

A REORDERING FOR THE PAGERANK PROBLEM*

AMY N. LANGVILLE[†] AND CARL D. MEYER[‡]

Abstract. We describe a reordering particularly suited to the PageRank problem, which reduces the computation of the PageRank vector to that of solving a much smaller system and then using forward substitution to get the full solution vector. We compare the theoretical rates of convergence of the original PageRank algorithm to that of the new reordered PageRank algorithm, showing that the new algorithm can do no worse than the original algorithm. We present results of an experimental comparison on five datasets, which demonstrate that the reordered PageRank algorithm can provide a speedup of as much as a factor of 6. We also note potential additional benefits that result from the proposed reordering.

Key words. Markov chains, PageRank, reorderings, power method, convergence, stationary vector, dangling nodes

AMS subject classifications. 65B99, 65F10, 65C40, 60J22, 65F15, 65F50

DOI. 10.1137/040607551

1. Introduction. It is well known that many subsets of the web contain a large proportion of dangling nodes, i.e., webpages with no outlinks. Dangling nodes can result from many sources: a page containing an image or a PostScript or pdf file; a page of data tables; or a page whose links have yet to be crawled by the search engine's spider. These dangling nodes can present philosophical, storage, and computational issues for a search engine such as Google that employs a ranking system for ordering retrieved webpages.

Let us introduce some notation. The hyperlink structure of the web defines a directed graph, which can be expressed as a sparse matrix. The founders of Google, Brin and Page, defined the elements of a hyperlink matrix \mathbf{H} as $h_{ij} = 1/|O_i|$ if there exists a hyperlink from page i to page j , and 0 otherwise. The scalar $|O_i|$ is the number of outlinks from page i . Thus, the nonzero rows of \mathbf{H} sum to 1. The matrix \mathbf{H} contains a $\mathbf{0}^T$ row for each dangling node. Brin and Page define the PageRank vector, a vector holding the global measure of importance for each page, to be the stationary vector for a Markov chain related to \mathbf{H} [3, 4]. This definition is intuitive, as the stationary vector gives the long-run proportion of time the chain will spend in each state. For Brin and Page's application, the stationary elements represent the long-run proportion of time a random web surfer will spend on a page in the web, and thus provide a measure of a page's importance or popularity. One problem with the hyperlink matrix \mathbf{H} created strictly from the web's structure is that it is not stochastic, and a Markov chain is defined only for stochastic matrices. To remedy this and create the stochastic matrix called \mathbf{S} , Brin and Page suggest replacing each $\mathbf{0}^T$ row (corresponding to a dangling node) of the sparse hyperlink matrix with a dense vector. The original fix for the dangling node problem used a uniform vector \mathbf{e}^T/n (\mathbf{e} is the vector of all ones), which was later replaced by the more general probability vector

*Received by the editors April 29, 2004; accepted for publication (in revised form) April 20, 2005; published electronically February 21, 2006.

<http://www.siam.org/journals/sisc/27-6/60755.html>

[†]Department of Mathematics, College of Charleston, Charleston, SC 29424 (langvillea@cofc.edu).

[‡]Center for Research in Scientific Computation, Department of Mathematics, North Carolina State University, Raleigh, NC 27695-8205 (meyer@math.ncsu.edu). This author's research was supported in part by NSF CCR-ITR-0113121 and NSF DMS 9714811.

\mathbf{v}^T , known as the personalization or teleportation vector. Of course, if this suggestion were to be implemented explicitly, storage requirements would increase dramatically. Instead, the stochasticity fix can be modeled implicitly with the construction of one vector denoted \mathbf{a} . Element $a_i = 1$ if row i of \mathbf{H} corresponds to a dangling node, and 0 otherwise. Then \mathbf{S} can be written as a rank-one update of \mathbf{H} ; $\mathbf{S} = \mathbf{H} + \mathbf{a}\mathbf{v}^T$.

