries_{\mathcal{P}_+}(y)} \frac{1}{\text{card}(q)} \cdot \frac{1}{\deg_{\mathcal{P}_+}(x)} \\
&= \frac{\text{weight}_{\mathcal{P}_+}(x, y)}{\deg_{\mathcal{P}_+}(x)}. \quad \square
\end{aligned}$$

8.6. CALCULATING THE ACCEPTANCE FUNCTIONS. Next, we address the issue of how to calculate the acceptance functions of the MH and the MD samplers.

The acceptance function of the MH algorithm is:

$$r_{\text{MH}}(x, y) = \min \left\{ \frac{\pi(y) \deg_{\mathcal{P}_+}(x)}{\pi(x) \deg_{\mathcal{P}_+}(y)}, 1 \right\}.$$

The acceptance function of the MD method is:

$$r_{\text{MD}}(x) = \frac{\deg_{\mathcal{P}_+}(x)}{C \hat{\pi}(x)}, \quad \text{where } C \geq \max_{x \in \mathcal{D}_{\mathcal{P}_+}} \frac{\deg_{\mathcal{P}_+}(x)}{\hat{\pi}(x)}.$$

The problem is that we cannot compute the degrees  $\deg_{\mathcal{P}_+}(x)$  and  $\deg_{\mathcal{P}_+}(y)$  efficiently, since  $\mathcal{P}_+$  is not an admissible query pool. What we do instead is apply perturbed acceptance functions:

$$r'_{\text{MH}}(x, y) = \min \left\{ \frac{\pi(y) \deg_{\mathcal{P}}(x)}{\pi(x) \deg_{\mathcal{P}}(y)}, 1 \right\}$$

and

$$r'_{\text{MD}}(x) = \frac{\deg_{\mathcal{P}}(x)}{C' \hat{\pi}(x)}, \quad \text{where } C' \geq \max_{x \in \mathcal{D}_{\mathcal{P}_+}} \frac{\deg_{\mathcal{P}}(x)}{\hat{\pi}(x)}.$$

(Note that when  $\pi$  is the uniform distribution,  $C'$  should be an upper limit on the maximum degree of documents.)

$r'_{\text{MH}}(x, y)$  and  $r'_{\text{MD}}(x)$  can be computed efficiently, because  $\mathcal{P}$  is an admissible query pool. The problem is that now the acceptance functions and the base Markov Chain  $P$  on which they are applied are mismatching, and thus the limit distributions are no longer guaranteed to equal the target distribution  $\pi$ . This scenario is the one captured by the approximate Metropolis–Hastings and the approximate Maximum Degree procedures, described in Section 5.

Before we analyze the sampling bias and sampling recall of the resulting samplers, let us identify the exact form of the target distribution, trial distribution, and approximate trial distribution employed by the approximate MH and MD procedures.

Let  $F$  be the connected component of  $G_{\mathcal{P}_+}$  to which the start vertex  $x_0$  of the random walk belongs. Let  $d_F$  be the restriction of the degree distribution  $d_{\mathcal{P}_+}$  to  $F$ :

$$d_F(x) = \frac{d_{\mathcal{P}_+}(x)}{d_{\mathcal{P}_+}(F)}, \quad \text{for all } x \in F. \quad (1)$$

As the random walk can never reach nodes outside  $F$ , then  $d_F$ , rather than  $d_{\mathcal{P}_+}$ , is the limit distribution of the simple random walk on  $G_{\mathcal{P}_+}$  that starts at  $x_0 \in F$ . Therefore, the trial distribution is  $d_F$ .

Similarly, as the random walk can never reach nodes outside  $F$ , the real target distribution when applying the MH and the MD samplers with  $x_0 \in F$  is the restriction of  $\pi$  to  $F$ , that is,  $\pi_F$ . Indeed, it can be easily verified that the restriction of the unnormalized form of  $\pi$  to  $F$  constitutes an unnormalized form of  $\pi_F$ .

Finally, the approximate trial distribution used by the two procedures is the restriction of the degree distribution  $d_{\mathcal{P}}$  to  $F$ :

$$q_F(x) = \frac{d_{\mathcal{P}}(x)}{d_{\mathcal{P}}(F)}, \quad \text{for all } x \in F.$$

Since  $\text{supp}(\pi) = \mathcal{D}$ ,  $\text{supp}(d_{\mathcal{P}_+}) = \mathcal{D}_{\mathcal{P}_+}$ ,  $\text{supp}(d_{\mathcal{P}}) = \mathcal{D}_{\mathcal{P}}$ , and  $F \subseteq \mathcal{D}_{\mathcal{P}_+} \subseteq \mathcal{D}_{\mathcal{P}} \subseteq \mathcal{D}$ , then  $\text{supp}(\pi_F) = \text{supp}(d_F) = \text{supp}(q_F) = F$ . Therefore,  $\pi_F, d_F, q_F$  satisfy the requirements of the approximate MH and MD procedures. Furthermore, the simple random walk on  $F$  constitutes a reversible Markov Chain, as  $G_{\mathcal{P}_+}$  is an undirected graph. Therefore, the necessary precondition of the approximate MH procedure is met.

Applying Theorems 5.6 and 5.5, we know that the random walks performed by the approximate MH and MD procedures have the following limit distribution  $\eta$ :

$$\eta(x) = \pi_F(x) \frac{\frac{d_F(x)}{q_F(x)}}{\mathbb{E}_{\pi_F} \left( \frac{d_F(X)}{q_F(X)} \right)}. \quad (2)$$

The following theorem bounds the sampling bias and the sampling recall of the MH and MD samplers:



```

1: Function OptimizedMDSampler( $SE, B, x_0, C'$ )
2:  $X := x_0$ 
3:  $t = 0$ 
4: while ( $t < B$ ) do
5:    $\text{delay} :=$  a geometric random variable whose success parameter is:  $\frac{\deg_P(x)}{C' \cdot \hat{\pi}(x)}$ 
6:    $t := t + \text{delay}$ 
7:   if ( $t \geq B$ ) break
8:    $Y := \text{sampleNeighbor}(SE, X)$ 
9:    $X := Y$ 
10:   $t := t + 1$ 
11: return  $X$ 

```

FIG. 11. The optimized Maximum Degree sampler.

**THEOREM 8.3.** *Let  $\varepsilon > 0$ . Suppose we run the MH sampler (respectively, the MD sampler) with a burn-in period  $B$  that guarantees the approximate MH Markov Chain (respectively, approximate MD Markov Chain) reaches a distribution, which is at distance of at most  $\varepsilon$  from the limit distribution. Then, the sampling bias of the MH sampler (respectively, MD sampler) is at most:*

$$\frac{1}{2} \text{ndev}_{\pi_F}(\text{vdensity}_P(X)) + \varepsilon.$$

*The sampling recall of the MH sampler (respectively, MD sampler) is at least:*

$$\pi(F) \cdot (1 - \frac{1}{2} \text{ndev}_{\pi_F}(\text{vdensity}_P(X)) - \varepsilon).$$

The proof appears in Appendix D.

**8.7. OPTIMIZED MD SAMPLER.** As discussed in Section 4.4.2, the special form of the acceptance function of the MD sampler allows for an optimized implementation. Since the acceptance function depends only on the current state  $x$  and not on the proposed state  $y$ , then rather than selecting a new proposed neighbor  $Y$  every time the acceptance–rejection procedure is invoked, we can select the neighbor only once, after acceptance is achieved. Further optimization is possible by selecting a priori the number of steps the random walk is going to spend at the state  $x$ , without actually performing the iterative coin tosses.

In this spirit, we describe in Figure 11 an optimized version of the MD sampler. Note that the fact the proposed neighbor is chosen only once results in significant savings in the number of search engine queries made. We analyze these savings below.

**8.8. COST ANALYSIS FOR THE MH AND MD SAMPLERS.** We now analyze the query and the fetch costs of the MH and MD samplers. We first analyze the query cost incurred by the random walk when it visits a particular node  $x$ :

**PROPOSITION 8.4.** *The expected number of queries the (MH or MD) random walk sampler submits to the search engine when visiting a node  $x$  is at most*

$$\frac{1}{\text{vdensity}_P(x)} + \text{qcost}(\hat{\pi}).$$

PROOF. Search engine queries are made in two places: (1) in the sample Neighbor procedure; and (2) (possibly) in the  $\text{getWeight}_{\hat{\pi}}$  procedure in order to compute the unnormalized target weights.

When  $\text{sampleNeighbor}$  is called with a document  $x$  as input, the procedure repeatedly selects random queries from  $\text{queries}_{\mathcal{P}}(x)$  until hitting a valid query. Therefore, the expected number of queries made in such a call is:

$$\frac{|\text{queries}_{\mathcal{P}}(x)|}{|\text{queries}_{\mathcal{P}_+}(x)|} = \frac{\deg_{\mathcal{P}}(x)}{\deg_{\mathcal{P}_+}(x)} = \frac{1}{\text{vdensity}_{\mathcal{P}}(x)}.$$

The procedure  $\text{getWeight}_{\hat{\pi}}$  is called once for every candidate neighbor  $y$ . Each invocation of the procedure makes at most  $\text{qcost}(\hat{\pi})$  search engine queries.

We conclude that the expected query cost of the step is at most:

$$\frac{1}{\text{vdensity}_{\mathcal{P}}(x)} + \text{qcost}(\hat{\pi}). \quad \square$$

The actual cost of the random walk at a certain step depends on the distribution of the random node  $X$  that is visited at this step. This cost may vary greatly between the early stages of the walk, where the distribution of visited nodes depends on the walk's starting point, and the later stages of the walk, where the distribution of visited nodes is close to the limit distribution. Note that the former cost is tricky to bound, because the distribution of the starting node may be arbitrary. We therefore present two bounds: (1) A bound on the worst-case cost of a step, which applies to all steps, regardless of the distribution of visited nodes; this bound is applicable in particular to the early steps of the random walk. (2) A bound on the expected cost of a step, assuming the distribution of visited nodes is close to the limit distribution; this bound is applicable mainly to the later stages of the random walk. We start with the worst-case bound:

**PROPOSITION 8.5.** *For any step of the random walk, the expected number of queries the (MH or MD) sampler submits to the search engine is at most*

$$\max_{x \in F} \deg_{\mathcal{P}}(x) + \text{qcost}(\hat{\pi}).$$

PROOF. By Proposition 8.4, the expected number of queries submitted to the search engine is at most:

$$\begin{aligned} & \max_{x \in \text{supp}(\eta_t)} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}(x)} + \text{qcost}(\hat{\pi}) \right) \\ & \leq \max_{x \in F} \left( \frac{\deg_{\mathcal{P}}(x)}{\deg_{\mathcal{P}_+}(x)} \right) + \text{qcost}(\hat{\pi}) \\ & \leq \max_{x \in F} \deg_{\mathcal{P}}(x) + \text{qcost}(\hat{\pi}). \end{aligned} \quad \square$$

At the later stages of the walk, where the distribution of  $X$  is known to be close to the limit distribution, we can achieve a much better bound:

**THEOREM 8.6.** *For any  $t \geq 0$ , the expected query cost of the  $t$ th step of the random walk sampler (whether MH or MD) is at most:*

$$\frac{1}{\mathbb{E}_{\pi_F}(\text{vdensity}_{\mathcal{P}}(\mathbf{X}))} + 3\varepsilon_t \sqrt{|F| \cdot \mathbb{E}_{u_F} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}^2(\mathbf{X})} \right)} + \text{qcost}(\hat{\pi}).$$

Here,  $\varepsilon_t$  is an upper bound on the total variation distance between the distribution of the node visited at the  $t$ th step and the limit distribution  $\eta$  of the random walk.  $u_F$  is the uniform distribution over  $F$ .

The proof appears in Appendix D. Using the spectral gap bound on the convergence rate of random walks (see the proof of Corollary 4.9),  $\varepsilon_t$  can be chosen as:

$$\varepsilon_t = \frac{(1 - \alpha)^t}{\eta_{\min}},$$

where  $\alpha$  is the spectral gap of the random walk's transition matrix and  $\eta_{\min} = \min_{x \in F} \eta(x)$ .

The above theorem therefore implies that for

$$t \geq \Omega \left( \log \frac{1}{\eta_{\min}} + \log |F| + \log \mathbb{E}_{u_F} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}^2(\mathbf{X})} \right) \right),$$

the query cost of the  $t$ th step is about  $1 / \mathbb{E}_{\pi_F}(\text{vdensity}_{\mathcal{P}}(\mathbf{X}))$ . In case  $\pi_F$  equals the uniform distribution  $u_F$ ,  $\text{qcost}(\hat{\pi}) = 0$ , and thus the query cost becomes roughly  $1 / \mathbb{E}_{u_F}(\text{vdensity}_{\mathcal{P}}(\mathbf{X}))$ .

The fetch cost of the random walk sampler does not depend on the distribution of visited nodes and is quantified in the following theorem.

**THEOREM 8.7.** *For any  $t \geq 0$ , the expected fetch cost of the  $t$ th step of the random walk sampler (whether MH or MD) is at most:*

$$1 + \text{fcost}(\hat{\pi}).$$

**PROOF.** Page fetches are made only in the `sampleNeighbor` procedure and in the `getWeight $\hat{\pi}$`  procedure. In `sampleNeighbor(x)` only  $x$  is fetched, in order to compute `queries $\mathcal{P}$ (x)`. In `getWeight $\hat{\pi}$ (x)`, at most `fcost( $\hat{\pi}$ )` pages are fetched in order to compute  `$\hat{\pi}(x)$` . So the fetch cost is at most:

$$1 + \text{fcost}(\hat{\pi}). \quad \square$$

**8.9. COST ANALYSIS FOR THE OPTIMIZED MD SAMPLER.** An optimized MD random walk can be viewed as performing a simple random walk on  $F$  and augmenting it with different “delays” at the different nodes visited along the walk. We call the steps of the simple random walk *real* steps, and the steps performed during the delay *free* steps. As free steps incur zero query and fetch costs, the amortized query and the fetch costs of a real step are reduced by a multiplicative factor of  $1 + \text{delay}$  compared to the unoptimized MD sampler (where all steps are real). Thus, in order to obtain upper bounds on the amortized costs of real steps of the optimized MD sampler, it will suffice to lower bound the delays associated with these steps.

The following proposition analyzes the expected delay when the simple random walk visits a node  $x$ .

PROPOSITION 8.8. *The expected delay of the optimized MD sampler when visiting a node  $x$  is*

$$\frac{C' \hat{\pi}(x)}{\deg_{\mathcal{P}}(x)}.$$

PROOF. In accordance with Figure 11, when visiting a node  $x$ , the delay is a geometric random variable whose success probability is  $\frac{\deg_{\mathcal{P}}(x)}{C' \hat{\pi}(x)}$ .  $\square$

For example, when  $\pi$  is the uniform distribution, then  $\hat{\pi}(x) = 1$  and  $C'$  is an upper limit on the maximum degree. Therefore, the expected delay is the ratio between the maximum degree and the degree of  $x$ .

Let  $X$  be a random node visited at the  $t$ th real step of the optimized MD random walk. It follows from the above proposition that the expected delay at this step is  $\frac{C' \hat{\pi}(X)}{\deg_{\mathcal{P}}(X)}$ . In order to lower bound the expectation of this delay, we need to analyze the distribution of  $X$ . When  $t$  is small, this distribution can be quite arbitrary, and therefore it is difficult to bound the corresponding expected delay. We note, however, that the delay is always non-negative, and therefore  $t$  real steps of the optimized MD sampler correspond to at least  $t$  steps of the unoptimized MD sampler. Since the upper bound on the expected query cost of the unoptimized MD sampler is monotonically nonincreasing with  $t$  (see the previous subsection), we conclude that for any  $t$ , the upper bound on the expected query cost of the  $t$ th step of the unoptimized MD random walk is applicable also to the optimized case.

Similarly to the unoptimized case, at the later stages of the walk, when the distribution of  $X$  is close to the limit distribution, we can achieve a much better bound.

