

# Scaling Personalized Web Search\*

Glen Jeh  
Stanford University  
glenj@db.stanford.edu

Jennifer Widom  
Stanford University  
widom@db.stanford.edu

## ABSTRACT

Recent web search techniques augment traditional text matching with a global notion of “importance” based on the linkage structure of the web, such as in Google’s *PageRank* algorithm. For more refined searches, this global notion of importance can be specialized to create personalized views of importance—for example, importance scores can be biased according to a user-specified set of initially-interesting pages. Computing and storing all possible personalized views in advance is impractical, as is computing personalized views at query time, since the computation of each view requires an iterative computation over the web graph. We present new graph-theoretical results, and a new technique based on these results, that encode personalized views as *partial vectors*. Partial vectors are shared across multiple personalized views, and their computation and storage costs scale well with the number of views. Our approach enables incremental computation, so that the construction of personalized views from partial vectors is practical at query time. We present efficient dynamic programming algorithms for computing partial vectors, an algorithm for constructing personalized views from partial vectors, and experimental results demonstrating the effectiveness and scalability of our techniques.

## Categories and Subject Descriptors

G.2.2 [Discrete Mathematics]: Graph Theory

## General Terms

Algorithms

## Keywords

web search, PageRank

## 1. INTRODUCTION AND MOTIVATION

General web search is performed predominantly through text queries to search engines. Because of the enormous size of the web, text alone is usually not selective enough to limit the number of query results to a manageable size. The *PageRank* algorithm [11], among others [9], has been proposed (and implemented in *Google* [1]) to exploit the linkage structure of the web to compute

global “importance” scores that can be used to influence the ranking of search results. To encompass different notions of importance for different users and queries, the basic *PageRank* algorithm can be modified to create “personalized views” of the web, redefining importance according to user preference. For example, a user may wish to specify his bookmarks as a set of preferred pages, so that any query results that are important with respect to his bookmarked pages would be ranked higher. While experimentation with the use of personalized *PageRank* has shown its utility and promise [5, 11], the size of the web makes its practical realization extremely difficult. To see why, let us review the intuition behind the *PageRank* algorithm and its extension for personalization.

The fundamental motivation underlying *PageRank* is the recursive notion that important pages are those linked-to by many important pages. A page with only two in-links, for example, may seem unlikely to be an important page, but it may be important if the two referencing pages are *Yahoo!* and *Netscape*, which themselves are important pages because they have numerous in-links. One way to formalize this recursive notion is to use the “random surfer” model introduced in [11]. Imagine that trillions of *random surfers* are browsing the web: if at a certain time step a surfer is looking at page  $p$ , at the next time step he looks at a random out-neighbor of  $p$ . As time goes on, the expected percentage of surfers at each page  $p$  converges (under certain conditions) to a limit  $r(p)$  that is independent of the distribution of starting points. Intuitively, this limit is the *PageRank* of  $p$ , and is taken to be an importance score for  $p$ , since it reflects the number of people expected to be looking at  $p$  at any one time.

The *PageRank* score  $r(p)$  reflects a “democratic” importance that has no preference for any particular pages. In reality, a user may have a set  $P$  of preferred pages (such as his bookmarks) which he considers more interesting. We can account for preferred pages in the random surfer model by introducing a “teleportation” probability  $c$ : at each step, a surfer jumps back to a random page in  $P$  with probability  $c$ , and with probability  $1 - c$  continues forth along a hyperlink. The limit distribution of surfers in this model would favor pages in  $P$ , pages linked-to by  $P$ , pages linked-to in turn, etc. We represent this distribution as a *personalized PageRank vector (PPV)* personalized on the set  $P$ . Informally, a PPV is a personalized view of the importance of pages on the web. Rankings of a user’s text-based query results can be biased according to a PPV instead of the global importance distribution.

Each PPV is of length  $n$ , where  $n$  is the number of pages on the web. Computing a PPV naively using a fixed-point iteration requires multiple scans of the web graph [11], which makes it impossible to carry out online in response to a user query. On the other hand, PPV’s for all preference sets, of which there are  $2^n$ , is far too large to compute and store offline. We present a method

\*This work was supported by the National Science Foundation under grant IIS-9817799. This is an abbreviated version of the full paper that omits appendices. The full version is available on the web at <http://dbpubs.stanford.edu/pub/2002-12>

Copyright is held by the author/owner(s).  
WWW2003, May 20–24, 2003, Budapest, Hungary.  
ACM 1-58113-680-3/03/0005.

for encoding PPV's as partially-computed, shared vectors that are practical to compute and store offline, and from which PPV's can be computed quickly at query time.

In our approach we restrict preference sets  $P$  to subsets of a set of *hub* pages  $H$ , selected as those of greater interest for personalization. In practice, we expect  $H$  to be a set of pages with high PageRank ("important pages"), pages in a human-constructed directory such as *Yahoo!* or *Open Directory* [2], or pages important to a particular enterprise or application. The size of  $H$  can be thought of as the available degree of personalization. We present algorithms that, unlike previous work [5, 11], scale well with the size of  $H$ . Moreover, the same techniques we introduce can yield approximations on the much broader set of all PPV's, allowing at least some level of personalization on arbitrary preference sets.

The main contributions of this paper are as follows.

- A method, based on new graph-theoretical results (listed next), of encoding PPV's as *partial quantities*, enabling an efficient, scalable computation that can be divided between pre-computation time and query time, in a customized fashion according to available resources and application requirements.
- Three main theorems: The *Linearity Theorem* allows every PPV to be represented as a linear combination of *basis vectors*, yielding a natural way to construct PPV's from shared components. The *Hubs Theorem* allows basis vectors to be encoded as *partial vectors* and a *hubs skeleton*, enabling basis vectors themselves to be constructed from common components. The *Decomposition Theorem* establishes a linear relationship among basis vectors, which is exploited to minimize redundant computation.
- Several algorithms for computing basis vectors, specializations of these algorithms for computing partial vectors and the hubs skeleton, and an algorithm for constructing PPV's from partial vectors using the hubs skeleton.
- Experimental results on real web data demonstrating the effectiveness and scalability of our techniques.

In Section 2 we introduce the notation used in this paper and formalize personalized PageRank mathematically. Section 3 presents basis vectors, the first step towards encoding PPV's as shared components. The full encoding is presented in Section 4. Section 5 discusses the computation of partial quantities. Experimental results are presented in Section 6. Related work is discussed in Section 7. Section 8 summarizes the contributions of this paper.

Due to space constraints, this paper omits proofs of the theorems and algorithms presented. These proofs are included as appendices in the full version of this paper [7].

## 2. PRELIMINARIES

Let  $G = (V, E)$  denote the *web graph*, where  $V$  is the set of all web pages and  $E$  contains a directed edge  $\langle p, q \rangle$  iff page  $p$  links to page  $q$ . For a page  $p$ , we denote by  $I(p)$  and  $O(p)$  the set of in-neighbors and out-neighbors of  $p$ , respectively. Individual in-neighbors are denoted as  $I_i(p)$  ( $1 \leq i \leq |I(p)|$ ), and individual out-neighbors are denoted analogously. For convenience, pages are numbered from 1 to  $n$ , and we refer to a page  $p$  and its associated number  $i$  interchangeably. For a vector  $\mathbf{v}$ ,  $v(p)$  denotes *entry*  $p$ , the  $p$ -th component of  $\mathbf{v}$ . We always typeset vectors in boldface and scalars (e.g.,  $v(p)$ ) in normal font. All vectors in this paper are  $n$ -dimensional and have nonnegative entries. They should be thought of as distributions rather than arrows. The *magnitude* of a vector  $\mathbf{v}$

is defined to be  $\sum_{i=1}^n v(i)$  and is written  $|\mathbf{v}|$ . In this paper, vector magnitudes are always in  $[0, 1]$ . In an implementation, a vector may be represented as a list of its nonzero entries, so another useful measure is the *size* of  $\mathbf{v}$ , the number of nonzero entries in  $\mathbf{v}$ .