Before the existence of the stationary PageRank vector could be guaranteed, Brin and Page had to make one final adjustment to the hyperlink matrix. Because the web is reducible, their original PageRank algorithm would often get caught in a rank sink, causing convergence problems. To remedy this, and thereby guarantee the existence and uniqueness of the PageRank vector, Brin and Page added another rank-one update, this time an irreducibility adjustment in the form of a dense perturbation matrix $\mathbf{e}\mathbf{v}^T$ that creates direct connections between each page. The stochastic, irreducible matrix called the Google matrix \mathbf{G} is given by

$$\begin{aligned}\mathbf{G} &= \alpha \mathbf{S} + (1 - \alpha) \mathbf{e}\mathbf{v}^T \\ &= \alpha \mathbf{H} + \alpha \mathbf{a}\mathbf{v}^T + (1 - \alpha) \mathbf{e}\mathbf{v}^T \\ &= \alpha \mathbf{H} + (\alpha \mathbf{a} + (1 - \alpha) \mathbf{e})\mathbf{v}^T,\end{aligned}$$

where $0 < \alpha < 1$. Brin and Page, and consequently Google, then use the power method to compute the unique stationary vector of this enormous Markov chain, thereby producing their famous PageRank vector.

We now turn to the philosophical issue of the presence of dangling nodes. In one of their early papers [2], Brin et al. report that they “often remove dangling nodes during the computation of PageRank, then add them back in after the PageRanks have converged.” From this vague statement it is hard to say exactly how Brin and Page compute PageRank for the dangling nodes. However, the removal of dangling nodes at any time during the power method does not make intuitive sense. Some dangling nodes should receive high PageRank. For example, a very authoritative pdf file could have many inlinks from respected sources and thus should receive a high PageRank. Simply removing the dangling nodes biases the PageRank vector unjustly. (See [5] and [13] for additional arguments against removal of dangling nodes.) Further, incorporating dangling nodes into the PageRank power method is very simple and inexpensive. The power method treats PageRank as an eigenvector problem and follows the iterative formula

$$\begin{aligned}\mathbf{x}^{(k)T} &= \mathbf{x}^{(k-1)T} \mathbf{G} \\ &= \alpha \mathbf{x}^{(k-1)T} \mathbf{S} + (1 - \alpha) \mathbf{x}^{(k-1)T} \mathbf{e}\mathbf{v}^T \\ &= \alpha \mathbf{x}^{(k-1)T} \mathbf{S} + (1 - \alpha) \mathbf{v}^T \\ (1) \quad &= \alpha \mathbf{x}^{(k-1)T} \mathbf{H} + \alpha \mathbf{x}^{(k-1)T} \mathbf{a}\mathbf{v}^T + (1 - \alpha) \mathbf{v}^T,\end{aligned}$$

since $\mathbf{x}^{(k-1)T}$ is a probability vector, and thus, $\mathbf{x}^{(k-1)T} \mathbf{e} = 1$. This shows that the power method applied to \mathbf{G} can be implemented with vector-matrix multiplications on the extremely sparse \mathbf{H} , and \mathbf{G} and \mathbf{S} are never formed or stored. Since the vector-matrix multiplications involve enormous entities, an important question to ask is how many iterations can we expect until the power method converges to the PageRank vector. It has been proven that the asymptotic rate of convergence of the power method applied to the Google matrix (after the stochasticity and irreducibility adjustments) is the rate at which $\alpha^k \rightarrow 0$. Since Google uses $\alpha = .85$, one can expect roughly 114 power iterations to give a convergence tolerance (as measured by the norm

of the difference of successive iterates) of less than 10^{-8} . Further, this means that on the order of $O(114 \cdot nnz(\mathbf{H}))$ operations must be performed to compute the PageRank vector with that prescribed tolerance, where $nnz(\mathbf{H})$ is the number of nonzeros in \mathbf{H} . Since \mathbf{H} is so sparse, this is not a prohibitive cost, but it does nevertheless take days of computation when n , the number of pages in the index, is billions.

At this point, we have arrived at the computational issue associated with the dangling nodes. Since the dangling nodes form identical rows in the \mathbf{H} matrix, they can be conveniently lumped into one class and the theory of lumpable and aggregated Markov chains can be used to compute the PageRank vector quickly. Lee, Golub, and Zenios were the first to notice and exploit the structure provided by the dangling nodes [13]. Their iterative power method implementation on aggregated Markov chains led to drastic improvements in the computation of PageRank. On some sample datasets dangling nodes make up over 80% of the webpages, meaning the Lee–Golub–Zenios algorithm can compute PageRank with a factor of 5 speedup, since the aggregated chain is roughly 1/5 the size of the full chain.