THEOREM 8.9. *For any  $t \geq 0$ , the expected delay of the  $t$ th real step of the optimized MD random walk is at least:*

$$C' \cdot Z_{\hat{\pi}_F} \cdot \pi(F) \cdot \left( \frac{\mathbb{E}_{\pi_F}(\text{vdensity}_{\mathcal{P}}(X))}{\deg_{\mathcal{P}_+}(F)} - 3\varepsilon_t \cdot \sqrt{\mathbb{E}_{u_F} \left( \frac{\pi_F^2(X)}{\deg_{\mathcal{P}}^2(X)} \right)} \right).$$

Here,  $\varepsilon_t$  is an upper bound on the total variation distance between the distribution of the node visited at the  $t$ th real step and the degree distribution  $d_F$ .  $u_F$  is the uniform distribution over  $F$ .

The proof appears in Appendix D.

The real steps of the optimized MD random walk induce a *simple* random walk (i.e., without acceptance–rejection) on  $F$ . As we saw before, the limit distribution of this random walk is the degree distribution  $d_F$ . Hence,  $\varepsilon_t$  is a bound on the distance between the  $t$ th step of the simple random walk and the limit distribution of this walk.

Using the spectral gap bound on the convergence rate of random walks (see the proof of Corollary 4.9),  $\varepsilon_t$  can be chosen as:

$$\varepsilon_t = \frac{(1 - \alpha)^t}{d_{F \min}},$$

where  $\alpha$  is the spectral gap of the simple random walk's transition matrix and  $d_{F \min} = \min_{x \in F} d_F(x)$ .

The above theorem therefore implies that for

$$t \geq \Omega \left( \log \frac{1}{d_{F\min}} + \log |F| + \log \mathbb{E}_{u_F} \left( \frac{\pi_F^2(X)}{\deg_P^2(X)} \right) \right),$$

the expected delay of the  $t$ th real step is about

$$C' \cdot Z_{\hat{\pi}_F} \cdot \pi(F) \cdot \frac{\mathbb{E}_{\pi_F}(\text{vdensity}_P(X))}{\deg_{P_+}(F)}.$$

When  $\pi$  is the uniform distribution and assuming  $\pi(F) = 1$ , this expression can be further simplified. In this case, (1)  $C'$  is an upper limit on the maximum degree  $\max_{x \in F} \deg_P(x)$ ; (2)  $Z_{\hat{\pi}_F} = |F|$ ; and (3)  $\text{qcost}(\hat{\pi}) = 0$ . For  $t$  large enough, the costs are reduced by a multiplicative factor of:

$$\begin{aligned} & 1 + \frac{\max_{x \in F} \deg_P(x) |F|}{\deg_{P_+}(F)} \mathbb{E}_{u_F}(\text{vdensity}_P(X)) \\ &= 1 + \frac{\max_{x \in F} \deg_P(x)}{\text{avg}_{x \in F} \deg_P(x)} \cdot \mathbb{E}_{u_F}(\text{vdensity}_P(X)). \end{aligned}$$

**8.10. CHOOSING THE QUERY POOL.** We next review the parameters of the implicit query pool that impact the random walk based sampler.

*Pool's Recall.* The considerations here are similar to the ones for the PB sampler (see Section 7.5). The main difference is that there is no need to crawl a large corpus in order to achieve satisfactory recall.

*Validity Density.* Validity density has a direct effect on the query cost of the random walk sampler. The lower is the expected validity density (under the target distribution), the higher is the query cost of the sampler, as it needs to submit more queries before encountering a valid one. Therefore, we would like to find pools having as little invalid queries as possible.

*Validity Density Deviation.* Similarly to the PB sampler (see Section 7.5), the recall and the bias of the random walk sampler is bounded by the normalized mean deviation of the validity density of documents, where documents are weighted by the target distribution.

*Average Degree.* Similarly to the PB sampler (see Section 7.5), the query and the fetch costs of the optimized MD sampler depend on the average degree of documents in  $F$ .

*Maximum Document Degree.* The query cost of the optimized MD sampler depends also on the ratio between the maximum document degree and the average document degree. We would like thus to limit maximum document degree as much as possible. A possible way to achieve this is to set a threshold on the number of queries we extract from a document. Of course, this may interfere with the desire to increase the pool's recall: if we set the threshold too low, the recall may decrease significantly.

*The Burn-In Period.* The query and the fetch costs of the random walk sampler are directly proportional to the random walk's burn-in period. The burn-in period depends on the structure of the graph induced by the query pool. Although we did

not analyze the dependence of the burn-in period on the query pool, the following intuition applies. In each step of the random walk, we want to move to a document, which is as unrelated to the current document as possible. This will ensure that we need only a few steps to obtain a sufficiently random document. If we choose the query pool to be “too specific” (e.g., a pool of long exact phrases), we are likely to spend many consecutive steps in small subsets of related documents, and thus increase the required burn-in period. Conversely, a query pool consisting mainly of single terms is likely to result in a shorter burn-in period, as two documents sharing a single term are not necessarily related to each other. On the other hand, single-term queries are more likely to overflow, thus reducing the validity density, which in turn incurs elevated costs as discussed above.

## 9. Experimental Results

We conducted four sets of experiments: (1) *pool measurements*: estimation of parameters of selected query pools; (2) *spectral gap estimations*: measurement of the spectral gaps of the transition matrices used by the random walk samplers; (3) *evaluation experiments*: evaluation of the bias of our samplers and the Bharat–Broder (BB) sampler; and (4) *exploration experiments*: measurements of real search engines.

**9.1. EXPERIMENTAL SETUP.** In order to conduct the first three sets of experiments, we built a home-made search engine over a corpus of 2.4 million documents obtained from the Open Directory Project (ODP). The ODP directory crawl consisted of 3 million pages, of which we kept only the ones that we could successfully fetch and parse, that were in text, HTML, or pdf format, and that were written in English. Each page was given a serial id, stored locally, and indexed by single terms and phrases. Only the first 10,000 terms in each page were considered. We used static ranking by document id to rank query results. Different experiments used different values of the result limit  $k$ . See more details below.

The first three sets of experiments were performed on a dual Intel Xeon 2.8-GHz processor workstation with 2-GB RAM and two 160-GB disks.

**9.2. POOL MEASUREMENTS.** In the first set of experiments, we wanted to measure the pool parameters that impact the quality and the efficiency of our samplers. In order to get a variety of results, we measured four different query pools: a single terms pool and three pools of exact phrases of lengths 3, 5, and 7. (We measured only four pools, because each measurement required substantial disk space and running time.)

In order to construct the pools, we extracted all the terms or the  $n$ -grams (with overlaps) of the corresponding length from the ODP documents.  $n$ -grams were not allowed to cross boundaries, such as paragraph boundaries. We split the ODP data set into two parts: a *training set*, consisting of every fifth page (when ordered by id), and a *test set*, consisting of the rest of the pages. The pools were built only from the training data, but the measurements were done only on the test data. In order to determine whether a query is overflowing, we set a result limit of  $k = 20$ .

All our measurements were made with respect to the uniform target distribution. We measured the following parameters: (1) the pool’s size (total number of



TABLE I. RESULTS OF POOL PARAMETER MEASUREMENTS

Parameter	Single terms	Phrases (3)	Phrases (5)	Phrases (7)
Training pool size	2.6M	97.5M	155.9M	151.1M
Fraction of overflowing queries	11.4%	3%	0.4%	0.1%
Fraction of underflowing queries	40.3%	56%	76.2%	82.1%
Average cardinality of valid queries	4.7	3.6	2.2	1.7
Recall of valid queries	61.9%	96.8%	88.6%	64.5%
Average document degree (valid queries)	5.1	77	47.1	37.1
Average document degree (all queries)	293.4	246.8	75.9	51.3
Normalized mean deviation of validity density	0.8	0.34	0.34	0.4
Overflow probability of query weight distribution	0.98	0.7	0.43	0.3

TABLE II. POOL PARAMETER MEASUREMENTS ON YAHOO!

Parameter	Single terms	Phrases (3)	Phrases (5)	Phrases (7)
Fraction of overflowing queries	49.3%	18%	3.3%	1%
Fraction of underflowing queries	5.4%	4.3%	10.3%	14.7%
Average cardinality of valid queries	104.1	107.6	47.9	20.6

queries); (2) the fraction of overflowing queries; (3) the fraction of underflowing queries; (4) the average cardinality of valid queries; (5) the recall of valid queries (i.e.,  $|\mathcal{D}_{p_+}|/|\mathcal{D}|$ ); (6) the average degree of documents relative to valid queries (i.e.,  $\text{avg}_{x \in \mathcal{D}_{p_+}} \text{deg}_{p_+}(x)$ ); (7) the average degree of documents relative to all queries (i.e.,  $\text{avg}_{x \in \mathcal{D}_{p_+}} \text{deg}_p(x)$ ); (8) the normalized mean deviation of the validity density (i.e.,  $\text{ndev}_{\pi_{p_+}}(\text{vdensity}_p(X))$ ); (9) the overflow probability of the query weight distribution (i.e.,  $\text{ovprob}(w_p)$ ). The results of our measurements are tabulated in Table I.

The measurements show that fraction of overflowing queries, average document degree, normalized mean deviation of validity density, and overflow probability improve as phrase length increases, while fraction of underflowing queries, recall, and average cardinality get worse. The results indicate that a phrase length of 5 achieves the best tradeoff among the parameters. It has very few overflowing queries (about 0.5%), while maintaining a recall of about 89%. The overflow probability of the query weight distribution on 3-term phrases is too high (70%), while the recall of the 7-term phrases is way too low (about 64%). The fraction of overflowing queries among single terms is surprisingly small. The explanation is that many of the terms are misspellings, technical terms, or digit strings. The overflow probability of single terms, though, is very high.

Since the ODP data set is presumably representative of the web, we expect most of these measurements to represent the situation on real search engines. The only exceptions are the fraction of overflowing and underflowing queries and the average cardinality of valid queries, which are distorted due to the small size of the ODP data set. We thus measured these parameters on the Yahoo! search engine. In this experiment we used real Yahoo!'s result limit  $k = 1,000$ . The results are given in Table II. It is encouraging to see that for 5-term phrases, the fraction of overflowing queries remains relatively low, while the fraction of underflowing queries goes down dramatically. The elevated fraction of overflowing queries among 3-term phrases is more evident here.

TABLE III. ESTIMATED BIAS, RECALL, AND COST OF THE POOL-BASED SAMPLER

Parameter	Single terms	Phrases (3)	Phrases (5)	Phrases (7)
Sampling bias	0.4	0.17	0.17	0.2
Sampling recall	61.9%	96.8%	88.6%	64.5%
Query cost	5407.4	3070	1996.6	3052
Fetch cost	255	256.7	82.6	53

Assuming the above measurements represent the situation on real search engines, we can use Theorems 7.5 and 7.9 to derive estimates of the sampling bias, sampling recall, query cost, and fetch cost for the pool-based sampler. The results are given in Table III. In the estimations we used the measurements on the Yahoo! search engine, and thus set  $k = 1,000$ .

**9.3. SPECTRAL GAP ESTIMATIONS.** We wanted to estimate the burn-in periods of the two random walk samplers: the MH sampler and the MD sampler. By Corollary 4.9, the main factor that determines the burn-in period is the spectral gap of the random walk's transition matrix. In order to obtain an accurate estimate of the spectral gaps with respect to a real search engine  $SE$ , we would have had to construct the adjacency matrix of the document graph  $G_{\mathcal{P}_+}$  corresponding to  $SE$ , transform this matrix into the two transition matrices  $P_{MH}$  and  $P_{MD}$ , and then estimate the spectral gaps of these matrices.<sup>7</sup> However, since we did not have full access to the corpus of real search engines, we could not explicitly construct the matrices  $P_{MH}$  and  $P_{MD}$ .

Our only option then was to estimate the spectral gaps of the matrices corresponding to our home-made ODP search engine. Since the ODP corpus is a diverse set of documents, presumably representing the whole web, we expect the document graph of the ODP search engine to have similar connectivity properties as the document graph of real search engines. Spectral gap is a “structural” property of a transition matrix, which depends only on the connectivity in the document graph, and not on the size of the graph. Therefore, estimations of the spectral gaps corresponding to the ODP search engine could be representative of the spectral gaps corresponding to real search engines.

In order to make the connectivity of the document graph of the ODP search engine as similar as possible to the connectivity of the document graphs in larger search engines, we set in this experiment  $k = 80$  as the result limit and used a pool of 3-term exact phrases. These settings made the graph more dense, similarly to document graphs that are based on very large corpora.

We constructed the two matrices  $P_{MH}$  and  $P_{MD}$  corresponding to the largest connected component of the document graph. This component contained 98.7% of the documents in the corpus, resulting in two matrices of dimensions 2.37 million by 2.37 million.

Since the matrices were so large, we could not compute their eigenvalues exactly. The largest eigenvalue of any transition matrix is always 1, since the matrix is stochastic. We thus had to estimate only the second largest eigenvalue. To this

<sup>7</sup>ptIn order to get more precise estimates of the required burn-in periods, the spectral gaps of the transition matrices  $P'_{MH}$  and  $P'_{MD}$ , rather than  $P_{MH}$  and  $P_{MD}$ , should have been calculated.

end, we applied the power iteration method (cf. Golub and van Loan [1996]). We obtained the following lower bounds on the two spectral gaps:

$$\alpha(P_{\text{MD}}) \geq \frac{1}{20,000}, \quad \alpha(P_{\text{MH}}) \geq \frac{1}{20,000}.$$

We note that the actual gaps could be larger, yet figuring them out exactly would have required many more iterations of the power method. The fact the two bounds are identical does not mean that the actual spectral gaps of the two Markov Chains are identical. In reality, one of them may be much higher than the other.

Plugging in the above estimates in the formula for the burn-in period (Corollary 4.9) and assuming  $\pi_{\min} = \frac{1}{2 \cdot 10^{10}}$ , as in major search engines, and  $\varepsilon = 0.1$ , we obtained:

$$T_{\varepsilon}(P_{\text{MH}}) = T_{\varepsilon}(P_{\text{MD}}) \leq \frac{1}{\alpha(T)} \left( \ln \frac{1}{\pi_{\min}} + \ln \frac{1}{\varepsilon} \right) \approx 520,000.$$

We can now use the cost analysis from Section 8.8 to estimate the query costs of the MH and the MD samplers. To this end, we: (1) assume that  $G_{\mathcal{P}_+}$  is connected (and thus  $F = \mathcal{D}_{\mathcal{P}_+}$ ); (2) employ the estimates for  $\mathbb{E}_{\pi_{\mathcal{P}_+}}(\text{vdensity}_{\mathcal{P}}(X)) = 1 - \text{ovprob}(w_{\mathcal{P}})$  and for  $\text{avg}_{x \in \mathcal{D}_{\mathcal{P}_+}} \deg_{\mathcal{P}}(x)$  derived from our pool measurements. Plugging in the numbers into the formulae given in Section 8.8, we expect the query cost of the MH sampler to be roughly 1,730,300 queries per sample, while the query cost of the optimized MD sampler (executed with  $C' = 10,000$ ) to be roughly 42,800 queries per sample.

We conclude that both the MH sampler and the MD sampler are significantly less efficient than the pool-based sampler. The cost of the MH sampler is prohibitive. The cost of the optimized MD sampler is much lower, and thus this sampler might be of practical use, if the desired number of samples is small. Recall the big advantage of this approach over the pool-based sampler: it does not need any pre-processing step to create an explicit query pool.

We stress that the above cost estimates are based on our theoretical analysis, which is pessimistic by nature. Our evaluation experiments described below indicate that in practice it may be possible to run the random walks for far fewer steps than is required by the theoretical bounds, yet obtain good samples.

**9.4. EVALUATION EXPERIMENTS.** In order to estimate the biases of our samplers and of the Bharat–Broder (BB) sampler, we ran them on the ODP search engine. In this controlled environment, we could compare the sampling results against the real data.

The corpus of our ODP search engine consisted of the test set only. A result limit of  $k = 5$  was used in order to have an overflow probability comparable to the one on Yahoo!.