We generalize the preference set  $P$  discussed in Section 1 to a *preference vector*  $\mathbf{u}$ , where  $|\mathbf{u}| = 1$  and  $u(p)$  denotes the amount of preference for page  $p$ . For example, a user who wants to personalize on his bookmarked pages  $P$  uniformly would have a  $\mathbf{u}$  where  $u(p) = \frac{1}{|P|}$  if  $p \in P$ , and  $u(p) = 0$  if  $p \notin P$ . We formalize personalized PageRank scoring using matrix-vector equations. Let  $\mathbf{A}$  be the matrix corresponding to the web graph  $G$ , where  $A_{ij} = \frac{1}{|O(j)|}$  if page  $j$  links to page  $i$ , and  $A_{ij} = 0$  otherwise. For simplicity of presentation, we assume that every page has at least one out-neighbor, as can be enforced by adding self-links to pages without out-links. The resulting scores can be adjusted to account for the (minor) effects of this modification, as specified in the appendices of the full version of this paper [7].

For a given  $\mathbf{u}$ , the personalized PageRank equation can be written as

$$\mathbf{v} = (1 - c)\mathbf{A}\mathbf{v} + c\mathbf{u} \quad (1)$$

where  $c \in (0, 1)$  is the "teleportation" constant discussed in Section 1. Typically  $c \approx 0.15$ , and experiments have shown that small changes in  $c$  have little effect in practice [11]. A solution  $\mathbf{v}$  to equation (1) is a steady-state distribution of random surfers under the model discussed in Section 1, where at each step a surfer teleports to page  $p$  with probability  $c \cdot u(p)$ , or moves to a random out-neighbor otherwise [11]. By a theorem of Markov Theory, a solution  $\mathbf{v}$  with  $|\mathbf{v}| = 1$  always exists and is unique [10].<sup>1</sup> The solution  $\mathbf{v}$  is the *personalized PageRank vector (PPV)* for preference vector  $\mathbf{u}$ . If  $\mathbf{u}$  is the uniform distribution vector  $\mathbf{u} = [1/n, \dots, 1/n]$ , then the corresponding solution  $\mathbf{v}$  is the *global PageRank vector* [11], which gives no preference to any pages.

For the reader's convenience, Table 1 on the next page lists terminology that will be used extensively in the coming sections.

## 3. BASIS VECTORS

We present the first step towards encoding PPV's as shared components. The motivation behind the encoding is a simple observation about the linearity<sup>2</sup> of PPV's, formalized by the following theorem.

**THEOREM 1 (LINEARITY).** *For any preference vectors  $\mathbf{u}_1$  and  $\mathbf{u}_2$ , if  $\mathbf{v}_1$  and  $\mathbf{v}_2$  are the two corresponding PPV's, then for any constants  $\alpha_1, \alpha_2 \geq 0$  such that  $\alpha_1 + \alpha_2 = 1$ ,*

$$\alpha_1 \mathbf{v}_1 + \alpha_2 \mathbf{v}_2 = (1 - c)\mathbf{A}(\alpha_1 \mathbf{v}_1 + \alpha_2 \mathbf{v}_2) + c(\alpha_1 \mathbf{u}_1 + \alpha_2 \mathbf{u}_2) \quad (2)$$

Informally, the Linearity Theorem says that the solution to a linear combination of preference vectors  $\mathbf{u}_1$  and  $\mathbf{u}_2$  is the same linear combination of the corresponding PPV's  $\mathbf{v}_1$  and  $\mathbf{v}_2$ . The proof is in the full version [7].

Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  be the unit vectors in each dimension, so that for each  $i$ ,  $\mathbf{x}_i$  has value 1 at entry  $i$  and 0 everywhere else. Let  $\mathbf{r}_i$  be the PPV corresponding to  $\mathbf{x}_i$ . Each *basis vector*  $\mathbf{r}_i$  gives the distribution of random surfers under the model that at each step, surfers teleport back to page  $i$  with probability  $c$ . It can be thought of as representing page  $i$ 's view of the web, where entry  $j$  of  $\mathbf{r}_i$  is

<sup>1</sup>Specifically,  $\mathbf{v}$  corresponds to the steady-state distribution of an ergodic, aperiodic Markov chain.

<sup>2</sup>More precisely, the transformation from personalization vectors  $\mathbf{u}$  to their corresponding solution vectors  $\mathbf{v}$  is linear.

Term	Description	Section
Hub Set $H$	A subset of web pages.	1
Preference Set $P$	Set of pages on which to personalize (restricted in this paper to subsets of $H$ ).	1
Preference Vector $\mathbf{u}$	Preference set with weights.	2
Personalized PageRank Vector (PPV)	Importance distribution induced by a preference vector.	2
Basis Vector $\mathbf{r}_p$ (or $\mathbf{r}_i$ )	PPV for a preference vector with a single nonzero entry at $p$ (or $i$ ).	3
Hub Vector $\mathbf{r}_p$	Basis vector for a hub page $p \in H$ .	3
Partial Vector $(\mathbf{r}_p - \mathbf{r}_p^H)$	Used with the hubs skeleton to construct a hub vector.	4.2
Hubs Skeleton $S$	Used with partial vectors to construct a hub vector.	4.3
Web Skeleton	Extension of the hubs skeleton to include pages not in $H$ .	4.4.3
Partial Quantities	Partial vectors and the hubs, web skeletons.	
Intermediate Results	Maintained during iterative computations.	5.2

**Table 1: Summary of terms.**

$j$ 's importance in  $i$ 's view. Note that the global PageRank vector is  $\frac{1}{n}(\mathbf{r}_1 + \dots + \mathbf{r}_n)$ , the average of every page's view.

An arbitrary personalization vector  $\mathbf{u}$  can be written as a weighted sum of the unit vectors  $\mathbf{x}_i$ :

$$\mathbf{u} = \sum_{i=1}^n \alpha_i \mathbf{x}_i \quad (3)$$

for some constants  $\alpha_1, \dots, \alpha_n$ . By the Linearity Theorem,

$$\mathbf{v} = \sum_{i=1}^n \alpha_i \mathbf{r}_i \quad (4)$$

is the corresponding PPV, expressed as a linear combination of the basis vectors  $\mathbf{r}_i$ .

Recall from Section 1 that preference sets (now preference vectors) are restricted to subsets of a set of hub pages  $H$ . If a *basis hub vector* (or hereafter *hub vector*) for each  $p \in H$  were computed and stored, then any PPV corresponding to a preference set  $P$  of size  $k$  (a preference vector with  $k$  nonzero entries) can be computed by adding up the  $k$  corresponding hub vectors  $\mathbf{r}_p$  with the appropriate weights  $\alpha_p$ .