In this paper, we present an analogous formulation of the Lee–Golub–Zenios algorithm (section 2). Our formulation uses a linear system formulation rather than a Markov chain formulation. Then we extend the idea of exploiting the dangling nodes to reduce computation, thereby producing a convenient reordering of the Google matrix, which has several advantageous properties (section 3). In section 4, we present the results of our reordering algorithm on several sample datasets.

2. A linear system formulation for exploiting dangling nodes. Although the size of the Google problem makes the power method with its eigenvector formulation ($\pi^T = \pi^T \mathbf{G}$ and $\pi^T \mathbf{e} = 1$) one of the few practical solution methods for obtaining the PageRank vector π^T , there are other formulations of the problem that are theoretically possible. The following theorem shows that a related linear system exists.

THEOREM 2.1 (linear system for Google problem). *Solving the linear system*

$$(2) \quad \mathbf{x}^T(\mathbf{I} - \alpha\mathbf{H}) = \mathbf{v}^T$$

and letting $\pi^T = \mathbf{x}^T / \mathbf{x}^T \mathbf{e}$ produces the PageRank vector.

Proof. π^T is the PageRank vector if it satisfies $\pi^T \mathbf{G} = \pi^T$ and $\pi^T \mathbf{e} = 1$. Clearly, $\pi^T \mathbf{e} = 1$. Showing $\pi^T \mathbf{G} = \pi^T$ is equivalent to showing $\pi^T(\mathbf{I} - \mathbf{G}) = \mathbf{0}^T$, which is equivalent to showing $\mathbf{x}^T(\mathbf{I} - \mathbf{G}) = \mathbf{0}^T$:

$$\begin{aligned} \mathbf{x}^T(\mathbf{I} - \mathbf{G}) &= \mathbf{x}^T(\mathbf{I} - \alpha\mathbf{H} - \alpha\mathbf{a}\mathbf{v}^T - (1 - \alpha)\mathbf{e}\mathbf{v}^T) \\ &= \mathbf{x}^T(\mathbf{I} - \alpha\mathbf{H}) - \mathbf{x}^T(\alpha\mathbf{a} + (1 - \alpha)\mathbf{e})\mathbf{v}^T \\ &= \mathbf{v}^T - \mathbf{v}^T = \mathbf{0}^T. \end{aligned}$$

The above line results from the fact that $\mathbf{x}^T(\alpha\mathbf{a} + (1 - \alpha)\mathbf{e})\mathbf{v}^T = 1$ because

$$\begin{aligned} 1 &= \mathbf{v}^T \mathbf{e} \\ &= \mathbf{x}^T(\mathbf{I} - \alpha\mathbf{H})\mathbf{e} \\ &= \mathbf{x}^T \mathbf{e} - \alpha\mathbf{x}^T \mathbf{H}\mathbf{e} \\ &= \mathbf{x}^T \mathbf{e} - \alpha\mathbf{x}^T(\mathbf{e} - \mathbf{a}) \\ &= (1 - \alpha)\mathbf{x}^T \mathbf{e} + \alpha\mathbf{x}^T \mathbf{a}. \quad \square \end{aligned}$$

Thus, PageRank is both the stationary distribution of a Markov chain and the solution of a linear system. While these two views are inherently linked, they suggest

different methods for finding the PageRank vector. Future advances in computation will likely require a thorough understanding and use of both views.

As a result, we examine the linear system view, which has received much less attention. The coefficient matrix $(\mathbf{I} - \alpha\mathbf{H})$ of the linear system has many nice properties, which were proven in [12]. Some that are relevant for this paper are as follows:

- $(\mathbf{I} - \alpha\mathbf{H})$ is nonsingular.
- The row sums of $(\mathbf{I} - \alpha\mathbf{H})^{-1}$ are equal to 1 for the dangling nodes and less than or equal to $1/(1 - \alpha)$ for the nondangling nodes.
- The row of $(\mathbf{I} - \alpha\mathbf{H})^{-1}$ corresponding to dangling node i is \mathbf{e}_i^T , where \mathbf{e}_i is the i th column of the identity matrix.