We ran five samplers: (1) the PB sampler with rejection sampling; (2) the PB sampler with importance sampling; (3) the MH sampler; (4) the MD sampler; and (5) the BB sampler. All the samplers used a query pool of 5-term phrases extracted from the ODP training set. In order to have a common basis, we allowed each sampler to submit exactly 5 million queries to the search engine.

Running the random walk samplers with the burn-in period dictated by the spectral gap estimations would have been useless, since we would have gotten very few samples from the 5 million queries submitted. We therefore opted for a short burn-in

TABLE IV. NUMBER OF SAMPLES GENERATED BY EACH SAMPLER WHEN RUN OVER THE ODP SEARCH ENGINE WITH THE BUDGET OF 5 MILLION QUERIES

Sampler	# of queries submitted	# of samples
PB + Rejection Sampling	5M	3,276
PB + Importance Sampling	5M	319,682
RW-MH	5M	4,234
RW-MD	5M	7,761
BB	5M	1,129,820

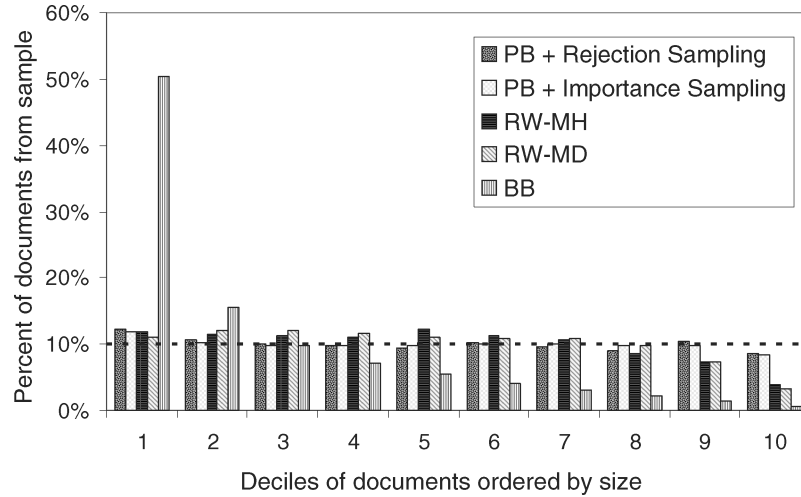


FIG. 12. Distribution of samples by document size.

period of 1,000 steps for the MH sampler and 10,000 steps for the MD sampler (these settings turned out to produce a comparable number of samples to what the PB sampler produced). This enabled us to evaluate whether the shortened burn-in period actually caused the random walk samplers to produce biased samples.

Since each sampler has a different query cost, the samplers produced varying number of samples using the 5 million queries. Table IV shows the actual number of samples generated by each sampler. The BB sampler generated the most samples, yet these samples are highly biased.

In Figure 12, we show the distribution of samples by document size. We ordered the documents in the corpus by size, from largest to smallest, and split them into 10 equally sized deciles. Truly uniform samples should distribute evenly among the deciles. The results show the overwhelming difference between our samplers and the BB sampler. The BB sampler generated a huge number of samples at the top decile (more than 50%! ). Our samplers, on the other hand, had no or little bias. The two PB samplers essentially show no bias, while the RW samplers have small negative bias towards very short documents, possibly due to the shortened burn-in period.

Figure 13 addresses bias towards highly ranked documents. We ordered the documents in the corpus by their static rank, from highest to lowest, and split into 10 equally sized deciles. The first decile corresponds to the most highly ranked documents. The results indicate that none of our samplers had any significant bias

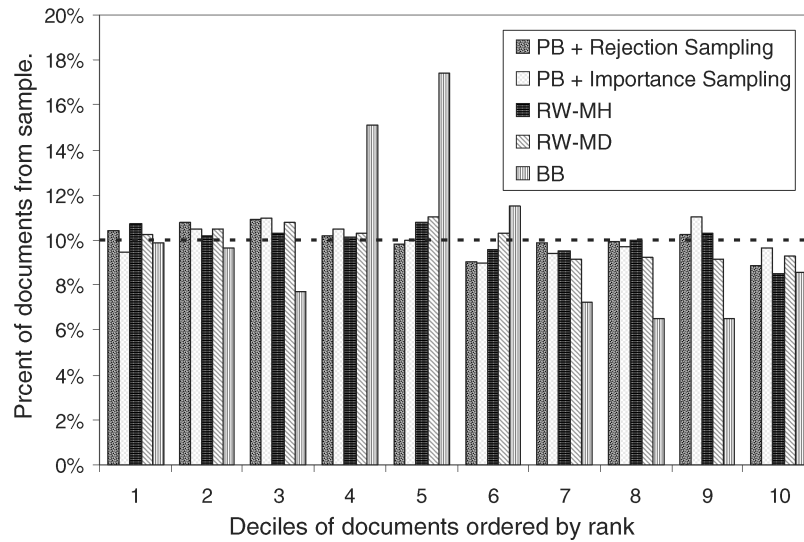


FIG. 13. Distribution of samples by document rank.

under the ranking criterion. Surprisingly, the BB sampler had bias towards the 4th, 5th, and 6th deciles. When digging into the data, we found that documents whose rank (i.e., id) belonged to these deciles had a higher average size than documents with lower or higher rank. Thus, the bias we see here is in fact an artifact of the bias towards long documents. A good explanation is that our 5-term exact phrases pool had a low overflow probability in the first place, so very few queries overflowed. Note that the ranking bias of the BB sampler is created only by overflowing queries.

We have several conclusions from the above experiments: (1) the 5-term phrases pool, which has small overflow probability, made an across-the-board improvement to all the samplers (including BB). This was evidenced by the lack of bias towards highly ranked documents. (2) The BB sampler suffers from a severe bias towards long documents, regardless of the query pool used. (3) Our pool-based samplers seem to give the best results, showing no bias in any of the experiments. (4) The random walk samplers have small negative bias towards short documents. Possibly by increasing the burn-in period of the random walk, this negative bias could be alleviated.

**9.5. EXPLORATION EXPERIMENTS.** We used our most successful sampler, the PB sampler, to generate uniform samples from three major search engines: Google, MSN Search, and Yahoo!.

During the experiments, conducted over a span of three weeks in April–May 2006, we submitted 395,000 queries to Google, 448,000 queries to MSN Search, and 370,000 queries to Yahoo!. Due to legal restrictions on automatic queries, we used the Google, MSN, and Yahoo! Web Search APIs. While automatic queries via APIs are, reportedly, served from slightly different corpora than the corpora used to serve human users, McCown and Nelson [2007] showed that the differences are quite minor. These APIs are limited to submitting only a few thousands of queries a day, which limited the scale of the experiments we could perform. The exploration experiments were conducted on seven machines of varying configurations.

Before describing the results of the experiment, we elaborate on the assumptions made in these experiments and on the procedure we used to test whether a given URL is indexed by a given search engine.

**9.5.1. Setup.** The limit  $k$  on the number of returned results was, at the time of the experiment, 1,000 for Google and Yahoo!, and 250 for MSN Search. As a query pool, we used 5-term phrases extracted from English pages of the ODP data set. The pool contained **666,641,510** phrases (we used the whole ODP corpus, not just the training set).

We discarded all the results that were not in text, HTML, or pdf format. When calculating the cardinality of a query, we did not rely on the number of results reported by the search engine, since this number is notoriously inaccurate. Instead, we explicitly fetched the entire list of available results, and calculated its size.

When submitting queries to the search engines, we specifically indicated not to filter out duplicate results and not to perform any other filtering (host, language, domain, date, etc . . . ). Note that duplicate results *are* filtered out by default, and thus our measurements do not directly reflect duplicate-free corpora but rather the complete “raw” corpora.<sup>8</sup> We choose to make our measurements against the raw corpora in order to prevent different filtering mechanisms employed by search engines from biasing our measurements. Comparing duplicate-free corpora proves even more difficult because each search engine may choose a different “representative” from a set of duplicate pages, thus defeating our URL membership test. Obviously, the raw corpus size should not be used as the only metric of search engine quality. Moreover, it may favor a less sophisticated search engine indexing more duplicate and near-duplicate documents.

Some documents are indexed by terms that do not appear in their content, for example, according to anchor text. Since our samplers have access only to the terms that appear in the document, we rejected samples whose content did not contain the query terms. Similarly, to avoid bias due to index depth, we rejected sampled documents for which the query terms were not found among their first 10,000 terms.

In order to compute queries <sub>$\mathcal{P}$</sub> ( $x$ ), for a given document  $x$ , we first tried to retrieve the cached version of  $x$  from the search engine. If the cached page was not available, we downloaded the page from the web and computed queries <sub>$\mathcal{P}$</sub> ( $x$ ) according to the downloaded version. When any error was encountered while downloading the page, the corresponding sample was dropped.

Note that as a result of dropping some of the samples (either due to fetching errors or due to the absence of the query terms at the first 10,000 terms of the document), the query cardinalities, which were calculated assuming all query results are valid, were not always accurate. Since the frequency of sample drops was low, we speculate that this did not have any significant effect on our measurements.

**9.5.2. URL Membership Testing.** A basic ingredient in our exploration experiments was a procedure, which given a URL  $u$  and search engine  $SE$ , determines whether  $SE$  indexes  $u$  or not. Implementing such a procedure that produces accurate results, while interacting only with the public interface of the search engine, turns

---

<sup>8</sup>Some search engines may filter duplicate documents at crawl/index time. In such cases, even “raw” corpora may not contain all the documents known to search engines.



TABLE V. RELATIVE OVERLAP OF GOOGLE, MSN SEARCH, AND YAHOO!  
(ENGLISH PAGES ONLY)

Samples from ↓ indexed by →	Google	MSN Search	Yahoo!	Both of the other engines
Google		<b>37.2%</b>	<b>45.8%</b>	<b>27.3%</b>
MSN Search	<b>51.2%</b>		<b>52.8%</b>	<b>37.1%</b>
Yahoo!	<b>35.6%</b>	<b>30.9%</b>		<b>20.3%</b>

out to be tricky. The procedure we employed, which combines ideas from previous studies, is described below.

First, we brought the given URL into a canonical form, by converting hex-encoded characters (%XX) to a standard iso-8859-1 characters, converting double slash (//) to a single one, dropping /index.html from the URL's tail, etc.

For each URL tested, we submitted up to seven different queries to the search engine, in order to determine whether the URL is indexed by the search engine. The first query was the URL itself. The second query was "inurl:URL". The last five queries were phrases, from 8 to 14 terms long, extracted randomly from the first 10,000 terms of the document. We sequentially submitted these seven queries to the search engine until we found the URL among the first 100 results (after canonization). If the URL was found by any of the queries, we assumed it is indexed by the search engine. Otherwise, it was assumed not to be indexed.

Obviously, this heuristic procedure is prone to some error. That is, its result is not guaranteed to correctly determine whether a URL is indexed or not.

If a URL was sampled from a search engine, but our procedure failed to find it in the same search engine, we dropped the sample. 5% of samples from Google, and less than 1% of samples from MSN Search and Yahoo! were dropped.

**9.5.3. Corpus Size.** We used importance sampling to estimate the overlaps among Google, MSN Search, and Yahoo!. Table V tabulates the measured relative overlap of the Google, MSN Search, and Yahoo! corpora. We note that since our query pool consisted mostly of English language phrases, our results refer mainly to the English portions of these corpora.

Using Liu's "rule of thumb" (see Section 4.3) and the Chernoff bound, we estimate the overlap results stated above to be within an absolute error of 4.5% from the real values with a confidence level of 95%. There could be an additional absolute error of up to 17%, due to the sampling bias of the PB sampler.

In Figure 14, we show the relative sizes of the English portions of the corpora of the three search engines, as derived from the relative overlap estimates (by dividing  $\text{overlap}[\text{SE1}][\text{SE2}]$  by  $\text{overlap}[\text{SE2}][\text{SE1}]$ ). Note that due to the division, the error in these results can be higher than in the overlap estimates.

**9.5.4. Top-Level Domain Name Distribution.** Figure 15 shows the domain name distributions in the three corpora. Here too, the results are within an absolute error of 4.5% from the real values with a confidence level of 95%. Note that the differences between the search engines are quite minor.

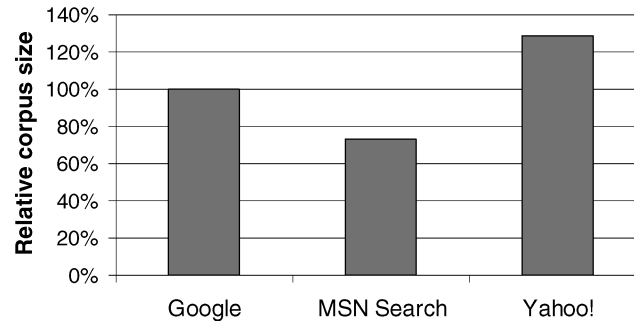


FIG. 14. Relative corpus size of Google, MSN Search, and Yahoo! (English pages only).

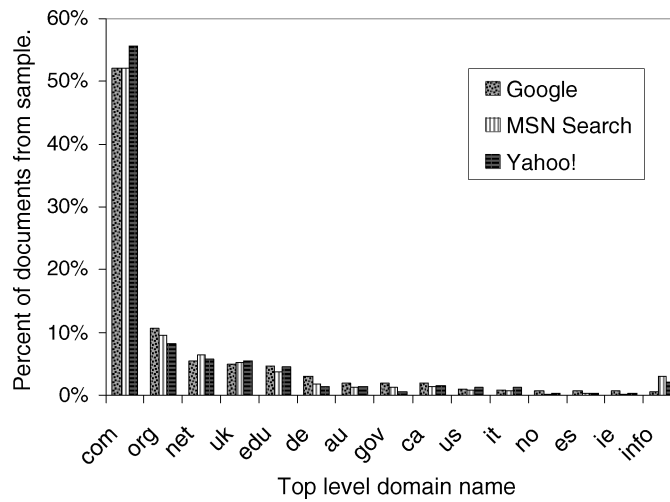


FIG. 15. Top-level domain name distribution of pages in Google, MSN Search, and Yahoo! corpora (English pages only).

**9.5.5. Corpus Freshness.** We evaluated the freshness of the three corpora. First, we checked the percentage of pages indexed that are still valid. Figure 16 shows the percentage of dead pages (ones returning a 4xx HTTP return code) in the three corpora (within an absolute error in each point of 4.5% from the real values with a confidence level of 95%). Note that the Google corpus contains significantly more dead pages than Yahoo! and MSN Search.

Next, we tried to determine to what degree each of the three corpora reflects the current content of indexed web pages. To this end, we assumed that when the cached copy of the document is up-to-date, the index is up-to-date too and vice versa. We compared the two versions of the sampled document: (1) the cached document, and (2) the actual document fetched from the web. These measurements were performed on samples of HTML documents only, for which we were able to fetch both the cached copy and the document itself. Fetching of both document versions was performed simultaneously.

Figure 17 shows, for each value  $0 \leq p \leq 100$ , the fraction of samples, for which at most  $p$  percent of the lines in the cached version are different from the web version. The Yahoo! results look rather bad especially for low  $p$  values. After

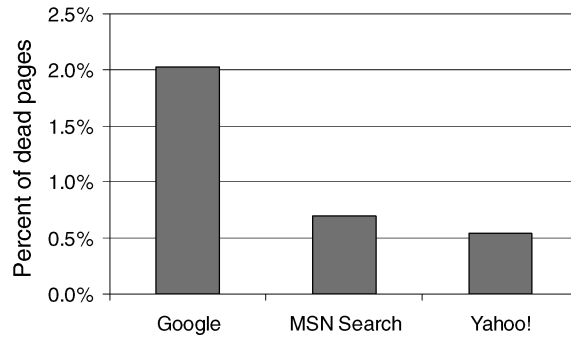


FIG. 16. Percentage of inaccessible pages in the Google, MSN Search, and Yahoo! corpora (English pages only).