Each hub vector can be computed naively using the fixed-point computation in [11]. However, each fixed-point computation is expensive, requiring multiple scans of the web graph, and the computation time (as well as storage cost) grows linearly with the number of hub vectors  $|H|$ . In the next section, we enable a more scalable computation by constructing hub vectors from shared components.

## 4. DECOMPOSITION OF BASIS VECTORS

In Section 3 we represented PPV's as a linear combination of  $|H|$  hub vectors  $\mathbf{r}_p$ , one for each  $p \in H$ . Any PPV based on hub pages can be constructed quickly from the set of precomputed hub vectors, but computing and storing all hub vectors is impractical. To compute a large number of hub vectors efficiently, we further decompose them into *partial vectors* and the *hubs skeleton*, components from which hub vectors can be constructed quickly at query time. The representation of hub vectors as partial vectors and the hubs skeleton saves both computation time and storage due to sharing of components among hub vectors. Note, however, that depending on available resources and application requirements, hub vectors can be constructed offline as well. Thus "query time" can be thought of more generally as "construction time".

We compute one partial vector for each hub page  $p$ , which essentially encodes the part of the hub vector  $\mathbf{r}_p$  unique to  $p$ , so that

components shared among hub vectors are not computed and stored redundantly. The complement to the partial vectors is the hubs skeleton, which succinctly captures the interrelationships among hub vectors. It is the "blueprint" by which partial vectors are assembled to form a hub vector, as we will see in Section 4.3.

The mathematical tools used in the formalization of this decomposition are presented next.<sup>3</sup>

### 4.1 Inverse P-distance

To formalize the relationship among hub vectors, we relate the personalized PageRank scores represented by PPV's to *inverse P-distances* in the web graph, a concept based on *expected-f distances* as introduced in [8].

Let  $p, q \in V$ . We define the *inverse P-distance*  $r'_p(q)$  from  $p$  to  $q$  as

$$r'_p(q) = \sum_{t: p \rightsquigarrow q} P[t] c(1-c)^{l(t)} \quad (5)$$

where the summation is taken over all *tours*  $t$  (paths that may contain cycles) starting at  $p$  and ending at  $q$ , possibly touching  $p$  or  $q$  multiple times. For a tour  $t = \langle w_1, \dots, w_k \rangle$ , the length  $l(t)$  is  $k-1$ , the number of edges in  $t$ . The term  $P[t]$ , which should be interpreted as "the probability of traveling  $t$ ", is defined as  $\prod_{i=1}^{k-1} \frac{1}{|O(w_i)|}$ , or 1 if  $l(t) = 0$ . If there is no tour from  $p$  to  $q$ , the summation is taken to be 0.<sup>4</sup> Note that  $r'_p(q)$  measures distances inversely: it is higher for nodes  $q$  "closer" to  $p$ . As suggested by the notation and proven in the full version [7],  $r'_p(q) = r_p(q)$  for all  $p, q \in V$ , so we will use  $r_p(q)$  to denote both the inverse P-distance and the personalized PageRank score. Thus PageRank scores can be viewed as an inverse measure of distance.

Let  $H \subseteq V$  be some nonempty set of pages. For  $p, q \in V$ , we define  $r_p^H(q)$  as a restriction of  $r_p(q)$  that considers only tours which pass through some page  $h \in H$  in equation (5). That is, a page  $h \in H$  must occur on  $t$  somewhere other than the endpoints.

<sup>3</sup>Note that while the mathematics and computation strategies in this paper are presented in the specific context of the web graph, they are general graph-theoretical results that may be applicable in other scenarios involving stochastic processes, of which PageRank is one example.

<sup>4</sup>The definition here of inverse P-distance differs slightly from the concept of expected- $f$  distance in [8], where tours are not allowed to visit  $q$  multiple times. Note that general expected- $f$  distances have the form  $\sum_t P[t] f(l(t))$ ; in our definition,  $f(x) = c(1-c)^x$ .

Precisely,  $r_p^H(q)$  is written as

$$r_p^H(q) = \sum_{t: p \rightsquigarrow H \rightsquigarrow q} P[t]c(1-c)^{l(t)} \quad (6)$$

where the notation  $t : p \rightsquigarrow H \rightsquigarrow q$  reminds us that  $t$  passes through some page in  $H$ . Note that  $t$  must be of length at least 2. In this paper,  $H$  is always the set of hub pages, and  $p$  is usually a hub page (until we discuss the web skeleton in Section 4.4.3).

## 4.2 Partial Vectors

Intuitively,  $r_p^H(q)$ , defined in equation (6), is the influence of  $p$  on  $q$  through  $H$ . In particular, if all paths from  $p$  to  $q$  pass through a page in  $H$ , then  $H$  separates  $p$  and  $q$ , and  $r_p^H(q) = r_p(q)$ . For well-chosen sets  $H$  (discussed in Section 4.4.2), it will be true that  $r_p(q) - r_p^H(q) = 0$  for many pages  $p, q$ . Our strategy is to take advantage of this property by breaking  $r_p$  into two components:  $(r_p - r_p^H)$  and  $r_p^H$ , using the equation

$$r_p = (r_p - r_p^H) + r_p^H \quad (7)$$

We first precompute and store the *partial vector*  $(r_p - r_p^H)$  instead of the full hub vector  $r_p$ . Partial vectors are cheaper to compute and store than full hub vectors, assuming they are represented as a list of their nonzero entries. Moreover, the size of each partial vector decreases as  $|H|$  increases, making this approach particularly scalable. We then add  $r_p^H$  back at query time to compute the full hub vector. However, computing and storing  $r_p^H$  explicitly could be as expensive as  $r_p$  itself. In the next section we show how to encode  $r_p^H$  so it can be computed and stored efficiently.

## 4.3 Hubs Skeleton

Let us briefly review where we are: In Section 3 we represented PPV's as linear combinations of hub vectors  $r_p$ , one for each  $p \in H$ , so that we can construct PPV's quickly at query time if we have precomputed the hub vectors, a relatively small subset of PPV's. To encode hub vectors efficiently, in Section 4.2 we said that instead of full hub vectors  $r_p$ , we first compute and store only partial vectors  $(r_p - r_p^H)$ , which intuitively account only for paths that do not pass through a page of  $H$  (i.e., the distribution is “blocked” by  $H$ ). Computing and storing the difference vector  $r_p^H$  efficiently is the topic of this section.

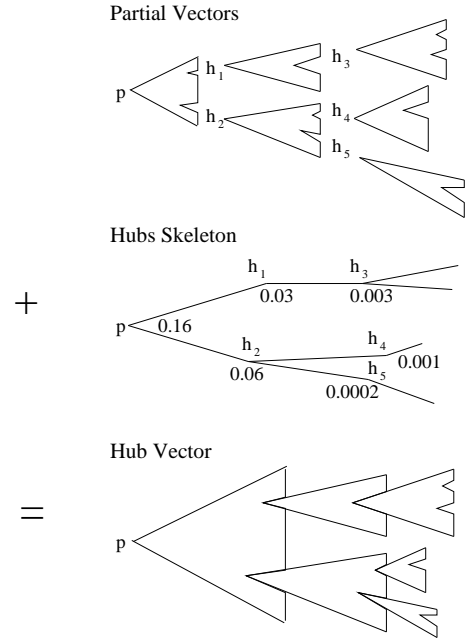
It turns out that the vector  $r_p^H$  can be expressed in terms of the partial vectors  $(r_h - r_h^H)$ , for  $h \in H$ , as shown by the following theorem. Recall from Section 3 that  $x_h$  has value 1 at  $h$  and 0 everywhere else.