The last property makes the computation of the PageRank vector especially efficient. Suppose the rows and columns of \mathbf{H} are permuted (i.e., the indices are reordered) so that the rows corresponding to dangling nodes are at the bottom of the matrix:

$$(3) \quad \mathbf{H} = \begin{array}{cc} & \begin{matrix} ND & D \end{matrix} \\ \begin{matrix} ND \\ D \end{matrix} & \begin{pmatrix} \mathbf{H}_{11} & \mathbf{H}_{12} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \end{array},$$

where ND is the set of nondangling nodes and D is the set of dangling nodes. Then the coefficient matrix in the sparse linear system formulation becomes

$$(\mathbf{I} - \alpha\mathbf{H}) = \begin{pmatrix} \mathbf{I} - \alpha\mathbf{H}_{11} & -\alpha\mathbf{H}_{12} \\ \mathbf{0} & \mathbf{I} \end{pmatrix},$$

and the inverse of this matrix is

$$(\mathbf{I} - \alpha\mathbf{H})^{-1} = \begin{pmatrix} (\mathbf{I} - \alpha\mathbf{H}_{11})^{-1} & \alpha(\mathbf{I} - \alpha\mathbf{H}_{11})^{-1}\mathbf{H}_{12} \\ \mathbf{0} & \mathbf{I} \end{pmatrix}.$$

Therefore, the unnormalized PageRank vector $\mathbf{x}^T = \mathbf{v}^T(\mathbf{I} - \alpha\mathbf{H})^{-1}$ can be written as

$$\mathbf{x}^T = (\mathbf{v}_1^T(\mathbf{I} - \alpha\mathbf{H}_{11})^{-1} \mid \alpha\mathbf{v}_1^T(\mathbf{I} - \alpha\mathbf{H}_{11})^{-1}\mathbf{H}_{12} + \mathbf{v}_2^T),$$

where the personalization vector \mathbf{v}^T has been partitioned into nondangling (\mathbf{v}_1^T) and dangling (\mathbf{v}_2^T) sections. In summary, we now have an algorithm that computes the PageRank vector using only the nondangling portion of the web, exploiting the rank-one structure of the dangling node fix.

ALGORITHM 1.

1. Reorder the states of the Markov chain so that the reordered matrix has the structure of (3).
 2. Solve for \mathbf{x}_1^T in $\pi_1^T(\mathbf{I} - \alpha\mathbf{H}_{11}) = \mathbf{v}_1^T$.
 3. Compute $\mathbf{x}_2^T = \alpha\pi_1^T\mathbf{H}_{12} + \mathbf{v}_2^T$.
 4. Normalize $\pi^T = [\mathbf{x}_1^T \mid \mathbf{x}_2^T] / \|\mathbf{x}_1^T \mid \mathbf{x}_2^T\|_1$.
-

This algorithm is much simpler and cleaner than, but equivalent to, the specialized iterative method proposed by Lee, Golub, and Zenios [13], which exploits the dangling nodes to reduce computation of the PageRank vector, sometimes by a factor of 5.

3. A PageRank algorithm based on a reordering of the Google matrix.

3.1. The reordered PageRank algorithm. The linear system formulation of section 2 leads to a deeper examination of the structure of the Google matrix \mathbf{H} . Since the presence of zero rows in the matrix is so advantageous, we hope to find more zero rows in the submatrix \mathbf{H}_{11} , which is needed to solve the system in step 1 of Algorithm 1. This process of locating zero rows can be repeated recursively on smaller and smaller submatrices of \mathbf{H} , continuing until a submatrix is created that has no zero rows. The result of this process is a decomposition of the \mathbf{H} matrix that looks like Figure 1. This process amounts to a simple reordering of the indices of the Markov chain. The top part shows the original \mathbf{H} matrix, and the bottom part is the reordered matrix according to the recursive dangling node idea. The dataset `CA.dat` (available from <http://www.cs.cornell.edu/Courses/cs685/2002fa/>) is a typical subset of the web. It contains 9,664 nodes and 16,773 links pertaining to the query topic of “california.”