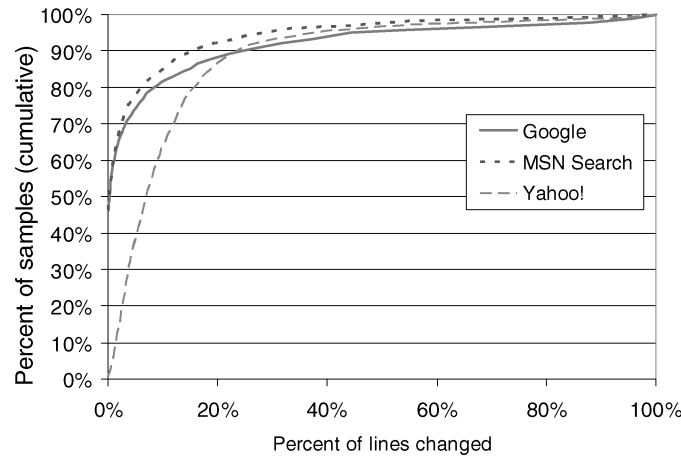


FIG. 17. Raw HTML freshness of the Google, MSN Search, and Yahoo! corpora (English pages only).

looking at the data, we found that some of the cached documents stored by Yahoo! were slightly pre-processed.

To provide a more objective measure of freshness, we measured text-only difference, which would ignore all formatting or other non-essential differences between the two document versions. We compared the bag-of-words representations of the two versions of each sample document. Figure 18 shows, for each value  $0 \leq p \leq 100$ , the fraction of samples, for which at most  $p$  percent of the words in the cached version are different from the web version.

We can see that about 55%–60% of the documents are “fresh” in all three search engines, while Google’s freshness is the lowest and MSN’s is the highest.

**9.5.6. Percentage of Dynamic Pages.** In this experiment, we calculated the percentage of dynamic pages in the three corpora. We deem a page “dynamic” if its URL contains the characters “?” or “&” or if the URL ends with one of the following filename extensions: .php, .php3, .asp, .cfm, .cgi, .pl, .jsp, .exe, .dll. Figure 19 shows substantial differences among the search engines in terms of the percentage of dynamic pages in their corpora: Google has the highest percentage, while MSN has the lowest. This may indicate Google is

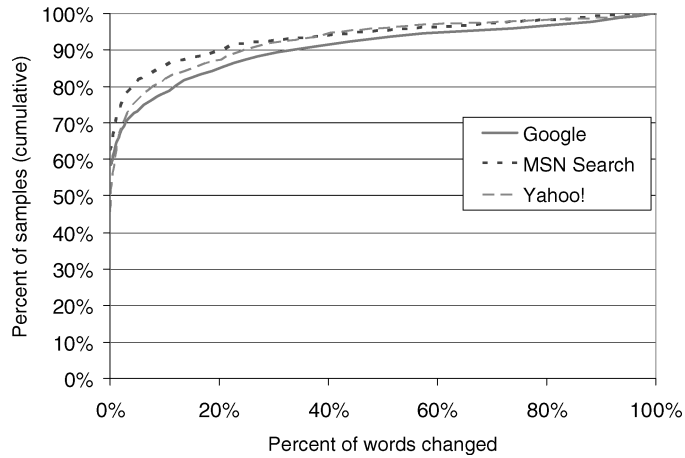


FIG. 18. Text freshness of the Google, MSN Search, and Yahoo! corpora (English pages only).

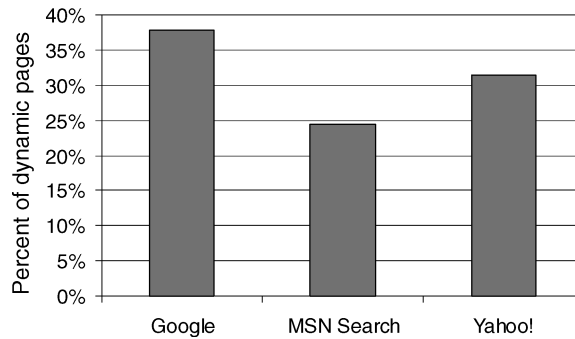


FIG. 19. Percentage of dynamic pages in the Google, MSN Search, and Yahoo! corpora (English pages only).

doing a better job in crawling “deep web” dynamic content. Here too, the results are within an absolute error of 4.5% from the real values with a confidence level of 95%.

## 10. Conclusions

We presented two novel search engine samplers. The first, the pool-based sampler, uses a pre-prepared pool of queries to generate random queries. Random documents are then selected from the results of these queries. The sampler employs a Monte Carlo method, like rejection sampling, to guarantee that the distribution of the samples is close to the target distribution. We show that the sampler works, even if the sampling “weights”, which are needed by the Monte Carlo method, are only approximate.

We provided full analysis of the sampler and identified the query pool parameters that impact its bias and performance. We then estimated these parameters on real data, consisting of the English pages from the ODP hierarchy, and showed

that a pool of 5-term phrases achieves the best tradeoff between accuracy and performance.

Our second sampler runs a random walk on a graph defined over the indexed documents. Its primary advantage is that it does not need a pre-prepared query pool. This sampler employs a Markov Chain Monte Carlo method, like the Metropolis–Hastings algorithm or the Maximum Degree method, to guarantee that the random walk converges to the target distribution. Theoretical bounds on the convergence rate of the random walk are quite high, yet practical evidence suggests that the random walk produces only slightly biased samples much before reaching the theoretical bounds.

Our experimental results bare bias towards pages in the English language, since the query pool we used consisted primarily of English phrase queries. We note that this bias is a limitation of our experimental setup and not of the techniques themselves. The bias could be eliminated by using a more comprehensive query pool.

During our experiments we performed our measurements against “raw” corpora. Duplicate-free measurements remain an interesting open problem.

Although we tested our samplers on web search engines only, we speculate that they may be applicable in a more general setting. For example, for sampling from databases, deep web sites, library records, medical data, etc.

As the query and fetch cost of our samplers are quite high, it remains as an open problem to design much more efficient search engine samplers.

## Appendix

### A. Monte Carlo Methods—Proofs

#### A.1. MARKOV CHAIN MONTE CARLO METHODS

**THEOREM 4.6 (RESTATED).** *The limit distribution of the Markov Chain defined by  $P_{MD}$  is  $\pi$ .*

**PROOF.** In order to show that  $\pi$  is the limit distribution of the Markov Chain defined by  $P_{MD}$ , we will show that

$$\pi P_{MD} = \pi.$$

That is, for each  $x \in \text{supp}(\pi)$ ,

$$\sum_{y \in \text{supp}(\pi)} \pi(y) P_{MD}(y, x) = \pi(x).$$

Substituting

$$P_{MD}(x, y) = \begin{cases} P(x, y) r_{MD}(x), & \text{if } x \neq y, \\ P(x, x) r_{MD}(x) + 1 - r_{MD}(x), & \text{if } x = y. \end{cases},$$

we get

$$\begin{aligned}
\sum_{y \in \text{supp}(\pi)} \pi(y) P_{\text{MD}}(y, x) &= \pi(x)(P(x, x)r_{\text{MD}}(x) + 1 - r_{\text{MD}}(x)) \\
&\quad + \sum_{y \neq x} \pi(y)P(y, x)r_{\text{MD}}(y) \\
&= \pi(x)(1 - r_{\text{MD}}(x)) + \sum_{y \in \text{supp}(\pi)} \pi(y)P(y, x)r_{\text{MD}}(y) \\
&= \pi(x) - \pi(x)r_{\text{MD}}(x) + \sum_{y \in \text{supp}(\pi)} \pi(y)P(y, x)r_{\text{MD}}(y).
\end{aligned}$$

Substituting  $r_{\text{MD}}(x) = \frac{\hat{p}(x)}{C \hat{\pi}(x)}$  we get  $\pi(x)r_{\text{MD}}(x) = \pi(x)\frac{\hat{p}(x)}{C \hat{\pi}(x)} = \frac{Z_{\hat{p}}}{C Z_{\hat{\pi}}} p(x)$ , so the above can be further simplified to:

$$\sum_{y \in \text{supp}(\pi)} \pi(y) P_{\text{MD}}(y, x) = \pi(x) - \frac{Z_{\hat{p}}}{C Z_{\hat{\pi}}} p(x) + \frac{Z_{\hat{p}}}{C Z_{\hat{\pi}}} \sum_{y \in \text{supp}(\pi)} p(y)P(y, x).$$

We know that  $p$  is the limit distribution of the Markov Chain defined by  $P$  and that  $\text{supp}(\pi) = \text{supp}(p)$ , thus  $\sum_{y \in \text{supp}(\pi)} p(y)P(y, x) = p(x)$ . Substituting into the above, we get

$$\sum_{y \in \text{supp}(\pi)} \pi(y) P_{\text{MD}}(y, x) = \pi(x). \quad \square$$

## B. Approximate Monte Carlo Methods—Proofs

### B.1. APPROXIMATE REJECTION SAMPLING

THEOREM 5.1 (RESTATED). *The sampling distribution of the approximate rejection sampling procedure is:*

$$\pi'(x) = \pi(x) \frac{\frac{p(x)}{q(x)}}{\mathbb{E}_{\pi} \left( \frac{p(X)}{q(X)} \right)}.$$

*The expected number of samples from  $p$  needed to generate each sample from  $\pi'$  is:*

$$\frac{C Z_{\hat{q}}}{Z_{\hat{\pi}} \mathbb{E}_{\pi} \left( \frac{p(X)}{q(X)} \right)}.$$

PROOF. We start by finding a closed form expression for  $\pi'$ . Let  $X$  denote a random sample from  $p$ . Let  $Z$  be a 0-1 random variable, which is 1 with probability  $r'_{\text{RS}}(X)$  and 0 otherwise.  $Z$  corresponds to the outcome of the coin toss in the acceptance–rejection procedure of approximate rejection sampling. The procedure returns  $X$  as a sample if and only if  $Z = 1$ . Therefore,  $\pi'$  is the distribution of  $X$  conditioned on the event “ $Z = 1$ ”. In other words, for all  $x \in \text{supp}(\pi)$ ,

$$\pi'(x) = \Pr(X = x | Z = 1).$$



By Bayes rule,

$$\Pr(X = x|Z = 1) = \frac{\Pr(Z = 1|X = x) \cdot \Pr(X = x)}{\Pr(Z = 1)} = \frac{r'_{\text{RS}}(x) p(x)}{\Pr(Z = 1)}.$$

Expanding over all possibilities for  $X$ , we have:

$$\Pr(Z = 1) = \sum_{y \in \text{supp}(\pi)} \Pr(Z = 1|X = y) \cdot \Pr(X = y) = \sum_{y \in \text{supp}(\pi)} r'_{\text{RS}}(y) p(y).$$

Hence,

$$\Pr(X = x|Z = 1) = \frac{r'_{\text{RS}}(x) p(x)}{\sum_{y \in \text{supp}(\pi)} r'_{\text{RS}}(y) p(y)}.$$

Recalling the definition of  $r'_{\text{RS}}$ , we have:

$$\begin{aligned} \pi'(x) = \Pr(X = x|Z = 1) &= \frac{\frac{\hat{\pi}(x)}{C\hat{q}(x)} p(x)}{\sum_{y \in \text{supp}(\pi)} \frac{\hat{\pi}(y)}{C\hat{q}(y)} p(y)} = \frac{\frac{\pi(x)Z_{\hat{\pi}}}{Cq(x)Z_{\hat{q}}} p(x)}{\sum_{y \in \text{supp}(\pi)} \frac{\pi(y)Z_{\hat{\pi}}}{Cq(y)Z_{\hat{q}}} p(y)} \\ &= \frac{\pi(x) \frac{p(x)}{q(x)}}{\sum_{y \in \text{supp}(\pi)} \pi(y) \frac{p(y)}{q(y)}} = \pi(x) \frac{\frac{p(x)}{q(x)}}{\mathbb{E}_{\pi} \left( \frac{p(X)}{q(X)} \right)}. \end{aligned}$$

We now turn to calculating the cost of the approximate rejection sampling procedure. The procedure generates samples from  $p$  until the acceptance–rejection procedure accepts. The probability of acceptance is  $\Pr(Z = 1)$ , and thus the number of samples generated from  $p$  is a geometric random variable whose probability of success is  $\Pr(Z = 1)$ . Its expectation is  $1/\Pr(Z = 1)$ . Let us then calculate this probability:

$$\begin{aligned} \Pr(Z = 1) &= \sum_{y \in \text{supp}(\pi)} r'_{\text{RS}}(y) p(y) = \sum_{y \in \text{supp}(\pi)} \frac{\pi(y)Z_{\hat{\pi}}}{Cq(y)Z_{\hat{q}}} p(y) \\ &= \frac{Z_{\hat{\pi}}}{CZ_{\hat{q}}} \cdot \sum_{y \in \text{supp}(\pi)} \pi(y) \frac{p(y)}{q(y)} \\ &= \frac{Z_{\hat{\pi}}}{CZ_{\hat{q}}} \cdot \mathbb{E}_{\pi} \left( \frac{p(X)}{q(X)} \right). \end{aligned} \quad \square$$

PROPOSITION 5.2 (RESTATED).

$$\|\pi' - \pi\| = \frac{1}{2} \text{ndev}_{\pi} \left( \frac{p(X)}{q(X)} \right).$$

PROOF

$$\begin{aligned}
\|\pi' - \pi\| &= \frac{1}{2} \sum_{x \in \text{supp}(\pi)} |\pi'(x) - \pi(x)| = \frac{1}{2} \sum_{x \in \text{supp}(\pi)} \left| \pi(x) \frac{\frac{p(x)}{q(x)}}{\mathbb{E}_\pi \left( \frac{p(X)}{q(X)} \right)} - \pi(x) \right| \\
&= \frac{1}{2} \cdot \frac{1}{\mathbb{E}_\pi \left( \frac{p(X)}{q(X)} \right)} \cdot \sum_{x \in \text{supp}(\pi)} \pi(x) \left| \frac{p(x)}{q(x)} - \mathbb{E}_\pi \left( \frac{p(X)}{q(X)} \right) \right| \\
&= \frac{1}{2} \cdot \frac{1}{\mathbb{E}_\pi \left( \frac{p(X)}{q(X)} \right)} \cdot \text{dev}_\pi \left( \frac{p(X)}{q(X)} \right) \\
&= \frac{1}{2} \text{ndev}_\pi \left( \frac{p(X)}{q(X)} \right). \quad \square
\end{aligned}$$

## B.2. APPROXIMATE IMPORTANCE SAMPLING

THEOREM 5.3 (RESTATED). *Let*

$$\hat{\mu}' = \frac{A}{B} = \frac{\frac{1}{n} \sum_{i=1}^n f(X_i) w'(X_i)}{\frac{1}{n} \sum_{i=1}^n w'(X_i)}$$

*be the estimator produced by the approximate importance sampling procedure for the parameter  $\mathbb{E}_\pi(f(X))$ . Then,*

$$\frac{\mathbb{E}_p(A)}{\mathbb{E}_p(B)} = \mathbb{E}_\pi(f(X)) + \frac{\text{cov}_\pi \left( f(X), \frac{p(X)}{q(X)} \right)}{\mathbb{E}_\pi \left( \frac{p(X)}{q(X)} \right)}.$$

PROOF. Since  $X_1, \dots, X_n$  are i.i.d. random variables, it suffices to analyze the expectations of  $f(X)w'(X)$  and  $w'(X)$ , where  $X \sim p$ .