**THEOREM 2 (HUBS).** *For any  $p \in V$ ,  $H \subseteq V$ ,*

$$r_p^H = \frac{1}{c} \sum_{h \in H} (r_p(h) - c \cdot x_p(h)) (r_h - r_h^H - cx_h) \quad (8)$$

In terms of inverse P-distances (Section 4.1), the Hubs Theorem says roughly that the distance from page  $p$  to any page  $q \in V$  through  $H$  is the distance  $r_p(h)$  from  $p$  to each  $h \in H$  times the distance  $r_h(q)$  from  $h$  to  $q$ , correcting for the paths among hubs by  $r_h^H(q)$ . The terms  $c \cdot x_p(h)$  and  $cx_h$  deal with the special cases when  $p$  or  $q$  is itself in  $H$ . The proof, which is quite involved, is in the full version [7].

The quantity  $(r_h - r_h^H)$  appearing on the right-hand side of (8) is exactly the partial vectors discussed in Section 4.2. Suppose we have computed  $r_p(H) = \{(h, r_p(h)) \mid h \in H\}$  for a hub page  $p$ . Substituting the Hubs Theorem into equation 7, we have the following *Hubs Equation* for constructing the hub vector  $r_p$  from



**Figure 1: Intuitive view of the construction of hub vectors from partial vectors and the hubs skeleton.**

partial vectors:

$$r_p = (r_p - r_p^H) + \frac{1}{c} \sum_{h \in H} (r_p(h) - c \cdot x_p(h)) [(r_h - r_h^H) - cx_h] \quad (9)$$

This equation is central to the construction of hub vectors from partial vectors.

The set  $r_p(H)$  has size at most  $|H|$ , much smaller than the full hub vector  $r_p$ , which can have up to  $n$  nonzero entries. Furthermore, the contribution of each entry  $r_p(h)$  to the sum is no greater than  $r_p(h)$  (and usually much smaller), so that small values of  $r_p(h)$  can be omitted with minimal loss of precision (Section 6). The set  $S = \{r_p(H) \mid p \in H\}$  forms the *hubs skeleton*, giving the interrelationships among partial vectors.

An intuitive view of the encoding and construction suggested by the Hubs Equation (9) is shown in Figure 1. At the top, each partial vector  $(r_h - r_h^H)$ , including  $(r_p - r_p^H)$ , is depicted as a notched triangle labeled  $h$  at the tip. The triangle can be thought of as representing paths starting at  $h$ , although, more accurately, it represents the distribution of importance scores computed based on the paths, as discussed in Section 4.1. A notch in the triangle shows where the computation of a partial vector “stopped” at another hub page. At the center, a part  $r_p(H)$  of the hubs skeleton is depicted as a tree so the “assembly” of the hub vector can be visualized. The hub vector is constructed by logically assembling the partial vectors using the corresponding weights in the hubs skeleton, as shown at the bottom.

## 4.4 Discussion

### 4.4.1 Summary

In summary, hub vectors are building blocks for PPV's corresponding to preference vectors based on hub pages. Partial vectors, together with the hubs skeleton, are building blocks for hub vec-

tors. Transitively, partial vectors and the hubs skeleton are building blocks for PPV's: they can be used to construct PPV's without first materializing hub vectors as an intermediate step (Section 5.4). Note that for preference vectors based on multiple hub pages, constructing the corresponding PPV from partial vectors directly can result in significant savings versus constructing from hub vectors, since partial vectors are shared across multiple hub vectors.

#### 4.4.2 Choice of $H$

So far we have made no assumptions about the set of hub pages  $H$ . Not surprisingly, the choice of hub pages can have a significant impact on performance, depending on the location of hub pages within the overall graph structure. In particular, the size of partial vectors is smaller when pages in  $H$  have higher PageRank, since high-PageRank pages are on average close to other pages in terms of inverse P-distance (Section 4.1), and the size of the partial vectors is related to the inverse P-distance between hub pages and other pages according to the Hubs Theorem. Our intuition is that high-PageRank pages are generally more interesting for personalization anyway, but in cases where the intended hub pages do not have high PageRank, it may be beneficial to include some high-PageRank pages in  $H$  to improve performance. We ran experiments confirming that the size of partial vectors is much smaller using high-PageRank pages as hubs than using random pages.

#### 4.4.3 Web Skeleton

The techniques used in the construction of hub vectors can be extended to enable at least approximate personalization on arbitrary preference vectors that are not necessarily based on  $H$ . Suppose we want to personalize on a page  $p \notin H$ . The Hubs Equation can be used to construct  $\mathbf{r}_p^H$  from partial vectors, given that we have computed  $\mathbf{r}_p(H)$ . As discussed in Section 4.3, the cost of computing and storing  $\mathbf{r}_p(H)$  is orders of magnitude less than  $\mathbf{r}_p$ . Though  $\mathbf{r}_p^H$  is only an approximation to  $\mathbf{r}_p$ , it may still capture significant personalization information for a properly-chosen hub set  $H$ , as  $\mathbf{r}_p^H$  can be thought of as a “projection” of  $\mathbf{r}_p$  onto  $H$ . For example, if  $H$  contains pages from *Open Directory*,  $\mathbf{r}_p^H$  can capture information about the broad topic of  $\mathbf{r}_p$ . Exploring the utility of the *web skeleton*  $W = \{\mathbf{r}_p(H) \mid p \in V\}$  is an area of future work.

## 5. COMPUTATION

In Section 4 we presented a way to construct hub vectors from partial vectors ( $\mathbf{r}_p - \mathbf{r}_p^H$ ), for  $p \in H$ , and the hubs skeleton  $S = \{\mathbf{r}_p(H) \mid p \in H\}$ . We also discussed the web skeleton  $W = \{\mathbf{r}_p(H) \mid p \in V\}$ . Computing these *partial quantities* naively using a fixed-point iteration [11] for each  $p$  would scale poorly with the number of hub pages. Here we present scalable algorithms that compute these quantities efficiently by using dynamic programming to leverage the interrelationships among them. We also show how PPV's can be constructed from partial vectors and the hubs skeleton at query time. All of our algorithms have the property that they can be stopped at any time (e.g., when resources are depleted), so that the current “best results” can be used as an approximation, or the computation can be resumed later for increased precision if resources permit.

We begin in Section 5.1 by presenting a theorem underlying all of the algorithms presented (as well as the connection between PageRank and inverse P-distance, as shown in the full version [7]). In Section 5.2, we present three algorithms, based on this theorem, for computing general basis vectors. The algorithms in Section 5.2 are not meant to be deployed, but are used as foundations for the algorithms in Section 5.3 for computing partial quantities. Section

5.4 discusses the construction of PPV's from partial vectors and the hubs skeleton.