In general, after this symmetric reordering, the hyperlink matrix \mathbf{H} has the following structure:

$$(4) \quad \mathbf{H} = \begin{pmatrix} \mathbf{H}_{11} & \mathbf{H}_{12} & \mathbf{H}_{13} & \cdots & \mathbf{H}_{1b} \\ & \mathbf{0} & \mathbf{H}_{23} & \cdots & \mathbf{H}_{2b} \\ & & \mathbf{0} & \cdots & \mathbf{H}_{3b} \\ & & & \ddots & \\ & & & & \mathbf{0} \end{pmatrix},$$

where $b \geq 2$ is the number of square diagonal blocks in the reordered matrix. Therefore, the coefficient matrix of the linear system formulation of the PageRank problem (2) has the following structure:

$$(5) \quad (\mathbf{I} - \alpha\mathbf{H}) = \begin{pmatrix} \mathbf{I} - \alpha\mathbf{H}_{11} & -\alpha\mathbf{H}_{12} & -\alpha\mathbf{H}_{13} & \cdots & -\alpha\mathbf{H}_{1b} \\ & \mathbf{I} & -\alpha\mathbf{H}_{23} & \cdots & -\alpha\mathbf{H}_{2b} \\ & & \mathbf{I} & \cdots & -\alpha\mathbf{H}_{3b} \\ & & & \ddots & \\ & & & & \mathbf{I} \end{pmatrix}.$$

As a result, the PageRank system in (2) after reordering can be solved by forward substitution. The only system that must be solved directly is the first subsystem, $\mathbf{x}_1^T(\mathbf{I} - \alpha\mathbf{H}_{11}) = \mathbf{v}_1^T$, where \mathbf{x}^T and \mathbf{v}^T have also been partitioned according to the number and size of the blocks. The remaining subvectors of \mathbf{x}^T are computed quickly and efficiently by forward substitution. In the `CA.dat` example, a $2,622 \times 2,622$ system can be solved instead of the full $9,664 \times 9,664$ system, or even the once-reduced $5,132 \times 5,132$ system of Algorithm 1. The reordered PageRank algorithm is an extension of the dangling node method of section 2. This reordered PageRank idea can also be written in a Markov chain formulation, thereby extending the ideas in [13]. The steps of the reordered PageRank (in its linear system formulation) are enumerated in Algorithm 2.

ALGORITHM 2.

1. Reorder the states of the Markov chain so that the reordered matrix has the structure of (4).
 2. Solve for \mathbf{x}_1^T in $\mathbf{x}_1^T(\mathbf{I} - \alpha\mathbf{H}_{11}) = \mathbf{v}_1^T$.
 3. For $i = 2$ to b , compute $\mathbf{x}_i^T = \alpha \sum_{j=1}^{i-1} \mathbf{x}_j^T \mathbf{H}_{ji} + \mathbf{v}_i^T$.
 4. Normalize $\boldsymbol{\pi}^T = [\mathbf{x}_1^T \ \mathbf{x}_2^T \ \cdots \ \mathbf{x}_b^T] / \|\mathbf{x}_1^T \ \mathbf{x}_2^T \ \cdots \ \mathbf{x}_b^T\|_1$.
-