$$\begin{aligned}
\mathbb{E}_p(f(X)w'(X)) &= \sum_{x \in \mathcal{U}} p(x) f(x) \frac{\hat{\pi}(x)}{\hat{q}(x)} \\
&= \frac{Z_{\hat{\pi}}}{Z_{\hat{q}}} \sum_{x \in \mathcal{U}} \pi(x) f(x) \frac{p(x)}{q(x)} \\
&= \frac{Z_{\hat{\pi}}}{Z_{\hat{q}}} \mathbb{E}_\pi \left( f(X) \frac{p(X)}{q(X)} \right).
\end{aligned}$$

$$\begin{aligned}
\mathbb{E}_p(w'(X)) &= \sum_{x \in \mathcal{U}} p(x) \frac{\hat{\pi}(x)}{\hat{q}(x)} = \frac{Z_{\hat{\pi}}}{Z_{\hat{q}}} \sum_{x \in \mathcal{U}} \pi(x) \frac{p(x)}{q(x)} \\
&= \frac{Z_{\hat{\pi}}}{Z_{\hat{q}}} \mathbb{E}_\pi \left( \frac{p(X)}{q(X)} \right).
\end{aligned}$$

Therefore,

$$\begin{aligned}
 \frac{\mathbb{E}_p(A)}{\mathbb{E}_p(B)} &= \frac{\mathbb{E}_p(f(X)w'(X))}{\mathbb{E}_p(w'(X))} = \frac{\frac{Z_{\pi}}{Z_{\hat{q}}} \cdot \mathbb{E}_{\pi} \left( f(X) \frac{p(X)}{q(X)} \right)}{\frac{Z_{\pi}}{Z_{\hat{q}}} \cdot \mathbb{E}_{\pi} \left( \frac{p(X)}{q(X)} \right)} \\
 &= \frac{\text{cov}_{\pi} \left( f(X), \frac{p(X)}{q(X)} \right) + \mathbb{E}_{\pi}(f(X)) \mathbb{E}_{\pi} \left( \frac{p(X)}{q(X)} \right)}{\mathbb{E}_{\pi} \left( \frac{p(X)}{q(X)} \right)} \\
 &= \mathbb{E}_{\pi}(f(X)) + \frac{\text{cov}_{\pi} \left( f(X), \frac{p(X)}{q(X)} \right)}{\mathbb{E}_{\pi} \left( \frac{p(X)}{q(X)} \right)}. \quad \square
 \end{aligned}$$

### B.3. APPROXIMATE METROPOLIS–HASTINGS

**THEOREM 5.5 (RESTATED).** *Let  $P'_{\text{MH}}$  be the transition matrix of the approximate Metropolis–Hastings algorithm. Then,  $P'_{\text{MH}}$  forms an ergodic Markov Chain and its unique limit distribution is  $\pi'$ .*

**PROOF.** The transition matrix of approximate MH is:

$$P'_{\text{MH}}(x, y) = \begin{cases} P(x, y) r'_{\text{MH}}(x, y), & \text{if } x \neq y, \\ P(x, x) r'_{\text{MH}}(x, x) + 1 - \sum_{z \in \mathcal{U}} P(x, z) r'_{\text{MH}}(x, z), & \text{if } x = y, \end{cases}$$

where

$$r'_{\text{MH}}(x, y) = \min \left\{ \frac{\pi(y) q(x)}{\pi(x) q(y)}, 1 \right\}.$$

Since  $\pi(x) > 0$  and  $q(x) > 0$  for all  $x \in \mathcal{U}$ , then  $r'_{\text{MH}}(x, y) > 0$  for all  $x, y \in \mathcal{U}$ . It follows that whenever  $P(x, y) > 0$ , then also  $P'_{\text{MH}}(x, y) > 0$ . This implies that the Markov Chain graph  $G_{P'_{\text{MH}}}$  corresponding to  $P'_{\text{MH}}$  contains all the edges of the Markov Chain graph  $G_P$  corresponding to  $P$ . Since  $P$  is ergodic, we know that  $G_P$  is strongly connected and aperiodic. As strong connectedness and aperiodicity are invariant under addition of edges, it follows that also  $G_{P'_{\text{MH}}}$  must be strongly connected and aperiodic, and thus  $P'_{\text{MH}}$  is ergodic.

As  $P'_{\text{MH}}$  is ergodic, the fundamental theorem of Markov Chains tell us that it has a unique limit distribution  $\eta$ . Furthermore,  $\eta$  is the unique stationary distribution of  $P'_{\text{MH}}$ . We use the next proposition to show that that  $\eta = \pi'$ .

**PROPOSITION B.1.** *If  $P$  is ergodic, and reversible with respect to  $\pi$ , then  $\pi$  is the unique stationary distribution of  $P$ .*

**PROOF.** For every  $y$ , we have

$$[\pi P](y) = \sum_{x \in \mathcal{U}} \pi(x) P(x, y) = \sum_{x \in \mathcal{U}} \pi(y) P(y, x) = \pi(y).$$

Hence  $\pi$  is stationary, and by the fundamental theorem of Markov Chains it is unique.  $\square$

In order to show that  $P'_{\text{MH}}$  satisfies the conditions of Proposition B.1, we are left to show that  $P'_{\text{MH}}$  is reversible with respect to  $\pi'$ . That is,

$$\pi'(x)P'_{\text{MH}}(x, y) = \pi'(y)P'_{\text{MH}}(y, x).$$

It will then follow that  $\pi'$  is the unique limit distribution of  $P'_{\text{MH}}$ .

For  $x = y$ , reversibility follows trivially. When  $x \neq y$ ,

$$\begin{aligned} & \pi'(x)P'_{\text{MH}}(x, y) \\ &= \pi'(x)P(x, y)r'_{\text{MH}}(x, y) \\ &= \pi(x) \frac{\frac{p(x)}{q(x)}}{\mathbb{E}_{\pi}\left(\frac{p(X)}{q(X)}\right)} P(x, y) \min \left\{ \frac{\pi(y) q(x)}{\pi(x) q(y)}, 1 \right\} \\ &= \frac{1}{\mathbb{E}_{\pi}\left(\frac{p(X)}{q(X)}\right)} p(x) P(x, y) \frac{\pi(x)}{q(x)} \min \left\{ \frac{\pi(y) q(x)}{\pi(x) q(y)}, 1 \right\} \\ &= \frac{1}{\mathbb{E}_{\pi}\left(\frac{p(X)}{q(X)}\right)} p(y) P(y, x) \frac{\pi(x)}{q(x)} \min \left\{ \frac{\pi(y) q(x)}{\pi(x) q(y)}, 1 \right\} \quad (\text{Since } P \text{ is reversible}) \\ &= \frac{1}{\mathbb{E}_{\pi}\left(\frac{p(X)}{q(X)}\right)} p(y) P(y, x) \frac{\pi(y)}{q(y)} \min \left\{ \frac{\pi(y) q(x)}{\pi(x) q(y)} \frac{\pi(x) q(y)}{q(x) \pi(y)}, \frac{\pi(x) q(y)}{q(x) \pi(y)} \right\} \\ &= \pi(y) \frac{\frac{p(y)}{q(y)}}{\mathbb{E}_{\pi}\left(\frac{p(X)}{q(X)}\right)} P(y, x) \min \left\{ 1, \frac{\pi(x) q(y)}{\pi(y) q(x)} \right\} \\ &= \pi'(y)P(y, x)r'_{\text{MH}}(y, x) = \pi'(y)P'_{\text{MH}}(y, x). \end{aligned}$$

#### B. 4. APPROXIMATE MAXIMUM DEGREE

**THEOREM 5.6 (RESTATED).** *Let  $P'_{\text{MD}}$  be the transition matrix of the approximate Maximum Degree procedure. Then,  $P'_{\text{MD}}$  forms an ergodic Markov Chain. The unique limit distribution of  $P'_{\text{MD}}$  is  $\pi'$ .*

**PROOF.** The transition matrix of approximate MD is:

$$P'_{\text{MD}}(x, y) = \begin{cases} P(x, y) r'_{\text{MD}}(x), & \text{if } x \neq y, \\ P(x, x) r'_{\text{MD}}(x) + 1 - r'_{\text{MD}}(x), & \text{if } x = y. \end{cases}$$

Recall that

$$r'_{\text{MD}}(x) = \frac{\hat{q}(x)}{C \hat{\pi}(x)}, \quad \text{where } C \geq \max_{x \in \mathcal{U}} \frac{\hat{q}(x)}{\hat{\pi}(x)}.$$

Since  $\hat{q}(x) > 0$  for all  $x \in \mathcal{U}$ , it can be seen from the above expression that for all  $x \neq y \in \mathcal{U}$ ,  $P'_{\text{MD}}(x, y) > 0$  if and only if  $P(x, y) > 0$ . Furthermore, if  $P'_{\text{MD}}(x, x) = 0$ , then also  $P(x, x) = 0$  (the converse is not necessarily true). We conclude that the Markov Chain graph  $G_{P'_{\text{MD}}}$  corresponding to  $P'_{\text{MD}}$  contains all the edges of the Markov Chain graph  $G_P$  corresponding to  $P$  (the only possible extra edges in  $G_{P'_{\text{MD}}}$  are self loops). As  $P$  is ergodic, and ergodicity is invariant under addition of edges to the Markov Chain graph, we conclude that also  $P'_{\text{MD}}$  is ergodic.

We next prove that  $P'_{\text{MD}}$  is reversible with respect to  $\pi'$ . This would imply, by Proposition B.1, that  $\pi'$  is also the unique limit distribution of  $P'_{\text{MD}}$ .

For  $x = y$ , reversibility follows trivially. When  $x \neq y$ ,

$$\begin{aligned}
 \pi'(x)P'_{\text{MD}}(x, y) &= \pi'(x)P(x, y)r'_{\text{MD}}(x, y) \\
 &= \pi(x) \frac{\frac{p(x)}{q(x)}}{\mathbb{E}_{\pi}\left(\frac{p(X)}{q(X)}\right)} P(x, y) \frac{\hat{q}(x)}{C \hat{\pi}(x)} \\
 &= \frac{1}{\mathbb{E}_{\pi}\left(\frac{p(X)}{q(X)}\right)} p(x)P(x, y) \frac{Z_{\hat{q}}}{C Z_{\hat{\pi}}} \\
 &= \frac{1}{\mathbb{E}_{\pi}\left(\frac{p(X)}{q(X)}\right)} p(y)P(y, x) \frac{Z_{\hat{q}}}{C Z_{\hat{\pi}}} \quad (\text{Since } P \text{ is reversible}) \\
 &= \pi(y) \frac{\frac{p(y)}{q(y)}}{\mathbb{E}_{\pi}\left(\frac{p(X)}{q(X)}\right)} P(y, x) \frac{\hat{q}(y)}{C \hat{\pi}(y)} \\
 &= \pi'(y)P(y, x)r'_{\text{MD}}(y, x) = \pi'(y)P'_{\text{MD}}(y, x). \quad \square
 \end{aligned}$$

### C. Pool-Based Sampler—Proofs

#### C.1. ANALYSIS OF SAMPLING RECALL AND SAMPLING BIAS

THEOREM 7.5 (RESTATED). *The sampling recall of the PB sampler is equal to:*

$$\text{recall}_{\pi}(\mathcal{P}_+) = \pi(\mathcal{D}_{\mathcal{P}_+}).$$

*The sampling bias of the PB sampler is at most:*

$$\frac{1}{2} \text{ndev}_{\pi_{\mathcal{P}_+}}(\text{vdensity}_{\mathcal{P}}(X)),$$

where  $\pi_{\mathcal{P}_+}$  is the restriction of  $\pi$  to  $\mathcal{D}_{\mathcal{P}_+} \cap \text{supp}(\pi)$ .

PROOF. The analysis of the recall of the PB sampler is identical to the recall analysis shown in the proof of Proposition 7.2. We therefore omit the details.

Let  $\eta$  be the sampling distribution of the PB sampler. Similarly to the proof of Proposition 7.2, it can be shown that:

$$\text{supp}(\eta) = \mathcal{D}_{\mathcal{P}_+} \cap \text{supp}(\pi).$$

Hence, the restriction of  $\pi$  to  $\text{supp}(\eta)$ , that is,  $\pi_{\text{supp}(\eta)}$ , is exactly  $\pi_{\mathcal{P}_+}$ . The sampling bias of the PB sampler is therefore the distance  $\|\eta - \pi_{\mathcal{P}_+}\|$ .

We note that the scenario we face here is exactly the one captured by Theorem 5.1: we have a rejection sampling procedure whose trial distribution does not match the unnormalized trial weights. The target distribution in this case is  $\hat{\pi}_{\mathcal{P}_+}$  and the trial distribution is  $d_{\mathcal{P}_+}$ . Let us denote the distribution induced by the trial weights by  $q$ . We next calculate a closed form expression for  $q$ .

The support of  $q$  is the same as the support of  $d_{\mathcal{P}_+}$ , because only documents that are sampled from  $d_{\mathcal{P}_+}$  are assigned an unnormalized trial weight. Hence,  $\text{supp}(q) = \text{supp}(d_{\mathcal{P}_+}) = \mathcal{D}_{\mathcal{P}_+}$ . For each  $x \in \mathcal{D}_{\mathcal{P}_+}$  the unnormalized trial weight is  $\deg_{\mathcal{P}}(x)$ .

Therefore, the normalization constant is:  $Z_{\hat{q}} = \sum_{x \in \mathcal{D}_{\mathcal{P}_+}} \deg_{\mathcal{P}}(x) = \deg_{\mathcal{P}}(\mathcal{D}_{\mathcal{P}_+})$ .  
Therefore, for every document  $x \in \mathcal{D}_{\mathcal{P}_+}$  we have:

$$q(x) = \frac{\deg_{\mathcal{P}}(x)}{\deg_{\mathcal{P}}(\mathcal{D}_{\mathcal{P}_+})}.$$

Applying Proposition 5.2, we bound the bias of the PB sampler as follows:

$$\|\eta - \pi_{\mathcal{P}_+}\| \leq \frac{1}{2} \text{ndev}_{\pi_{\mathcal{P}_+}} \left( \frac{d_{\mathcal{P}_+}(X)}{q(X)} \right).$$

For each  $x \in \mathcal{D}_{\mathcal{P}_+}$ , we have:

$$\frac{d_{\mathcal{P}_+}(x)}{q(x)} = \frac{\frac{\deg_{\mathcal{P}_+}(x)}{\deg_{\mathcal{P}_+}(\mathcal{D}_{\mathcal{P}_+})}}{\frac{\deg_{\mathcal{P}}(x)}{\deg_{\mathcal{P}}(\mathcal{D}_{\mathcal{P}_+})}} = \text{vdensity}_{\mathcal{P}}(x) \frac{\deg_{\mathcal{P}}(\mathcal{D}_{\mathcal{P}_+})}{\deg_{\mathcal{P}_+}(\mathcal{D}_{\mathcal{P}_+})}.$$

Therefore,

$$\begin{aligned} \text{ndev}_{\pi_{\mathcal{P}_+}} \left( \frac{d_{\mathcal{P}_+}(X)}{q(X)} \right) &= \text{ndev}_{\pi_{\mathcal{P}_+}} \left( \text{vdensity}_{\mathcal{P}}(X) \cdot \frac{\deg_{\mathcal{P}}(\mathcal{D}_{\mathcal{P}_+})}{\deg_{\mathcal{P}_+}(\mathcal{D}_{\mathcal{P}_+})} \right) \\ &= \text{ndev}_{\pi_{\mathcal{P}_+}} (\text{vdensity}_{\mathcal{P}}(X)), \end{aligned}$$

where the latter equality follows from the fact the normalized mean deviation is invariant under multiplication by a scalar.  $\square$

**THEOREM 7.6 (RESTATED).** *The sampling bias of the PB sampler is at most:*

$$\frac{\text{ovprob}(w_{\mathcal{P}})}{1 - \text{ovprob}(w_{\mathcal{P}})}.$$

**PROOF.** In order to prove the theorem, we prove that

$$\frac{1}{2} \cdot \text{ndev}_{\pi_{\mathcal{P}_+}} (\text{vdensity}_{\mathcal{P}}(X)) \leq \frac{\text{ovprob}(w_{\mathcal{P}})}{1 - \text{ovprob}(w_{\mathcal{P}})}.$$