### 5.1 Decomposition Theorem

Recall the random surfer model of Section 1, instantiated for preference vector  $\mathbf{u} = \mathbf{x}_p$  (for page  $p$ 's view of the web). At each step, a surfer  $s$  teleports to page  $p$  with some probability  $c$ . If  $s$  is at  $p$ , then at the next step,  $s$  with probability  $1 - c$  will be at a random out-neighbor of  $p$ . That is, a fraction  $(1 - c) \frac{1}{|O(p)|}$  of the time, surfer  $s$  will be at any given out-neighbor of  $p$  one step after teleporting to  $p$ . This behavior is strikingly similar to the model instantiated for preference vector  $\mathbf{u}' = \frac{1}{|O(p)|} \sum_{i=1}^{|O(p)|} \mathbf{x}_{O_i(p)}$ , where surfers teleport directly to each  $O_i(p)$  with equal probability  $\frac{1}{|O(p)|}$ . The similarity is formalized by the following theorem.

**THEOREM 3 (DECOMPOSITION).** *For any  $p \in V$ ,*

$$\mathbf{r}_p = \frac{(1 - c)}{|O(p)|} \sum_{i=1}^{|O(p)|} \mathbf{r}_{O_i(p)} + c\mathbf{x}_p \quad (10)$$

The Decomposition Theorem says that the basis vector  $\mathbf{r}_p$  for  $p$  is an average of the basis vectors  $\mathbf{r}_{O_i(p)}$  for its out-neighbors, plus a compensation factor  $c\mathbf{x}_p$ . The proof is in the full version [7].

The Decomposition Theorem gives another way to think about PPV's. It says that  $p$ 's view of the web ( $\mathbf{r}_p$ ) is the average of the views of its out-neighbors, but with extra importance given to  $p$  itself. That is, pages important in  $p$ 's view are either  $p$  itself, or pages important in the view of  $p$ 's out-neighbors, which are themselves “endorsed” by  $p$ . In fact, this recursive intuition yields an equivalent way of formalizing personalized PageRank scoring: basis vectors can be defined as vectors satisfying the Decomposition Theorem.

While the Decomposition Theorem identifies relationships among basis vectors, a division of the computation of a basis vector  $\mathbf{r}_p$  into related subproblems for dynamic programming is not inherent in the relationships. For example, it is possible to compute some basis vectors first and then to compute the rest using the former as solved subproblems. However, the presence of cycles in the graph makes this approach ineffective. Instead, our approach is to consider as a subproblem the computation of a vector to less precision. For example, having computed  $\mathbf{r}_{O_i(p)}$  to a certain precision, we can use the Decomposition Theorem to combine the  $\mathbf{r}_{O_i(p)}$ 's to compute  $\mathbf{r}_p$  to greater precision. This approach has the advantage that precision needs not be fixed in advance: the process can be stopped at any time for the current best answer.

### 5.2 Algorithms for Computing Basis Vectors

We present three algorithms in the general context of computing full basis vectors. These algorithms are presented primarily to develop our algorithms for computing partial quantities, presented in Section 5.3. All three algorithms are iterative fixed-point computations that maintain a set of *intermediate results* ( $\mathbf{D}_k[*]$ ,  $\mathbf{E}_k[*]$ ). For each  $p$ ,  $\mathbf{D}_k[p]$  is a lower-approximation of  $\mathbf{r}_p$  on iteration  $k$ , i.e.,  $\mathbf{D}_k[p](q) \leq \mathbf{r}_p(q)$  for all  $q \in V$ . We build solutions  $\mathbf{D}_k[p]$  ( $k = 0, 1, 2, \dots$ ) that are successively better approximations to  $\mathbf{r}_p$ , and simultaneously compute the error components  $\mathbf{E}_k[p]$ , where  $\mathbf{E}_k[p]$  is the “projection” of the vector  $(\mathbf{r}_p - \mathbf{D}_k[p])$  onto the (actual) basis vectors. That is, we maintain the invariant that for all  $k \geq 0$  and all  $p \in V$ ,

$$\mathbf{D}_k[p] + \sum_{q \in V} \mathbf{E}_k[p](q) \mathbf{r}_q = \mathbf{r}_p \quad (11)$$

Thus,  $D_k[p]$  is a lower-approximation of  $r_p$  with error

$$\left| \sum_{q \in V} E_k[p](q) r_q \right| = |E_k[p]|$$

We begin with  $D_0[p] = \mathbf{0}$  and  $E_0[p] = x_p$ , so that logically, the approximation is initially  $\mathbf{0}$  and the error is  $r_p$ . To store  $E_k[p]$  and  $D_k[p]$  efficiently, we can represent them in an implementation as a list of their nonzero entries. While all three algorithms have in common the use of these intermediate results, they differ in how they use the Decomposition Theorem to refine intermediate results on successive iterations.

It is important to note that the algorithms presented in this section and their derivatives in Section 5.3 compute vectors to arbitrary precision; they are not approximations. In practice, the precision desired may vary depending on the application. Our focus is on algorithms that are efficient and scalable with the number of hub vectors, regardless of the precision to which vectors are computed.

### 5.2.1 Basic Dynamic Programming Algorithm

In the *basic dynamic programming algorithm*, a new basis vector for each page  $p$  is computed on each iteration using the vectors computed for  $p$ 's out-neighbors on the previous iteration, via the Decomposition Theorem. On iteration  $k$ , we derive  $(D_{k+1}[p], E_{k+1}[p])$  from  $(D_k[p], E_k[p])$  using the equations:

$$D_{k+1}[p] = \frac{1-c}{|O(p)|} \sum_{i=1}^{|O(p)|} D_k[O_i(p)] + c x_p \quad (12)$$

$$E_{k+1}[p] = \frac{1-c}{|O(p)|} \sum_{i=1}^{|O(p)|} E_k[O_i(p)] \quad (13)$$

A proof of the algorithm's correctness is given in the full version [7], where the error  $|E_k[p]|$  is shown to be reduced by a factor of  $1 - c$  on each iteration.

Note that although the  $E_k[*]$  values help us to see the correctness of the algorithm, they are not used here in the computation of  $D_k[*]$  and can be omitted in an implementation (although they will be used to compute partial quantities in Section 5.3). The sizes of  $D_k[p]$  and  $E_k[p]$  grow with the number of iterations, and in the limit they can be up to the size of  $r_p$ , which is the number of pages reachable from  $p$ . Intermediate scores  $(D_k[*], E_k[*])$  will likely be much larger than available main memory, and in an implementation  $(D_k[*], E_k[*])$  could be read off disk and  $(D_{k+1}[*], E_{k+1}[*])$  written to disk on each iteration. When the data for one iteration has been computed, data from the previous iteration may be deleted. Specific details of our implementation are discussed in Section 6.

### 5.2.2 Selective Expansion Algorithm

The *selective expansion algorithm* is essentially a version of the naive algorithm that can readily be modified to compute partial vectors, as we will see in Section 5.3.1.

We derive  $(D_{k+1}[p], E_{k+1}[p])$  by "distributing" the error at each page  $q$  (that is,  $E_k[p](q)$ ) to its out-neighbors via the Decomposition Theorem. Precisely, we compute results on iteration- $k$  us-

ing the equations:

$$D_{k+1}[p] = D_k[p] + \sum_{q \in Q_k(p)} c \cdot E_k[p](q) x_q \quad (14)$$

$$E_{k+1}[p] = E_k[p] - \sum_{q \in Q_k(p)} E_k[p](q) x_q + \sum_{q \in Q_k(p)} \frac{1-c}{|O(q)|} \sum_{i=1}^{|O(q)|} E_k[p](q) x_{O_i(q)} \quad (15)$$

for a subset  $Q_k(p) \subseteq V$ . If  $Q_k(p) = V$  for all  $k$ , then the error is reduced by a factor of  $1 - c$  on each iteration, as in the basic dynamic programming algorithm. However, it is often useful to choose a selected subset of  $V$  as  $Q_k(p)$ . For example, if  $Q_k(p)$  contains the  $m$  pages  $q$  for which the error  $E_k[p](q)$  is highest, then this *top- $m$*  scheme limits the number of expansions and delays the growth in size of the intermediate results while still reducing much of the error. In Section 5.3.1, we will compute the hub vectors by choosing  $Q_k(p) = H$ . The correctness of selective expansion is proven in the full version [7].