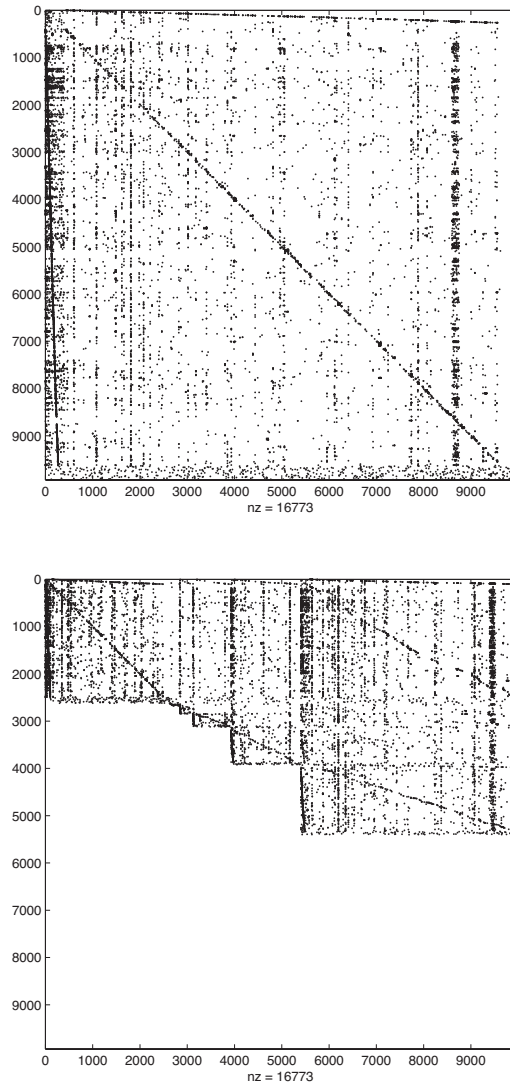


FIG. 1. Original \mathbf{H} (top) and reordered \mathbf{H} matrix (bottom) for *CA.dat*.

3.2. Analysis of the reordered PageRank algorithm. We discuss each step of Algorithm 2. The reordering, step 1 of Algorithm 2, can be accomplished with a simple recursive scheme that successively locates zero rows of submatrices within \mathbf{H} , and thus amounts to the work of a few power iterations on the full matrix. However, the memory manipulation required to successively reorder the submatrices can make this step time consuming and potentially of less value than Algorithm 1 (see the experiments of section 4). The forward substitution step of Algorithm 2, step 3, amounts to slightly less than the work of one power iteration. Of course, step 4 is trivial. Thus, the small system solve of step 2 is the most time-consuming step. Table 1 in section 4 gives an enlightening numerical comparison of the amount of work spent on each step.

Some hyperlink matrices \mathbf{H} can be particularly suited to the reordered PageRank algorithm. In this case, large sets of zero rows are consecutively extracted, and the remaining \mathbf{H}_{11} block is very small. In the worst case, $b = 2$ and $\text{nnz}(\mathbf{H}_{11})/\text{nnz}(\mathbf{H})$ is close to 1. Regardless, the smaller system $\mathbf{x}_1^T(\mathbf{I} - \alpha\mathbf{H}_{11}) = \mathbf{v}_1^T$ can be solved by any appropriate direct or iterative method. However, we have found that the sparsity of the \mathbf{H}_{11} matrix makes a simple iterative method very attractive. The Jacobi method [14] is especially simple due to the properties of the \mathbf{H}_{11} matrix. The diagonal elements of \mathbf{H}_{11} are 0, since by convention no webpage links to itself. This means that the diagonal matrix \mathbf{D} involved in the splitting of the coefficient matrix $(\mathbf{I} - \alpha\mathbf{H}_{11})$ is the identity matrix \mathbf{I} . Therefore, the Jacobi iterates are

$$(6) \quad \mathbf{x}^{(k)T} = \alpha\mathbf{x}^{(k-1)T}\mathbf{H}_{11} + \mathbf{v}_1^T.$$

This stationary iterative process converges to a unique solution for all starting iterates because $\rho(\alpha\mathbf{H}_{11}) \leq 1$. (\mathbf{H}_{11} is substochastic.) The asymptotic rate of convergence of the iterative reordered PageRank process is the same as or better than the asymptotic rate of convergence of the standard PageRank power method. In fact, the reordered PageRank method will require about $O(\frac{\log(\tau)}{\log(\alpha)} \cdot \text{nnz}(\mathbf{H}_{11}))$ operations to reach a convergence tolerance of τ . In summary, Algorithm 2, the reordered PageRank method, requires roughly $\frac{\text{nnz}(\mathbf{H}_{11})}{\text{nnz}(\mathbf{H})}\%$ of the time that the original PageRank method requires. The Jacobi iterates of (6), like the power iterates of the original PageRank algorithm, also allow for easy parallelization, which is essential due to the scale of the matrices involved in practical applications. Further, acceleration methods [1, 6, 8, 9, 10], such as extrapolation and preconditioners, can be applied to the small \mathbf{H}_{11} system to achieve even greater speedups.

4. Experiments with the reordering method. We compare the reordered PageRank algorithm to the original PageRank algorithm, experimenting with four datasets. The first two datasets are small, containing fewer than 10,000 pages. The remaining two datasets are much larger, with one containing over 325,000 pages and the other over 450,000 pages. Each dataset is a different type of subset of the web. The **CA.dat** dataset contains 9,664 pages pertaining to the query topic of “california.” **NCSU.dat** is a crawl that started at the North Carolina State homepage and stopped once 10,000 pages were reached. Most of the pages are internal to the NCSU domain, but some are external. The dataset **ND.dat** contains 325,729 pages within the Notre Dame domain. Finally, **SU450k.dat** contains 451,237 pages from a 2001 crawl of the Stanford domain by the Stanford WebBase project.