Let  $\mu = \mathbb{E}_{\pi_{\mathcal{P}_+}} (\text{vdensity}_{\mathcal{P}}(X))$ . Then,

$$\frac{1}{2} \cdot \text{ndev}_{\pi_{\mathcal{P}_+}} (\text{vdensity}_{\mathcal{P}}(X)) = \frac{\mathbb{E}_{\pi_{\mathcal{P}_+}} (|\text{vdensity}_{\mathcal{P}}(X) - \mu|)}{2\mu}.$$

For a document  $x \in \mathcal{D}_{\mathcal{P}_+}$ , we define the *invalidity density* of  $x$  as:

$$\text{invdensity}_{\mathcal{P}}(x) = 1 - \text{vdensity}_{\mathcal{P}}(x).$$

Then,

$$\frac{\mathbb{E}_{\pi_{\mathcal{P}_+}} (|\text{vdensity}_{\mathcal{P}}(X) - \mu|)}{2\mu} = \frac{\mathbb{E}_{\pi_{\mathcal{P}_+}} (|1 - \text{invdensity}_{\mathcal{P}}(X) - \mu|)}{2\mu}.$$

By the triangle inequality, for every  $x \in \mathcal{D}_{\mathcal{P}_+}$ ,

$$|1 - \text{invdensity}_{\mathcal{P}}(x) - \mu| \leq 1 - \mu + \text{invdensity}_{\mathcal{P}}(x) = 1 - \mu + 1 - \text{vdensity}_{\mathcal{P}}(x).$$



Hence,

$$\begin{aligned}
 \frac{\mathbb{E}_{\pi_{\mathcal{P}_+}}(|1 - \text{invdensity}_{\mathcal{P}}(X) - \mu|)}{2\mu} &\leq \frac{\mathbb{E}_{\pi_{\mathcal{P}_+}}(1 - \mu + 1 - \text{vdensity}_{\mathcal{P}}(X))}{2\mu} \\
 &= \frac{1 - \mu + 1 - \mathbb{E}_{\pi_{\mathcal{P}_+}}(\text{vdensity}_{\mathcal{P}}(X))}{2\mu} \\
 &= \frac{1 - \mu + 1 - \mu}{2\mu} = \frac{1 - \mu}{\mu}. \tag{3}
 \end{aligned}$$

Next, we relate  $\mu = \mathbb{E}_{\pi_{\mathcal{P}_+}}(\text{vdensity}_{\mathcal{P}}(X))$  to  $\text{ovprob}(w_{\mathcal{P}})$ .

$$\begin{aligned}
 \mathbb{E}_{\pi_{\mathcal{P}_+}}(\text{vdensity}_{\mathcal{P}}(X)) &= \sum_{x \in \mathcal{D}_{\mathcal{P}_+}} \pi_{\mathcal{P}_+}(x) \cdot \frac{\deg_{\mathcal{P}_+}(x)}{\deg_{\mathcal{P}}(x)} \\
 &= \sum_{x \in \mathcal{D}_{\mathcal{P}_+}} \frac{\pi_{\mathcal{P}_+}(x)}{\deg_{\mathcal{P}}(x)} \sum_{q \in \text{queries}_{\mathcal{P}_+}(x)} 1 \\
 &= \sum_{q \in \mathcal{P}_+} \sum_{x \in \text{results}(q)} \frac{\pi_{\mathcal{P}_+}(x)}{\deg_{\mathcal{P}}(x)} \\
 &= \sum_{q \in \mathcal{P}_+} w_{\mathcal{P}}(q) = w_{\mathcal{P}}(\mathcal{P}_+)
 \end{aligned}$$

Since underflowing queries in  $\mathcal{P}$  have zero weight, then

$$\begin{aligned}
 w_{\mathcal{P}}(\mathcal{P}_+) &= 1 - w_{\mathcal{P}}(\mathcal{P}_-) = 1 - \Pr(Q \in \mathcal{P}_- | w_{\mathcal{P}}) \\
 &= 1 - \Pr(\text{card}(Q) > k | w_{\mathcal{P}}) = 1 - \text{ovprob}(w_{\mathcal{P}}).
 \end{aligned}$$

Substituting back in Eq. (3), we have:

$$\frac{1}{2} \cdot \text{ndev}_{\pi_{\mathcal{P}_+}}(\text{vdensity}_{\mathcal{P}}(X)) \leq \frac{1 - (1 - \text{ovprob}(w_{\mathcal{P}}))}{1 - \text{ovprob}(w_{\mathcal{P}})} = \frac{\text{ovprob}(w_{\mathcal{P}})}{1 - \text{ovprob}(w_{\mathcal{P}})}. \quad \square$$

## C.2. ANALYSIS OF QUERY AND FETCH COSTS

**THEOREM 7.9 (RESTATED).** *The query cost of the PB sampler is at most:*

$$\frac{C \cdot \deg_{\mathcal{P}_+}(\mathcal{D}_{\mathcal{P}_+})}{Z_{\hat{\pi}} \cdot \pi(\mathcal{D}_{\mathcal{P}_+}) \cdot (1 - \text{ovprob}(w_{\mathcal{P}}))} \cdot \left( \frac{|\mathcal{P}|}{|\mathcal{P}_+|} \cdot \frac{k}{\text{avg}_{q \in \mathcal{P}_+} \text{card}(q)} + \text{qcost}(\hat{\pi}) \right).$$

*The fetch cost of the PB sampler is at most:*

$$\frac{C \cdot \deg_{\mathcal{P}_+}(\mathcal{D}_{\mathcal{P}_+})}{Z_{\hat{\pi}} \cdot \pi(\mathcal{D}_{\mathcal{P}_+}) \cdot (1 - \text{ovprob}(w_{\mathcal{P}}))} \cdot (1 + \text{fcost}(\hat{\pi})).$$

**PROOF.** Search engine queries are made in two of the subroutines of the PB sampler: (1) In DDSampler, which selects a random document from the results of the query  $Q$  that it gets from QCSampler; and (2) In QCSampler, in order to determine the cardinality of the selected random query  $Q$ . Queries can also be possibly made in the  $\text{getWeight}_{\hat{\pi}}(x)$  procedure, when calculating the target weight  $\hat{\pi}(x)$  of a document  $x$ . Note that DDSampler needs the results of a query  $Q$  only after  $Q$  has already been processed by QCSampler and thus has already been submitted to

the search engine. Therefore, by careful data management, we can keep the results already returned for  $Q$  in memory, and hence avoid the additional search engine queries in DDSampler. We assume, then, from now on that DDSampler does not submit queries to the search engine.

In order to analyze the total number of queries made by PBSampler, we define the following random variables: (1)  $T$  is the number of iterations made in the loop of the outer procedure. Note that  $T$  determines exactly the number of calls to the  $\text{getWeight}_{\hat{\pi}}(x)$  procedure as well as the number of calls to the QCSampler procedure. (2)  $S_i$  (for  $i = 1, 2, \dots$ ) is the number of iterations made in the loop of QCSampler during its  $i$ -th invocation. The total number of queries submitted to the search engine by the PB sampler is at most  $\sum_{i=1}^T S_i + T \cdot \text{qcost}(\hat{\pi})$ .

The query cost of PBSampler is then  $\mathbb{E}(\sum_{i=1}^T S_i) + \mathbb{E}(T) \cdot \text{qcost}(\hat{\pi})$ . In order to analyze the first term, we resort to Wald's identity. Note that the conditions of Wald's identity are met:  $S_1, S_2, \dots$  are i.i.d. random variables and the event  $\{T = i\}$  is independent of  $S_{i+1}, S_{i+2}, \dots$ . Hence, the query cost of the PB sampler is  $\mathbb{E}(T) \cdot \mathbb{E}(S) + \mathbb{E}(T) \cdot \text{qcost}(\hat{\pi}) = \mathbb{E}(T) \cdot (\mathbb{E}(S) + \text{qcost}(\hat{\pi}))$ , where  $S$  is the random number of iterations in a single invocation of QCSampler.

In order to bound  $\mathbb{E}(T)$ , we analyze the parameters of the rejection sampling applied at the outer procedure:

- (1) The target distribution is  $\pi_{\mathcal{P}_+}$ . As shown in the proof of Proposition 7.2, the corresponding normalization constant is:  $Z_{\hat{\pi}_{\mathcal{P}_+}} = Z_{\hat{\pi}} \cdot \pi(\mathcal{D}_{\mathcal{P}_+})$ .
- (2) The trial distribution is  $d_{\mathcal{P}_+}$ .
- (3) As shown in the proof of Theorem 7.5, the trial weight distribution is  $q$ . Its normalization constant is:  $Z_{\hat{q}} = \deg_{\mathcal{P}}(\mathcal{D}_{\mathcal{P}_+})$ .

As shown in the proof of Theorem 7.5,

$$\mathbb{E}_{\pi_{\mathcal{P}_+}} \left( \frac{d_{\mathcal{P}_+}(X)}{q(X)} \right) = \frac{\deg_{\mathcal{P}}(\mathcal{D}_{\mathcal{P}_+})}{\deg_{\mathcal{P}_+}(\mathcal{D}_{\mathcal{P}_+})} \cdot \mathbb{E}_{\pi_{\mathcal{P}_+}}(\text{vdensity}_{\mathcal{P}}(X)).$$

In Theorem 7.6, we showed:

$$\mathbb{E}_{\pi_{\mathcal{P}_+}}(\text{vdensity}_{\mathcal{P}}(X)) = 1 - \text{ovprob}(w_{\mathcal{P}}).$$

Therefore, applying Theorem 5.1, we have:

$$\begin{aligned} \mathbb{E}(T) &= \frac{C \cdot Z_{\hat{q}}}{Z_{\hat{\pi}_{\mathcal{P}_+}} \cdot \mathbb{E}_{\pi_{\mathcal{P}_+}} \left( \frac{d_{\mathcal{P}_+}(X)}{q(X)} \right)} \\ &= \frac{C \cdot \deg_{\mathcal{P}}(\mathcal{D}_{\mathcal{P}_+})}{Z_{\hat{\pi}} \cdot \pi(\mathcal{D}_{\mathcal{P}_+}) \cdot \frac{\deg_{\mathcal{P}}(\mathcal{D}_{\mathcal{P}_+})}{\deg_{\mathcal{P}_+}(\mathcal{D}_{\mathcal{P}_+})} \cdot (1 - \text{ovprob}(w_{\mathcal{P}}))} \\ &= \frac{C \cdot \deg_{\mathcal{P}_+}(\mathcal{D}_{\mathcal{P}_+})}{Z_{\hat{\pi}} \cdot \pi(\mathcal{D}_{\mathcal{P}_+}) \cdot (1 - \text{ovprob}(w_{\mathcal{P}}))}. \end{aligned} \tag{4}$$

Next, we bound the expected number of iterations made in the loop of QCSampler. To this end, we analyze the parameters of the rejection sampling applied at QCSampler:

- (1) The target distribution is  $c_{\mathcal{P}_+}$  and its normalization constant is:  $Z_{\hat{c}_{\mathcal{P}_+}} = \text{card}(\mathcal{P}_+)$ .

- (2) The trial distribution is the uniform distribution on  $\mathcal{P}$  and its normalization constant is:  $Z_{\hat{u}_P} = |\mathcal{P}|$ . The trial weight distribution is identical to the trial distribution in this case.
- (3) The envelope constant is  $C = k$ .

By the analysis of rejection sampling, we have:

$$\mathbb{E}(S) = \frac{k \cdot |\mathcal{P}|}{\text{card}(\mathcal{P}_+)} = \frac{|\mathcal{P}|}{|\mathcal{P}_+|} \cdot \frac{k}{\text{avg}_{q \in \mathcal{P}_+} \text{card}(q)}. \quad (5)$$

Combining the expressions derived at Eqs. (4) and (5), the total query cost of the PB sampler is at most:

$$\begin{aligned} & \mathbb{E}(T) \cdot (\mathbb{E}(S) + \text{qcost}(\hat{\pi})) \\ &= \frac{C \cdot \deg_{\mathcal{P}_+}(\mathcal{D}_{\mathcal{P}_+})}{Z_{\hat{\pi}} \cdot \pi(\mathcal{D}_{\mathcal{P}_+}) \cdot (1 - \text{ovprob}(w_{\mathcal{P}}))} \cdot \left( \frac{|\mathcal{P}|}{|\mathcal{P}_+|} \cdot \frac{k}{\text{avg}_{q \in \mathcal{P}_+} \text{card}(q)} + \text{qcost}(\hat{\pi}) \right). \end{aligned}$$

Next, we analyze the fetch cost of the PB sampler. Pages are fetched in two of the sampler's procedures: (1) In the outer procedure, in order to compute the degrees of documents; and (2) In the  $\text{getWeight}_{\hat{\pi}}(x)$  procedure, in order to compute the target weight of documents. The number of fetches is determined directly by the number of iterations in the outer procedure, which we denoted by  $T$ . Therefore, the fetch cost is at most  $\mathbb{E}(T) \cdot (1 + \text{fcost}(\hat{\pi}))$ . Using the analysis of  $\mathbb{E}(T)$  above, the fetch cost is at most:

$$\frac{C \cdot \deg_{\mathcal{P}_+}(\mathcal{D}_{\mathcal{P}_+})}{Z_{\hat{\pi}} \cdot \pi(\mathcal{D}_{\mathcal{P}_+}) \cdot (1 - \text{ovprob}(w_{\mathcal{P}}))} \cdot (1 + \text{fcost}(\hat{\pi})). \quad \square$$

#### D. Random-Walk-Based Sampler—Proofs

**THEOREM 8.3 (RESTATED).** *Let  $\varepsilon > 0$ . Suppose we run the MH sampler (respectively, the MD sampler) with a burn-in period  $B$  that guarantees the approximate MH Markov Chain (respectively, approximate MD Markov Chain) reaches a distribution, which is at distance of at most  $\varepsilon$  from the limit distribution. Then, the sampling bias of the MH sampler (respectively, MD sampler) is at most:*

$$\frac{1}{2} \text{ndev}_{\pi_F}(\text{vdensity}_{\mathcal{P}}(X)) + \varepsilon.$$

*The sampling recall of the MH sampler (respectively, MD sampler) is at least:*

$$\pi(F) \cdot (1 - \frac{1}{2} \text{ndev}_{\pi_F}(\text{vdensity}_{\mathcal{P}}(X)) - \varepsilon).$$

**PROOF.** By Theorems 5.6 and 5.5, the limit distributions of the MH and the MD samplers are equal to:

$$\eta(x) = \pi_F(x) \frac{\frac{d_F(x)}{q_F(x)}}{\mathbb{E}_{\pi_F} \left( \frac{d_F(X)}{q_F(X)} \right)}.$$

By Proposition 5.2, the distance between this limit distribution and the target distribution  $\pi_F$  is bounded as follows:

$$\|\eta - \pi_F\| \leq \frac{1}{2} \text{ndev}_{\pi_F} \left( \frac{d_F(X)}{q_F(X)} \right).$$

Now, for every  $x \in F$ ,

$$\frac{d_F(x)}{q_F(x)} = \frac{\frac{\deg_{\mathcal{P}_+}(x)}{\deg_{\mathcal{P}_+}(F)}}{\frac{\deg_{\mathcal{P}}(x)}{\deg_{\mathcal{P}}(F)}} = \text{vdensity}_{\mathcal{P}}(x) \cdot \frac{\deg_{\mathcal{P}}(F)}{\deg_{\mathcal{P}_+}(F)}.$$

Since normalized mean deviation is invariant under multiplication by scalars, we therefore have:

$$\|\eta - \pi_F\| \leq \frac{1}{2} \text{ndev}_{\pi_F}(\text{vdensity}_{\mathcal{P}}(X)).$$

Consider now the distribution  $\eta'$  obtained after running the approximate MH Markov Chain (respectively, MD Markov Chain) for  $B$  steps.  $\eta'$  is the distribution of samples generated by the procedure. Since  $B \geq T_\varepsilon(P'_{\text{MH}})$  (respectively,  $B \geq T_\varepsilon(P'_{\text{MD}})$ ), then

$$\|\eta' - \eta\| \leq \varepsilon.$$

Therefore, by the triangle inequality,

$$\|\eta' - \pi_F\| \leq \frac{1}{2} \text{ndev}_{\pi_F}(\text{vdensity}_{\mathcal{P}}(X)) + \varepsilon. \quad (6)$$