### 5.2.3 Repeated Squaring Algorithm

The *repeated squaring algorithm* is similar to the selective expansion algorithm, except that instead of extending  $(D_{k+1}[*], E_{k+1}[*])$  one step using equations (14) and (15), we compute what are essentially iteration- $2k$  results using the equations

$$D_{2k}[p] = D_k[p] + \sum_{q \in Q_k(p)} E_k[p](q) D_k[q] \quad (16)$$

$$E_{2k}[p] = E_k[p] - \sum_{q \in Q_k(p)} E_k[p](q) x_q + \sum_{q \in Q_k(p)} E_k[p](q) E_k[q] \quad (17)$$

where  $Q_k(p) \subseteq V$ . For now we can assume that  $Q_k(p) = V$  for all  $p$ ; we will set  $Q_k(p) = H$  to compute the hubs skeleton in Section 5.3.2. The correctness of these equations is proven in the full version [7], where it is shown that repeated squaring reduces the error much faster than the basic dynamic programming or selective expansion algorithms. If  $Q_k(p) = V$ , the error is squared on each iteration, as equation (17) reduces to:

$$E_{2k}[p] = \sum_{q \in V} E_k[p](q) E_k[q] \quad (18)$$

As an alternative to taking  $Q_k(p) = V$ , we can also use the top- $m$  scheme of Section 5.2.2.

Note that while all three algorithms presented can be used to compute the set of all basis vectors, they differ in their requirements on the computation of other vectors when computing  $r_p$ : the basic dynamic programming algorithm requires the vectors of out-neighbors of  $p$  to be computed as well, repeated squaring requires results  $(D_k[q], E_k[q])$  to be computed for  $q$  such that  $E_k[p](q) > 0$ , and selective expansion computes  $r_p$  independently.

## 5.3 Computing Partial Quantities

In Section 5.2 we presented iterative algorithms for computing full basis vectors to arbitrary precision. Here we present modifications to these algorithms to compute the partial quantities:

- Partial vectors  $(r_p - r_p^H), p \in H$ .
- The hubs skeleton  $S = \{r_p(H) \mid p \in H\}$  (which can be computed more efficiently by itself than as part of the entire web skeleton).

- The web skeleton  $W = \{r_p(H) \mid p \in V\}$ .

Each partial quantity can be computed in time no greater than its size, which is far less than the size of the hub vectors.

### 5.3.1 Partial Vectors

Partial vectors can be computed using a simple specialization of the selective expansion algorithm (Section 5.2.2): we take  $Q_0(p) = V$  and  $Q_k(p) = V - H$  for  $k > 0$ , for all  $p \in V$ . That is, we never “expand” hub pages after the first step, so tours passing through a hub page  $H$  are never considered. Under this choice of  $Q_k(p)$ ,  $D_k[p] + cE_k[p]$  converges to  $(r_p - r_p^H)$  for all  $p \in V$ . Of course, only the intermediate results  $(D_k[p], E_k[p])$  for  $p \in H$  should be computed. A proof is presented in the full version [7].

This algorithm makes it clear why using high-PageRank pages as hub pages improves performance: from a page  $p$  we expect to reach a high-PageRank page  $q$  sooner than a random page, so the expansion from  $p$  will stop sooner and result in a shorter partial vector.

### 5.3.2 Hubs Skeleton

While the hubs skeleton is a subset of the complete web skeleton and can be computed as such using the technique to be presented in Section 5.3.3, it can be computed much faster by itself if we are not interested in the entire web skeleton, or if higher precision is desired for the hubs skeleton than can be computed for the entire web skeleton.

We use a specialization of the repeated squaring algorithm (Section 5.2.3) to compute the hubs skeleton, using the intermediate results from the computation of partial vectors. Suppose  $(D_k[p], E_k[p])$ , for  $k \geq 1$ , have been computed by the algorithm of Section 5.3.1, so that  $\sum_{q \notin H} E_k[p](q) < \epsilon$ , for some error  $\epsilon$ . We apply the repeated squaring algorithm on these results using  $Q_k(p) = H$  for all successive iterations. As shown in the full version [7], after  $i$  iterations of repeated squaring, the total error  $|E_i[p]|$  is bounded by  $(1 - c)^{2^i} + \epsilon/c$ . Thus, by varying  $k$  and  $i$ ,  $r_p(H)$  can be computed to arbitrary precision.

Notice that only the intermediate results  $(D_k[h], E_k[h])$  for  $h \in H$  are ever needed to update scores for  $D_k[p]$ , and of the former, only the entries  $D_k[h](q), E_k[h](q)$ , for  $q \in H$ , are used to compute  $D_k[p](q)$ . Since we are only interested in the hub scores  $D_k[p](q)$ , we can simply drop all non-hub entries from the intermediate results. The running time and storage would then depend only on the size of  $r_p(H)$  and not on the length of the entire hub vectors  $r_p$ . If the restricted intermediate results fit in main memory, it is possible to defer the computation of the hubs skeleton to query time.

### 5.3.3 Web Skeleton

To compute the entire web skeleton, we modify the basic dynamic programming algorithm (Section 5.2.1) to compute only the hub scores  $r_p(H)$ , with corresponding savings in time and memory usage. We restrict the computation by eliminating entries  $q \notin H$  from the intermediate results  $(D_k[p], E_k[p])$ , similar to the technique used in computing the hubs skeleton.

The justification for this modification is that the hub score  $D_{k+1}[p](h)$  is affected only by the hub scores  $D_k[*](h)$  of the previous iteration, so that  $D_{k+1}[p](h)$  in the modified algorithm is equal to that in the basic algorithm. Since  $|H|$  is likely to be orders of magnitude less than  $n$ , the size of the intermediate results is reduced significantly.

## 5.4 Construction of PPV's

Finally, let us see how a PPV for preference vector  $u$  can be constructed directly from partial vectors and the hubs skeleton using the Hubs Equation. (Construction of a single hub vector is a specialization of the algorithm outlined here.) Let  $u = \alpha_1 p_1 + \dots + \alpha_z p_z$  be a preference vector, where  $p_i \in H$  for  $1 \leq i \leq z$ . Let  $Q \subseteq H$ , and let

$$r_u(h) = \sum_{i=1}^z \alpha_i (r_{p_i}(h) - c \cdot x_{p_i}(h)) \quad (19)$$

which can be computed from the hubs skeleton. Then the PPV  $v$  for  $u$  can be constructed as

$$v = \sum_{i=1}^z \alpha_i (r_{p_i} - r_{p_i}^H) + \frac{1}{c} \sum_{\substack{h \in Q \\ r_u(h) > 0}} r_u(h) [(r_h - r_h^H) - cx_h] \quad (20)$$

Both the terms  $(r_{p_i} - r_{p_i}^H)$  and  $(r_h - r_h^H)$  are partial vectors, which we assume have been precomputed. The term  $cx_h$  represents a simple subtraction from  $(r_h - r_h^H)$ . If  $Q = H$ , then (20) represents a full construction of  $v$ . However, for some applications, it may suffice to use only parts of the hubs skeleton to compute  $v$  to less precision. For example, we can take  $Q$  to be the  $m$  hubs  $h$  for which  $r_u(h)$  is highest. Experimentation with this scheme is discussed in Section 6.3. Alternatively, the result can be improved incrementally (e.g., as time permits) by using a small subset  $Q$  each time and accumulating the results.