The experiments in this section were conducted on a 867 MHz Mac G4 with 1.5 GB of RAM. We used $\alpha = .9$ as the scaling factor and $\tau = 10^{-10}$ as the convergence tolerance. The original PageRank algorithm and both reordered PageRank algorithms (Algorithms 1 and 2) were run on each dataset. The results are summarized in Table 1.

The reordered PageRank algorithm beats the original PageRank algorithm by a factor of almost 6 on some datasets. The speedup depends on the ratio $\frac{\text{nnz}(\mathbf{H})}{\text{nnz}(\mathbf{H}_{11})}$, which depends on the properties of a particular dataset and the number of zero rows that can be successively squeezed out of the submatrices. The reordered PageRank algorithm is guaranteed to outperform the original PageRank method as long as some dangling nodes are present, i.e., $b \geq 2$.

Unfortunately, Table 1 clearly shows that Algorithm 2 is not necessarily an improvement over Algorithm 1 because the cost of memory management required by the reordering step often offsets the decrease in the size of the system solve. Nevertheless,

TABLE 1
Comparison of original PageRank and reordered PageRank algorithms.

		CA.dat	NCSU.dat	ND.dat	SU450k.dat
Original PageRank	$n(\mathbf{H})$	10K	10K	325K	450K
	$nnz(\mathbf{H})$	16K	101K	1,497K	1,082K
	Iterations	178	162	166	164
	Time (sec.)	3.30	5.28	126.71	175.55
Reordered PR Algorithm 1 ($b = 2$)	b	2	2	2	2
	$n(\mathbf{H}_{11})$	5K	7K	138K	137K
	$nnz(\mathbf{H}_{11})$	8K	81K	1,208K	307K
	Iterations	170	162	166	145
	Time: Step 1	.14	.32	2.37	2.84
	Step 2	.78	3.96	62.47	26.84
	Step 3	.10	.13	.66	1.16
	Total	1.02	4.41	65.50	30.89
	Speedup	3.2	1.2	1.9	5.7
Reordered PR Algorithm 2 ($b > 2$)	b	7	5	18	12
	$n(\mathbf{H}_{11})$	2.6K	6.8K	127K	84K
	$nnz(\mathbf{H}_{11})$	5K	79K	1,191K	267K
	Iterations	169	160	170	145
	Time: Step 1	.24	.58	22.74	12.42
	Step 2	.57	4.17	64.62	20.76
	Step 3	.11	.13	2.53	2.16
	Total	.92	4.88	89.89	35.34
	Speedup	3.6	1.1	1.4	5.0

compared to the original PageRank algorithm, both algorithms do drastically improve the time required to compute the PageRank vector.

5. Future work. The relationship between the spectrum of \mathbf{G} and that of \mathbf{S} is known [7, 12]. However, their relationship to the spectrum of \mathbf{H} is not known. Recall that \mathbf{S} and \mathbf{G} are web matrices that have been artificially enhanced to generate advantageous mathematical properties in the matrices. The true web matrix is \mathbf{H} , yet little is known about how far \mathbf{S} and \mathbf{G} and their corresponding ranking vectors stray from the web's true nature. The reordering presented here may illuminate this unknown relationship. The symmetric reordering shows that the eigenvalues of \mathbf{H} are the eigenvalues of \mathbf{H}_{11} plus many 0 eigenvalues for the $\mathbf{0}$ diagonal blocks. The size of \mathbf{H}_{11} may be small enough that its subdominant eigenvalues and eigenvectors can be computed. These subdominant eigenvectors can provide beneficial information about the community structure of a subset of the web [11]. Further, since \mathbf{S} is created from a rank-one update to \mathbf{H} , perturbation theory for eigenvalues and eigenvectors can provide additional information about the PageRank vector with respect to its original hyperlink structure (\mathbf{H} , not \mathbf{S} or \mathbf{G}).