We next show that this gives also an upper bound on the sampling bias, that is, on  $\|\eta'_{\text{supp}(\eta')} - \pi_{\text{supp}(\eta')}\|$ .

$$\begin{aligned} \|\eta'_{\text{supp}(\eta')} - \pi_{\text{supp}(\eta')}\| &= \frac{1}{2} \sum_{x \in \text{supp}(\eta')} \left| \eta'(x) - \frac{\pi(x)}{\pi(\text{supp}(\eta'))} \right| \\ &\leq \frac{1}{2} \sum_{x \in \text{supp}(\eta')} |\eta'(x) - \pi_F(x)| + \frac{1}{2} \sum_{x \in \text{supp}(\eta')} \left| \pi_F(x) - \frac{\pi(x)}{\pi(\text{supp}(\eta'))} \right| \end{aligned} \quad (7)$$

Let us bound each of the two above sums separately. To bound the first one, we resort to Eq. (6):

$$\begin{aligned} \frac{1}{2} \sum_{x \in \text{supp}(\eta')} |\eta'(x) - \pi_F(x)| &= \|\eta' - \pi_F\| - \frac{1}{2} \sum_{x \in F \setminus \text{supp}(\eta')} \pi_F(x) \\ &= \|\eta' - \pi_F\| - \frac{1}{2}(1 - \pi_F(\text{supp}(\eta'))). \end{aligned} \quad (8)$$

As for the second sum,

$$\begin{aligned}
& \frac{1}{2} \sum_{x \in \text{supp}(\eta')} \left| \pi_F(x) - \frac{\pi(x)}{\pi(\text{supp}(\eta'))} \right| \\
&= \frac{1}{2} \sum_{x \in \text{supp}(\eta')} \left| \pi_F(x) - \frac{\pi_F(x)}{\pi_F(\text{supp}(\eta'))} \right| \\
&= \frac{1}{2} \left( \frac{1}{\pi_F(\text{supp}(\eta'))} - 1 \right) \sum_{x \in \text{supp}(\eta')} \pi_F(x) \\
&= \frac{1}{2} \left( \frac{1}{\pi_F(\text{supp}(\eta'))} - 1 \right) \cdot \pi_F(\text{supp}(\eta')) \\
&= \frac{1}{2} (1 - \pi_F(\text{supp}(\eta'))) \tag{9}
\end{aligned}$$

Combining Eqs. (6)–(9), we have:

$$\begin{aligned}
\|\eta'_{\text{supp}(\eta')} - \pi_{\text{supp}(\eta')}\| &\leq \|\eta' - \pi_F\| - \frac{1}{2}(1 - \pi_F(\text{supp}(\eta'))) + \frac{1}{2}(1 - \pi_F(\text{supp}(\eta'))) \\
&= \|\eta' - \pi_F\| \leq \frac{1}{2} \text{ndev}_{\pi_F}(\text{vdensity}_{\mathcal{P}}(\mathbf{X})) + \varepsilon.
\end{aligned}$$

This gives us the desired upper bound on the sampling bias. We now turn to bounding the recall of the MH and MD samplers. Using Eq. (6) and the characterization of total variation distance given in Lemma 3.2, we have:

$$|\eta'(\text{supp}(\eta')) - \pi_F(\text{supp}(\eta'))| \leq \frac{1}{2} \text{ndev}_{\pi_F}(\text{vdensity}_{\mathcal{P}}(\mathbf{X})) + \varepsilon.$$

Since  $\eta'(\text{supp}(\eta')) = 1$ , then this implies:

$$\pi(\text{supp}(\eta')) = \pi(F) \cdot \pi_F(\text{supp}(\eta')) \geq \pi(F) \cdot (1 - \frac{1}{2} \text{ndev}_{\pi_F}(\text{vdensity}_{\mathcal{P}}(\mathbf{X})) - \varepsilon).$$

This gives us the lower bound on the recall of the MH and MD samplers.  $\square$

**THEOREM 8.6 (RESTATED).** *For any  $t \geq 0$ , the expected query cost of the  $t$ th step of the random walk sampler (whether MH or MD) is at most:*

$$\frac{1}{\mathbb{E}_{\pi_F}(\text{vdensity}_{\mathcal{P}}(\mathbf{X}))} + 3\varepsilon_t \sqrt{|F| \cdot \mathbb{E}_{u_F} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}^2(\mathbf{X})} \right)} + \text{qcost}(\hat{\pi}).$$

Here,  $\varepsilon_t$  is an upper bound on the total variation distance between the distribution of the node visited at the  $t$ th step and the limit distribution  $\eta$  of the random walk.  $u_F$  is the uniform distribution over  $F$ .

**PROOF.** Let  $\eta_t$  be the distribution of the node visited at the  $t$ th step of the random walk. By Proposition 8.4, the expected query cost of the  $t$ th step of the random walk is:

$$\mathbb{E}_{\eta_t} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}(\mathbf{X})} \right) + \text{qcost}(\hat{\pi}).$$

As the term  $\text{qcost}(\hat{\pi})$  is independent of  $X$ , we will focus from now on bounding the expected inverse validity density.

Our strategy for bounding  $\mathbb{E}_{\eta_t}(1/\text{vdensity}_{\mathcal{P}}(X))$  will be the following. We will first show an upper bound on  $\mathbb{E}_{\eta}(1/\text{vdensity}_{\mathcal{P}}(X))$ , where  $\eta$  is the limit distribution of the random walk. We will then bound the difference between  $\mathbb{E}_{\eta}(1/\text{vdensity}_{\mathcal{P}}(X))$  and  $\mathbb{E}_{\eta_t}(1/\text{vdensity}_{\mathcal{P}}(X))$ .

The following proposition analyzes the expected inverse validity density relative to the limit distribution:

PROPOSITION D.1

$$\mathbb{E}_{\eta} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}(X)} \right) = \frac{1}{\mathbb{E}_{\pi_F} \text{vdensity}_{\mathcal{P}}(X)}.$$

PROOF. Recall (Eq. (2)) that for both the MH and MD random walks,

$$\eta(x) = \pi_F(x) \cdot \frac{\frac{d_F(x)}{q_F(x)}}{\mathbb{E}_{\pi_F} \left( \frac{d_F(X)}{q_F(X)} \right)}.$$

Here,  $F$  is the connected component on which the random walk runs,  $\pi_F$  is the restriction of the target distribution  $\pi$  to  $F$ ,  $d_F$  is the restriction of the degree distribution  $d_{\mathcal{P}_+}$  to  $F$ , and  $q_F$  is the restriction of the approximate degree distribution  $d_{\mathcal{P}}$  to  $F$ . Therefore,

$$\begin{aligned} \mathbb{E}_{\eta} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}(X)} \right) &= \sum_{x \in F} \eta(x) \cdot \frac{1}{\text{vdensity}_{\mathcal{P}}(x)} \\ &= \sum_{x \in F} \pi_F(x) \cdot \frac{\frac{d_F(x)}{q_F(x)}}{\mathbb{E}_{\pi_F} \left( \frac{d_F(X)}{q_F(X)} \right)} \cdot \frac{1}{\text{vdensity}_{\mathcal{P}}(x)} \\ &= \sum_{x \in F} \pi_F(x) \cdot \frac{\text{vdensity}_{\mathcal{P}}(x)}{\mathbb{E}_{\pi_F}(\text{vdensity}_{\mathcal{P}}(X))} \cdot \frac{1}{\text{vdensity}_{\mathcal{P}}(x)} \\ &= \frac{1}{\mathbb{E}_{\pi_F}(\text{vdensity}_{\mathcal{P}}(X))} \sum_{x \in F} \pi_F(x) \\ &= \frac{1}{\mathbb{E}_{\pi_F}(\text{vdensity}_{\mathcal{P}}(X))}. \quad \square \end{aligned}$$

In order to bound the difference

$$\left| \mathbb{E}_{\eta} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}(x)} \right) - \mathbb{E}_{\eta_t} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}(x)} \right) \right|,$$

we will resort to the following general lemma that bounds the difference between expectations of the same function of two random variables in terms of the distance between their distributions.

LEMMA D.2. *Let  $\rho$  and  $\rho'$  be two distributions on the same finite probability space  $\mathcal{U}$  and let  $f : \mathcal{U} \rightarrow \mathbb{R}$  be some real-valued function. Then,*

$$|\mathbb{E}_{\rho'}(f(X)) - \mathbb{E}_{\rho}(f(X))| \leq 2\varepsilon \cdot \mathbb{E}_u(|f(X)|) + |\mathcal{U}| \cdot \text{cov}_u(|f(X)|, |\rho'(X) - \rho(X)|),$$

where  $\varepsilon \geq \|\rho - \rho'\|$  and  $u$  is the uniform distribution on  $\mathcal{U}$ .



PROOF

$$\begin{aligned}
& |\mathbb{E}_{\rho'}(f(\mathbf{X})) - \mathbb{E}_{\rho}(f(\mathbf{X}))| \\
&= \left| \sum_{\mathbf{x} \in \mathcal{U}} \rho'(\mathbf{x})f(\mathbf{x}) - \sum_{\mathbf{x} \in \mathcal{U}} \rho(\mathbf{x})f(\mathbf{x}) \right| \\
&= \left| \sum_{\mathbf{x} \in \mathcal{U}} f(\mathbf{x})(\rho'(\mathbf{x}) - \rho(\mathbf{x})) \right| \\
&\leq |\mathcal{U}| \cdot \sum_{\mathbf{x} \in \mathcal{U}} \frac{1}{|\mathcal{U}|} |f(\mathbf{x})| \cdot |\rho'(\mathbf{x}) - \rho(\mathbf{x})| \\
&= |\mathcal{U}| \cdot \mathbb{E}_u(|f(\mathbf{X})| \cdot |\rho'(\mathbf{X}) - \rho(\mathbf{X})|) \\
&= |\mathcal{U}| \cdot (\mathbb{E}_u(|f(\mathbf{X})|) \cdot \mathbb{E}_u(|\rho'(\mathbf{X}) - \rho(\mathbf{X})|) + \text{cov}_u(|f(\mathbf{X})|, |\rho'(\mathbf{X}) - \rho(\mathbf{X})|)).
\end{aligned}$$

It is easy to see that  $|\mathcal{U}| \cdot \mathbb{E}_u(|\rho'(\mathbf{X}) - \rho(\mathbf{X})|) = 2\|\rho' - \rho\| \leq 2\varepsilon$ . So,

$$\begin{aligned}
|\mathbb{E}_{\rho'}(f(\mathbf{X})) - \mathbb{E}_{\rho}(f(\mathbf{X}))| &\leq 2\varepsilon \cdot \mathbb{E}_u(|f(\mathbf{X})|) \\
&\quad + |\mathcal{U}| \cdot \text{cov}_u(|f(\mathbf{X})|, |\rho'(\mathbf{X}) - \rho(\mathbf{X})|). \quad \square
\end{aligned}$$

We conclude from Lemma D.2 that:

$$\begin{aligned}
& \left| \mathbb{E}_{\eta} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}(\mathbf{x})} \right) - \mathbb{E}_{\eta_t} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}(\mathbf{x})} \right) \right| \\
&\leq 2\varepsilon_t \cdot \mathbb{E}_{u_F} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}(\mathbf{X})} \right) \\
&\quad + |F| \cdot \text{cov}_{u_F} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}(\mathbf{X})}, |\eta_t(\mathbf{X}) - \eta(\mathbf{X})| \right). \quad (10)
\end{aligned}$$

Recall that we use the upper bound  $\varepsilon_t = (1 - \alpha)^t / \eta_{\min}$  on  $\|\eta_t - \eta\|$ . We next show that also the covariance term can be bounded in terms of  $\varepsilon_t$ .

**LEMMA D.3.** *Consider a Markov Chain on a finite space  $\mathcal{U}$ . Let  $\alpha$  be the spectral gap of this Markov Chain's transition matrix, let  $\eta$  be its limit distribution, and let  $\eta_t$  be the distribution of the node visited at the  $t$ th step. Let  $f$  be any nonnegative real-valued function  $f : \mathcal{U} \rightarrow \mathbb{R}^+$ . Then,*

$$\text{cov}_u(f(\mathbf{X}), |\eta_t(\mathbf{X}) - \eta(\mathbf{X})|) \leq \frac{(1 - \alpha)^t}{\sqrt{\eta_{\min}}} \cdot \sqrt{\mathbb{E}_u(f^2(\mathbf{X}))},$$

where  $u$  is the uniform distribution on  $\mathcal{U}$ .

PROOF

$$\begin{aligned}
& \text{cov}_u(f(\mathbf{X}), |\eta_t(\mathbf{X}) - \eta(\mathbf{X})|) \\
&= \mathbb{E}_u(|\eta_t(\mathbf{X}) - \eta(\mathbf{X})| \cdot f(\mathbf{X})) - \mathbb{E}_u(|\eta_t(\mathbf{X}) - \eta(\mathbf{X})|) \cdot \mathbb{E}_u(f(\mathbf{X})).
\end{aligned}$$

In accordance with Theorem 4.8, for every  $\mathbf{x} \in \mathcal{U}$ ,

$$|\eta_t(\mathbf{x}) - \eta(\mathbf{x})| \leq \sqrt{\frac{\eta(\mathbf{x})}{\eta_{\min}}} (1 - \alpha)^t.$$

Thus,

$$\begin{aligned}
& \text{cov}_u(f(\mathbf{X}), |\eta_t(\mathbf{X}) - \eta(\mathbf{X})|) \\
& \leq \frac{(1 - \alpha)^t}{\sqrt{\eta_{\min}}} \cdot \left( \mathbb{E}_u(\sqrt{\eta(\mathbf{X})} \cdot f(\mathbf{X})) - \mathbb{E}_u(\sqrt{\eta(\mathbf{X})}) \cdot \mathbb{E}_u(f(\mathbf{X})) \right) \\
& \leq \frac{(1 - \alpha)^t}{\sqrt{\eta_{\min}}} \cdot \mathbb{E}_u(\sqrt{\eta(\mathbf{X})} \cdot f(\mathbf{X})) \quad (\text{Since } f(\mathbf{x}) \geq 0, \forall \mathbf{x}) \\
& \leq \frac{(1 - \alpha)^t}{\sqrt{\eta_{\min}}} \cdot \sqrt{\mathbb{E}_u(\eta(\mathbf{X}))} \cdot \sqrt{\mathbb{E}_u(f^2(\mathbf{X}))} \quad (\text{Cauchy-Schwartz}) \\
& = \frac{(1 - \alpha)^t}{\sqrt{\eta_{\min}}} \cdot \sqrt{\mathbb{E}_u(f^2(\mathbf{X}))} \quad (\text{Since } \eta \text{ is a distribution}). \quad \square
\end{aligned}$$

We conclude from Lemma D.3 that:

$$\begin{aligned}
& |F| \cdot \text{cov}_{u_F} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}(\mathbf{X})}, |\eta_t(\mathbf{X}) - \eta(\mathbf{X})| \right) \\
& \leq |F| \cdot \frac{(1 - \alpha)^t}{\sqrt{\eta_{\min}}} \cdot \sqrt{\mathbb{E}_{u_F} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}^2(\mathbf{X})} \right)} \\
& \leq \sqrt{|F|} \cdot \frac{(1 - \alpha)^t}{\eta_{\min}} \cdot \sqrt{\mathbb{E}_{u_F} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}^2(\mathbf{X})} \right)} \quad (\text{Since } \eta_{\min} \leq 1/|F|) \\
& = \varepsilon_t \cdot \sqrt{|F| \cdot \mathbb{E}_{u_F} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}^2(\mathbf{X})} \right)}. \tag{11}
\end{aligned}$$

Combining Proposition D.1, Eq. (10), and Eq. (11), we have:

$$\begin{aligned}
& \mathbb{E}_{\eta_t} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}(\mathbf{X})} \right) \\
& \leq \mathbb{E}_{\eta} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}(\mathbf{X})} \right) + \left| \mathbb{E}_{\eta_t} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}(\mathbf{X})} \right) - \mathbb{E}_{\eta} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}(\mathbf{X})} \right) \right| \\
& \leq \frac{1}{\mathbb{E}_{\pi_F}(\text{vdensity}_{\mathcal{P}}(\mathbf{X}))} + 2\varepsilon_t \cdot \mathbb{E}_{u_F} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}(\mathbf{X})} \right) \\
& \quad + \varepsilon_t \cdot \sqrt{|F| \cdot \mathbb{E}_{u_F} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}^2(\mathbf{X})} \right)} \\
& \leq \frac{1}{\mathbb{E}_{\pi_F}(\text{vdensity}_{\mathcal{P}}(\mathbf{X}))} + 3\varepsilon_t \sqrt{|F| \cdot \mathbb{E}_{u_F} \left( \frac{1}{\text{vdensity}_{\mathcal{P}}^2(\mathbf{X})} \right)} \\
& \quad (\text{Since } \mathbb{E}(Z) \leq \sqrt{\mathbb{E}(Z^2)}).
\end{aligned}$$

**THEOREM 8.9 (RESTATED).** *For any  $t \geq 0$ , the expected delay of the  $t$ th real step of the optimized MD random walk is at least:*

$$C' \cdot Z_{\hat{\pi}_F} \cdot \pi(F) \cdot \left( \frac{\mathbb{E}_{\pi_F}(\text{vdensity}_{\mathcal{P}}(\mathbf{X}))}{\deg_{\mathcal{P}_+}(F)} - 3\varepsilon_t \cdot \sqrt{\mathbb{E}_{u_F} \left( \frac{\pi_F^2(\mathbf{X})}{\deg_{\mathcal{P}}^2(\mathbf{X})} \right)} \right).$$

Here,  $\varepsilon_t$  is an upper bound on the total variation distance between the distribution of the node visited at the  $t$ th real step and the degree distribution  $d_F$ .  $u_F$  is the uniform distribution over  $F$ .