## 6. EXPERIMENTS

We performed experiments using real web data from Stanford's *WebBase* [6], a crawl of the web containing 120 million pages. Since the iterative computation of PageRank is unaffected by *leaf pages* (i.e., those with no out-neighbors), they can be removed from the graph and added back in after the computation [11]. After removing leaf pages, the graph consisted of 80 million pages.

Both the web graph and the intermediate results  $(D_k[*], E_k[*])$  were too large to fit in main memory, and a partitioning strategy, based on that presented in [4], was used to divide the computation into portions that can be carried out in memory. Specifically, the set of pages  $V$  was partitioned into  $k$  arbitrary sets  $P_1, \dots, P_k$  of equal size ( $k = 10$  in our experiments). The web graph, represented as an edge-list  $E$ , is partitioned into  $k$  chunks  $E_i$  ( $1 \leq i \leq k$ ), where  $E_i$  contains all edges  $\langle p, q \rangle$  for which  $p \in P_i$ . Intermediate results  $D_k[p]$  and  $E_k[p]$  were represented together as a list  $L_k[p] = \langle (q_1, d_1, e_1), (q_2, d_2, e_2), \dots \rangle$  where  $D_k[p](q_z) = d_z$  and  $E_k[p](q_z) = e_z$ , for  $z = 1, 2, \dots$ . Only pages  $q_z$  for which either  $d_z > 0$  or  $e_z > 0$  were included. The set of intermediate results  $L_k[*]$  was partitioned into  $k^2$  chunks  $L_k^{i,j}[*]$ , so that  $L_k^{i,j}[p]$  contains triples  $(q_z, d_z, e_z)$  of  $L_k[p]$  for which  $p \in P_i$  and  $q_z \in P_j$ . In each of the algorithms for computing partial quantities, only a single column  $L_k^{*,j}[*]$  was kept in memory at any one time, and part of the next-iteration results  $L_{k+1}[*]$  were computed by successively reading in individual blocks of the graph or intermediate results as appropriate. Each iteration requires only one linear scan of the intermediate results and web graph, except for repeated squaring, which does not use the web graph explicitly.

### 6.1 Computing Partial Vectors

For comparison, we computed both (full) hub vectors and partial vectors for various sizes of  $H$ , using the selective expansion al-

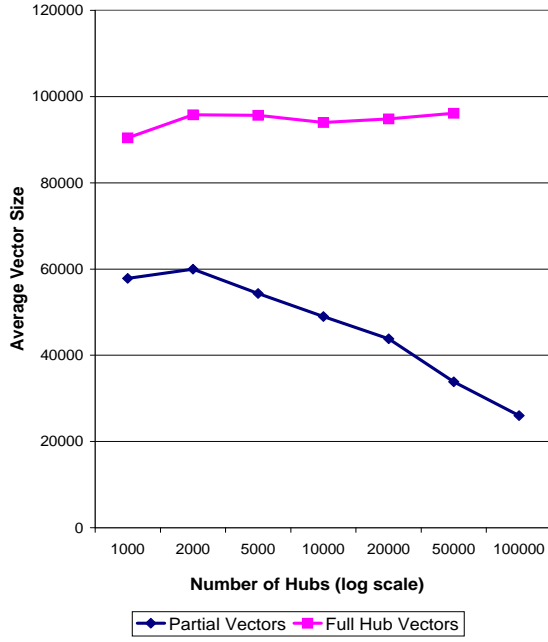


Figure 2: Average Vector Size vs. Number of Hubs

gorithm with  $Q_k(p) = V$  (full hub vectors) and  $Q_k(p) = V - H$  (partial vectors). As discussed in Section 4.4.2, we found the partial vectors approach to be much more effective when  $H$  contains high-PageRank pages rather than random pages. In our experiments  $H$  ranged from the top 1000 to top 100,000 pages with the highest PageRank. The constant  $c$  was set to 0.15.

To evaluate the performance and scalability of our strategy independently of implementation and platform, we focus on the size of the results rather than computation time, which is linear in the size of the results. Because of the number of trials we had to perform and limitations on resources, we computed results only up to 6 iterations, for  $|H|$  up to 100,000. Figure 2 plots the average size of (full) hub vectors and partial vectors (recall that size is the number of nonzero entries), as computed after 6 iterations of the selective expansion algorithm, which for computing full hub vectors is equivalent to the basic dynamic programming algorithm. Note that the x-axis plots  $|H|$  in logarithmic scale.

Experiments were run using a 1.4 gigahertz CPU on a machine with 3.5 gigabytes of memory. For  $|H| = 50,000$ , the computation of full hub vectors took about 2.8 seconds per vector, and about 0.33 seconds for each partial vector. We were unable to compute full hub vectors for  $|H| = 100,000$  due to the time required, although the average vector size is expected not to vary significantly with  $|H|$  for full hub vectors. In Figure 2 we see that the reduction in size from using our technique becomes more significant as  $|H|$  increases, suggesting that our technique scales well with  $|H|$ .

## 6.2 Computing the Hubs Skeleton

We computed the hubs skeleton for  $|H| = 10,000$  by running the selective expansion algorithm for 6 iterations using  $Q_k(p) = H$ , and then running the repeated squaring algorithm for 10 iterations (Section 5.3.2), where  $Q_k(p)$  is chosen to be the top 50 entries under the top- $m$  scheme (Section 5.2.2). The average size of the hubs skeleton is 9021 entries. Each iteration of the repeated squaring algorithm took about an hour, a cost that depends only on  $|H|$  and is constant with respect to the precision to which the partial

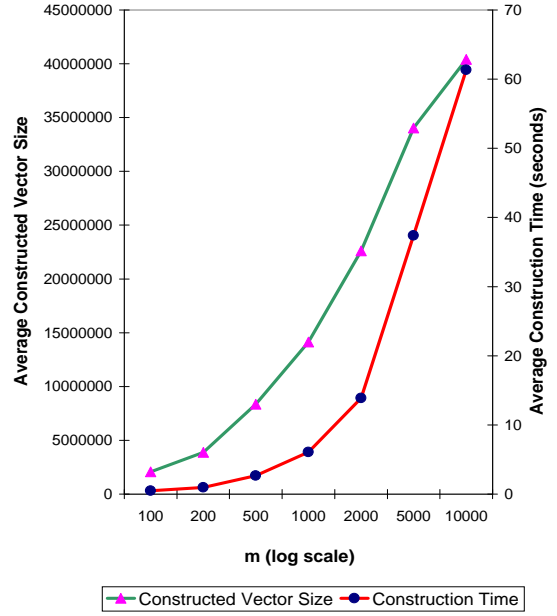


Figure 3: Construction Time and Size vs. Hubs Skeleton Portion ( $m$ )

vectors are computed.