6. Conclusions. Reorderings for linear systems have been well studied. Some popular reorderings, such as the minimum degree reordering, the reverse Cuthill–McKee reordering, and the Dulmage–Mendelson reordering, do not give the nice structure of (5), which requires only one small system solve plus a quick forward substitution that amounts to the work of one power iteration. This reordering naturally produces an efficient alternative to the original PageRank algorithm, which we have called the reordered PageRank algorithm. The reordered algorithm has an asymptotic rate of convergence that is as good as or better than the asymptotic rate of convergence of the original method. The reordered algorithm produces a dataset-dependent speedup over the original algorithm. We presented experimental results

comparing the reordered PageRank algorithm to the original PageRank algorithm on five datasets, showing a factor of nearly 6 speedup on one dataset. This reordering and its linear system formulation also open the door for new PageRank acceleration techniques beyond the usual power method acceleration techniques.

Acknowledgments. The paper [13] by Chris Lee of Stanford University and discussions with Michele Benzi of Emory University sparked the reordering idea presented in this paper. We thank Cleve Moler (Mathworks), Ronny Lempel (IBM Research Labs), and Jon Kleinberg (Cornell) for providing us with several small datasets that we used for testing. Our large dataset from the Stanford web crawl was provided by Chris Lee.

REFERENCES

- [1] A. ARASU, J. NOVAK, A. TOMKINS, AND J. TOMLIN, *PageRank computation and the structure of the Web: Experiments and algorithms*, in The Eleventh International WWW Conference, ACM Press, New York, 2002.
- [2] S. BRIN, R. MOTWANI, L. PAGE, AND T. WINOGRAD, *What can you do with a Web in your pocket?*, Data Engrg. Bull., 21 (1998), pp. 37–47.
- [3] S. BRIN AND L. PAGE, *The anatomy of a large-scale hypertextual Web search engine*, Comput. Networks and ISDN Syst., 33 (1998), pp. 107–117.
- [4] S. BRIN, L. PAGE, R. MOTWANI, AND T. WINOGRAD, *The PageRank Citation Ranking: Bringing Order to the Web*, Technical Report 1999-0120, Computer Science Department, Stanford University, Stanford, CA, 1999.
- [5] N. EIRON, K. S. MCCURLEY, AND J. A. TOMLIN, *Ranking the Web frontier*, in The Thirteenth International World Wide Web Conference, ACM Press, New York, 2004.
- [6] G. H. GOLUB AND C. GREIF, *Arnoldi-Type Algorithms for Computing Stationary Distribution Vectors, with Application to PageRank*, Technical Report SCCM-2004-15, Scientific Computation and Computational Mathematics, Stanford University, Stanford, CA, 2004.
- [7] T. H. HAVELIWALA AND S. D. KAMVAR, *The Second Eigenvalue of the Google Matrix*, Technical Report 2003-20, Stanford University, Stanford, CA, 2003.
- [8] S. D. KAMVAR, T. H. HAVELIWALA, AND G. H. GOLUB, *Adaptive Methods for the Computation of PageRank*, Technical Report 2003-26, Stanford University, Stanford, CA, 2003.
- [9] S. D. KAMVAR, T. H. HAVELIWALA, C. D. MANNING, AND G. H. GOLUB, *Exploiting the Block Structure of the Web for Computing PageRank*, Technical Report 2003-17, Stanford University, Stanford, CA, 2003.
- [10] S. D. KAMVAR, T. H. HAVELIWALA, C. D. MANNING, AND G. H. GOLUB, *Extrapolation methods for accelerating PageRank computations*, in The Twelfth International World Wide Web Conference, ACM Press, New York, 2003.
- [11] J. KLEINBERG, *Authoritative sources in a hyperlinked environment*, J. ACM, 46 (1999), pp. 604–632.
- [12] A. N. LANGVILLE AND C. D. MEYER, *Deeper inside PageRank*, Internet Math. J., 1 (2005), pp. 335–380.
- [13] C. P.-C. LEE, G. H. GOLUB, AND S. A. ZENIOS, *A Fast Two-Stage Algorithm for Computing PageRank and Its Extensions*, Technical Report SCCM-2003-15, Scientific Computation and Computational Mathematics, Stanford University, Stanford, CA, 2003.
- [14] C. D. MEYER, *Matrix Analysis and Applied Linear Algebra*, SIAM, Philadelphia, 2000.