PROOF. Let  $d_t$  be the distribution of the node visited at the  $t$ th real step of the random walk. By Proposition 8.8, and since  $\text{supp}(d_t) \subseteq F$ , the expected delay of the  $t$ th step is:

$$\mathbb{E}_{d_t} \left( \frac{C' \hat{\pi}(x)}{\deg_{\mathcal{P}}(x)} \right) = C' \cdot Z_{\hat{\pi}_F} \cdot \pi(F) \cdot \mathbb{E}_{d_t} \left( \frac{\pi_F(x)}{\deg_{\mathcal{P}}(x)} \right). \quad (12)$$

The structure of this proof is similar to that of Theorem 8.6, with the difference that here we are looking for a lower bound instead of an upper bound. We will first show a lower bound on  $\mathbb{E}_{d_F}(\frac{\pi_F(x)}{\deg_{\mathcal{P}}(x)})$ , where  $d_F$  is the limit distribution of the simple random walk. We will then bound the difference between  $\mathbb{E}_{d_F}(\frac{\pi_F(x)}{\deg_{\mathcal{P}}(x)})$  and  $\mathbb{E}_{d_t}(\frac{\pi_F(x)}{\deg_{\mathcal{P}}(x)})$ .

The following proposition analyzes the expected delay relative to the limit distribution:

PROPOSITION D.4

$$\mathbb{E}_{d_F} \left( \frac{\pi_F(X)}{\deg_{\mathcal{P}}(X)} \right) = \frac{\mathbb{E}_{\pi_F}(\text{vdensity}_{\mathcal{P}}(X))}{\deg_{\mathcal{P}_+}(F)}.$$

PROOF. Recall (Eq. (1)) that the limit distribution of the real steps of the random walk is the degree distribution restricted to  $F$ :

$$d_F(x) = \frac{d_{\mathcal{P}_+}(x)}{d_{\mathcal{P}_+}(F)} = \frac{\deg_{\mathcal{P}_+}(x)}{\deg_{\mathcal{P}_+}(F)}.$$

Therefore,

$$\begin{aligned} \mathbb{E}_{d_F} \left( \frac{\pi_F(X)}{\deg_{\mathcal{P}}(X)} \right) &= \sum_{x \in F} \frac{\deg_{\mathcal{P}_+}(x)}{\deg_{\mathcal{P}_+}(F)} \cdot \frac{\pi_F(x)}{\deg_{\mathcal{P}}(x)} \\ &= \frac{1}{\deg_{\mathcal{P}_+}(F)} \sum_{x \in F} \pi_F(x) \cdot \frac{\deg_{\mathcal{P}_+}(x)}{\deg_{\mathcal{P}}(x)} \\ &= \frac{\mathbb{E}_{\pi_F}(\text{vdensity}_{\mathcal{P}}(X))}{\deg_{\mathcal{P}_+}(F)} \quad \square \end{aligned}$$

By Lemma D.2,

$$\begin{aligned} &\left| \mathbb{E}_{d_F} \left( \frac{\pi_F(x)}{\deg_{\mathcal{P}}(x)} \right) - \mathbb{E}_{d_t} \left( \frac{\pi_F(x)}{\deg_{\mathcal{P}}(x)} \right) \right| \\ &\leq 2\varepsilon_t \cdot \mathbb{E}_{u_F} \left( \frac{\pi_F(x)}{\deg_{\mathcal{P}}(x)} \right) + |F| \cdot \text{cov}_{u_F} \left( \frac{\pi_F(x)}{\deg_{\mathcal{P}}(x)}, |d_t(X) - d_F(X)| \right). \quad (13) \end{aligned}$$

Recall that we use the upper bound  $\varepsilon_t = (1 - \alpha)^t / d_{F \min}$  on  $\|d_t - d_F\|$ . We next show that also the covariance term can be bounded in terms of  $\varepsilon_t$ .

From Lemma D.3 we conclude that:

$$\begin{aligned}
& |F| \cdot \text{cov}_{u_F} \left( \frac{\pi_F(x)}{\deg_P(x)}, |d_t(X) - d_F(X)| \right) \\
& \leq |F| \cdot \frac{(1-\alpha)^t}{\sqrt{d_{F\min}}} \cdot \sqrt{\mathbb{E}_{u_F} \left( \frac{\pi_F^2(x)}{\deg_P^2(x)}, |d_t(X) - d_F(X)| \right)} \\
& \leq \sqrt{|F|} \cdot \frac{(1-\alpha)^t}{d_{F\min}} \cdot \sqrt{\mathbb{E}_{u_F} \left( \frac{\pi_F^2(x)}{\deg_P^2(x)} \right)} \quad (\text{Since } d_{F\min} \leq 1/|F|) \\
& = \varepsilon_t \cdot \sqrt{|F| \cdot \mathbb{E}_{u_F} \left( \frac{\pi_F^2(x)}{\deg_P^2(x)} \right)}. \tag{14}
\end{aligned}$$

Combining Proposition D.1, Eq. (13), and Eq. (14), we have:

$$\begin{aligned}
& \mathbb{E}_{d_t} \left( \frac{\pi_F(x)}{\deg_P(x)} \right) \\
& \geq \mathbb{E}_{d_F} \left( \frac{\pi_F(x)}{\deg_P(x)} \right) - \left| \mathbb{E}_{d_t} \left( \frac{\pi_F(x)}{\deg_P(x)} \right) - \mathbb{E}_{d_F} \left( \frac{\pi_F(x)}{\deg_P(x)} \right) \right| \\
& \geq \frac{\mathbb{E}_{\pi_F}(\text{vdensity}_P(X))}{\deg_{P_+}(F)} - 2\varepsilon_t \cdot \mathbb{E}_{u_F} \left( \frac{\pi_F(x)}{\deg_P(x)} \right) - \varepsilon_t \cdot \sqrt{|F| \cdot \mathbb{E}_{u_F} \left( \frac{\pi_F^2(x)}{\deg_P^2(x)} \right)} \\
& \geq \frac{\mathbb{E}_{\pi_F}(\text{vdensity}_P(X))}{\deg_{P_+}(F)} - 3\varepsilon_t \sqrt{|F| \cdot \mathbb{E}_{u_F} \left( \frac{\pi_F^2(x)}{\deg_P^2(x)} \right)} \quad (\text{Since } \mathbb{E}(Z) \leq \sqrt{\mathbb{E}(Z^2)}).
\end{aligned}$$

Substituting the above result into Eq. (12) we obtain the theorem.

**ACKNOWLEDGMENTS.** We thank Sara Cohen, Idit Keidar, Moshe Sidi, and the anonymous reviewers of WWW2006 and JACM for their valuable suggestions that helped to improve our article.

## REFERENCES

- ALDOUS, D. 1987. On the Markov chain simulation method for uniform combinatorial distributed and simulated annealing. *Prob. Eng. Inf. Sci.* 1, 33–46.
- ANAGNOSTOPOULOS, A., BRODER, A., AND CARMEL, D. 2005. Sampling search-engine results. In *Proceedings of the 14th International World Wide Web Conference (WWW)*. 245–256.
- BAR-YOSSEF, Z., BERG, A., CHIEN, S., FAKCHAROENPHOL, J., AND WEITZ, D. 2000. Approximating aggregate queries about Web pages via random walks. In *Proceedings of the 26th International Conference on Very Large Databases (VLDB)*. ACM, New York, 535–544.
- BAR-YOSSEF, Z., AND GUREVICH, M. 2007. Efficient search engine measurements. In *Proceedings of the 16th International World Wide Web Conference (WWW)*. 401–410.
- BATTELLE, J. 2005. John Battelle’s searchblog. <http://battellemedia.com/archives/001889.php>.
- BHARAT, K., AND BRODER, A. 1998. A technique for measuring the relative size and overlap of public Web search engines. In *Proceedings of the 7th International World Wide Web Conference (WWW7)*. 379–388.
- BONDAR, S. 2006. Search engine indexing limits: Where do the bots stop? (Available at <http://www.sitepoint.com/article/indexing-limits-where-bots-stop>.)

- BOYD, S., DIACONIS, P., AND XIAO, L. 2004. Fastest mixing Markov chain on a graph. *SIAM Rev.* 46, 4, 667–689.
- BRADLOW, E. T., AND SCHMITTEIN, D. C. 2000. The little engines that could: Modeling the performance of World Wide Web search engines. *Market. Sci.* 19, 43–62.
- BRODER, A., FONTOURA, M., JOSIFOVSKI, V., KUMAR, R., MOTWANI, R., NABAR, S., PANIGRAHY, R., TOMKINS, A., AND XU, Y. 2006. Estimating corpus size via queries. In *Proceedings of the 15th ACM International Conference on Information and Knowledge Management* (Arlington, VA). ACM, New York, 594–603.
- CAN, F., NURAY, R., AND SEVDIK, A. B. 2004. Automatic performance evaluation of Web search engines. *Inf. Proc. Manage.* 40, 495–514.
- CHENEY, M., AND PERRY, M. 2005. A comparison of the size of the Yahoo! and Google indices. Available at <http://vburton.ncsa.uiuc.edu/indexsize.html>.
- DAVISON, B. D. 2004. The potential of the metasearch engine. In *Proceedings of the Annual Meeting of the American Society for Information Science and Technology (ASIST)*. Vol. 41. 393–402.
- DIACONIS, P., AND SALOFF-COSTE, L. 1998. What do we know about the Metropolis algorithm? *J. Comput. Syst. Sci.* 57, 20–36.
- DOBRA, A., AND FIENBERG, S. E. 2004. How large is the World Wide Web? *Web Dyn.* 23–44.
- GILLMAN, D. 1998. A Chernoff bound for random walks on expander graphs. *SIAM J. Comput.* 27, 4, 1203–1220.
- GOLUB, G. H., AND VAN LOAN, C. F. 1996. *Matrix Computations*, 3rd ed. The Johns Hopkins University Press, Baltimore, MD.
- GORDON, M., AND PATHAK, P. 1999. Finding information on the World Wide Web: the retrieval effectiveness of search engines. *Inf. Proc. Manage.* 35, 2, 141–180.
- GULLI, A., AND SIGNORINI, A. 2005. The indexable Web is more than 11.5 billion pages. In *Proceedings of the 14th International Conference on World Wide Web (WWW)—Special Interest Tracks and Posters*. 902–903.
- HASTINGS, W. K. 1970. Monte Carlo sampling methods using Markov chains and their applications. *Biometrika* 57, 1, 97–109.
- HAWKING, D., CRASWEL, N., BAILEY, P., AND GRIFFITHS, K. 2001. Measuring search engine quality. *Inf. Retrieval* 4, 1, 33–59.
- HENZINGER, M. R., HEYDON, A., MITZENMACHER, M., AND NAJORK, M. 1999. Measuring index quality using random walks on the Web. In *Proceedings of the 8th International World Wide Web Conference*. 213–225.
- HENZINGER, M. R., HEYDON, A., MITZENMACHER, M., AND NAJORK, M. 2000. On near-uniform URL sampling. In *Proceedings of the 9th International World Wide Web Conference*. 295–308.
- HESTERBERG, T. C. 1988. Advances in importance sampling. Ph.D. dissertation, Stanford University.
- KAHALE, N. 1997. Large deviation for Markov chains. *Combinatorics, Probability and Computing* 6, 465–474.
- LAWRENCE, S., AND GILES, C. L. 1998. Searching the World Wide Web. *Science* 5360, 280, 98.
- LAWRENCE, S., AND GILES, C. L. 1999. Accessibility of information on the Web. *Nature* 400, 107–109.
- LIU, J. S. 1996. Metropolized independent sampling with comparisons to rejection sampling and importance sampling. *Stat. Comput.* 6, 113–119.
- LIU, J. S. 2001. *Monte Carlo Strategies in Scientific Computing*. Springer-Verlag, New York.
- MARSHAL, A. 1956. The use of multi-stage sampling schemes in Monte Carlo computations. In *Proceedings of the Symposium on Monte Carlo Methods*, M. Meyer, Ed. Vol. 21. Wiley, New York, 123–140.
- MAYER, T. 2005. Our blog is growing up - and so has our index. Available at <http://www.ysearchblog.com/archives/000172.html>.
- MCCOWN, F., AND NELSON, M. L. 2007. Agreeing to disagree: Search engines and their public interfaces. In *JCDL '07: Proceedings of the 2007 Conference on Digital Libraries*. 309–318.
- METROPOLIS, N., ROSENBLUTH, A., ROSENBLUTH, M., TELLER, A., AND TELLER, E. 1953. Equations of state calculations by fast computing machines. *J. Chem. Phys.* 21, 1087–1091.
- PRICE, G. 2005. More on the total database size battle and Googlehacking with Yahoo. Available at <http://blog.searchenginewatch.com/blog/050811-231448>.
- RUSMEVICHIENTONG, P., PENNOCK, D., LAWRENCE, S., AND GILES, C. L. 2001. Methods for sampling pages uniformly from the World Wide Web. In *Proceedings of AAAI Fall Symposium on Using Uncertainty within Computation*.

- SIEGMUND, D. 1985. *Sequential Analysis - Tests and Confidence Intervals*. Springer-Verlag, New York.
- SINCLAIR, A. 1993. *Algorithms for Random Generation and Counting: A Markov Chain Approach*. Birkhauser Verlag, Basel, Switzerland.
- STUART, A. and Ord, J.-K. 1994. *Kendall's Advanced Theory of Statistics*. Vol. 1, Distribution Theory. Haddon Arnold, London, UK.
- SULLIVAN, D. 2005. Search engine sizes. <http://searchenginewatch.com/showPage.html?page=2156481>.
- VÉRONIS, J. 2005. Yahoo: Missing pages? (2). <http://aixtal.blogspot.com/2005/08/yahoo-missing-pages-2.html>.
- VON NEUMANN, J. 1963. Various techniques used in connection with random digits. In *John von Neumann, Collected Works*. Vol. V. Oxford.

RECEIVED AUGUST 2006; REVISED MARCH 2008; ACCEPTED AUGUST 2008