## 6.3 Constructing Hub Vectors from Partial Vectors

Next we measured the construction of (full) hub vectors from partial vectors and the hubs skeleton. Note that in practice we may construct PPV's directly from partial vectors, as discussed in Section 5.4. However, performance of the construction would depend heavily on the user's preference vector. We consider hub vector computation because it better measures the performance benefits of our partial vectors approach.

As suggested in Section 4.3, the precision of the hub vectors constructed from partial vectors can be varied at query time according to application and performance demands. That is, instead of using the entire set  $r_p(H)$  in the construction of  $r_p$ , we can use only the highest  $m$  entries, for  $m \leq |H|$ . Figure 3 plots the average size and time required to construct a full hub vector from partial vectors in memory versus  $m$ , for  $|H| = 10,000$ . Results are averaged over 50 randomly-chosen hub vectors. Note that the x-axis is in logarithmic scale.

Recall from Section 6.1 that the partial vectors from which the hubs vector is constructed were computed using 6 iterations, limiting the precision. Thus, the error values in Figure 3 are roughly 16% (ranging from 0.166 for  $m = 100$  to 0.163 for  $m = 10,000$ ). Nonetheless, this error is much smaller than that of the iteration-6 full hub vectors computed in Section 6.1, which have error  $(1 - c)^6 = 38\%$ . Note, however, that the size of a vector is a better indicator of precision than the magnitude, since we are usually most interested in the number of pages with nonzero entries in the distribution vector. An iteration-6 full hub vector (from Section 6.1) for page  $p$  contains nonzero entries for pages at most 6 links away from  $p$ , 93,993 pages on average. In contrast, from Figure 3 we see that a hub vector containing 14 million nonzero entries can be constructed from partial vectors in 6 seconds.



## 7. RELATED WORK

The use of personalized PageRank to enable personalized web search was first proposed in [11], where it was suggested as a modification of the global PageRank algorithm, which computes a universal notion of importance. The computation of (personalized) PageRank scores was not addressed beyond the naive algorithm.

In [5], personalized PageRank scores were used to enable “topic-sensitive” web search. Specifically, precomputed hub vectors corresponding to broad categories in *Open Directory* were used to bias importance scores, where the vectors and weights were selected according to the text query. Experiments in [5] concluded that the use of personalized PageRank scores can improve web search, but the number of hub vectors used was limited to 16 due to the computational requirements, which were not addressed in that work. Scaling the number of hub pages beyond 16 for finer-grained personalization is a direct application of our work.

Another technique for computing web-page importance, *HITS*, was presented in [9]. In *HITS*, an iterative computation similar in spirit to PageRank is applied at query time on a subgraph consisting of pages matching a text query and those “nearby”. Personalizing based on user-specified web pages (and their linkage structure in the web graph) is not addressed by *HITS*. Moreover, the number of pages in the subgraphs used by *HITS* (order of thousands) is much smaller than that we consider in this paper (order of millions), and the computation from scratch at query time makes the *HITS* approach difficult to scale.

Another algorithm that uses query-dependent importance scores to improve upon a global version of importance was presented in [12]. Like *HITS*, it first restricts the computation to a subgraph derived from text matching. (Personalizing based on user-specified web pages is not addressed.) Unlike *HITS*, [12] suggested that importance scores be precomputed offline for every possible text query, but the enormous number of possibilities makes this approach difficult to scale.

The concept of using “hub nodes” in a graph to enable partial computation of solutions to the shortest-path problem was used in [3] in the context of database search. That work deals with searches within databases, and on a scale far smaller than that of the web.

Some system aspects of (global) PageRank computation were addressed in [4]. The disk-based data-partitioning strategy used in the implementation of our algorithm is adopted from that presented therein.

Finally, the concept of inverse  $P$ -distance used in this paper is based on the concept of expected- $f$  distance introduced in [8], where it was presented as an intuitive model for a similarity measure in graph structures.

## 8. SUMMARY

We have addressed the problem of scaling personalized web search:

- We started by identifying a linear relationship that allows personalized PageRank vectors to be expressed as a linear combination of *basis vectors*. Personalized vectors corresponding to arbitrary preference sets drawn from a *hub set*  $H$  can be constructed quickly from the set of precomputed basis *hub vectors*, one for each hub  $h \in H$ .
- We laid the mathematical foundations for constructing hub vectors efficiently by relating personalized PageRank scores to *inverse  $P$ -distances*, an intuitive notion of distance in arbitrary directed graphs. We used this notion of distance to identify interrelationships among basis vectors.
- We presented a method of encoding hub vectors as *partial*

*vectors* and the *hubs skeleton*. Redundancy is minimized under this representation: each partial vector for a hub page  $p$  represents the part of  $p$ 's hub vector unique to itself, while the skeleton specifies how partial vectors are assembled into full vectors.

- We presented algorithms for computing basis vectors, and showed how they can be modified to compute partial vectors and the hubs skeleton efficiently.
- We ran experiments on real web data showing the effectiveness of our approach. Results showed that our strategy results in significant resource reduction over full vectors, and scales well with  $|H|$ , the degree of personalization.

## 9. ACKNOWLEDGMENT

The authors thank Taher Haveliwala for many useful discussions and extensive help with implementation.

## 10. REFERENCES

- [1] <http://www.google.com>.
- [2] <http://dmoz.org>.
- [3] R. Goldman, N. Shivakumar, S. Venkatasubramanian, and H. Garcia-Molina. Proximity search in databases. In *Proceedings of the Twenty-Fourth International Conference on Very Large Databases*, New York, New York, Aug. 1998.
- [4] T. H. Haveliwala. Efficient computation of PageRank. Technical report, Stanford University Database Group, 1999. <http://dbpubs.stanford.edu/pub/1999-31>.
- [5] T. H. Haveliwala. Topic-sensitive PageRank. In *Proceedings of the Eleventh International World Wide Web Conference*, Honolulu, Hawaii, May 2002.
- [6] J. Hirai, S. Raghavan, A. Paepcke, and H. Garcia-Molina. WebBase: A repository of web pages. In *Proceedings of the Ninth International World Wide Web Conference*, Amsterdam, Netherlands, May 2000. <http://www-diglib.stanford.edu/~testbed/doc2/WebBase/>.
- [7] G. Jeh and J. Widom. Scaling personalized web search. Technical report, Stanford University Database Group, 2002. <http://dbpubs.stanford.edu/pub/2002-12>.
- [8] G. Jeh and J. Widom. SimRank: A measure of structural-context similarity. In *Proceedings of the Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, Edmonton, Alberta, Canada, July 2002.
- [9] J. M. Kleinberg. Authoritative sources in a hyperlinked environment. In *Proceedings of the Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, San Francisco, California, Jan. 1998.
- [10] R. Motwani and P. Raghavan. *Randomized Algorithms*. Cambridge University Press, United Kingdom, 1995.
- [11] L. Page, S. Brin, R. Motwani, and T. Winograd. The PageRank citation ranking: Bringing order to the Web. Technical report, Stanford University Database Group, 1998. <http://citeseer.nj.nec.com/368196.html>.
- [12] M. Richardson and P. Domingos. The intelligent surfer: Probabilistic combination of link and content information in PageRank. In *Proceedings of Advances in Neural Information Processing Systems 14*, Cambridge, Massachusetts, Dec. 2002.