

PAGERANK BEYOND THE WEB

DAVID F. GLEICH

Abstract. Google’s PageRank method was developed to evaluate the importance of web-pages via their link structure. The mathematics of PageRank, however, are entirely general and apply to any graph or network in any domain. Thus, PageRank is now regularly used in bibliometrics, social and information network analysis, and for link prediction and recommendation. It’s even used for systems analysis of road networks, as well as biology, chemistry, neuroscience, and physics. We’ll see the mathematics and ideas that unite these diverse applications.

Key words. PageRank, Markov chain

1. Google’s PageRank. Google created PageRank to address a problem they encountered with their search engine for the world wide web [Brin and Page, 1998; Page et al., 1999]. Given a search query from a user, they could immediately find an immense set of web pages that contained virtually the exact same words as the user entered. Yet, they wanted to incorporate a measure of a page’s importance into these results to distinguish highly recognizable and relevant pages from those that were less well known. To do this, Google designed a system of scores called PageRank that used the link structure of the web to determine which pages are important. While there are many derivations of the PageRank equations [Langville and Meyer, 2006; Pan et al., 2004; Higham, 2005], we will derive it based on a hypothetical random web surfer. Upon visiting a page on the web, our random surfer tosses a coin. If it comes up heads, the surfer randomly clicks a link on the current page and transitions to the new page. If it comes up tails, the surfer *teleports* to a – possibly random – page independent of the current page’s identity. Pages where the random surfer is more likely to appear based on the web’s structure are more important in a PageRank sense.

More generally, we can consider random surfer models on a graph with an arbitrary set of nodes, instead of pages, and transition probabilities, instead of randomly clicked links. The teleporting step is designed to model an external influence on the importance of each node and can be far more nuanced than a simple random choice. Teleporting *is* the essential distinguishing feature of the PageRank random walk that had not appeared in the literature before [Vigna, 2009]. It ensures that the resulting importance scores always exist and are unique. It also makes the PageRank importance scores easy to compute.

These features: simplicity, generality, guaranteed existence, uniqueness, and fast computation are the reasons that PageRank is used in applications far beyond its origins in Google’s web-search. (Although, the success that Google achieved no doubt contributed to additional interest in PageRank!) In biology, for instance, new microarray experiments churn out thousands of genes relevant to a particular experimental condition. Models such as GeneRank [Morrison et al., 2005] deploy the exact same motivation as Google, and almost identical mathematics in order to assist biologists in finding and ordering genes related to a microarray experiment or related to a disease. Throughout our review, we will see applications of PageRank to biology, chemistry, ecology, neuroscience, physics, sports, and computer systems.

Two uses underlie the majority of PageRank applications. In the first, PageRank is used as a network centrality measure [Koschützki et al., 2005]. A network centrality score yields the importance of each node in light of the entire graph structure. And the goal is to use PageRank to help understand the graph better by focusing on what

PageRank reveals as important. It is often compared or contrasted with a host of other centrality or graph theoretic measures. These applications tend to use global, near-uniform teleportation behaviors.

In the second type of use, PageRank is used to illuminate a region of a large graph around a target set of interest; for this reason, we call the second use a localized measure. It is also called personalized PageRank based on PageRank’s origins in the web. Consider a random surfer in a large graph that periodically teleports back to a single start node. If the teleportation is sufficiently frequent, the surfer will never move far from the start node, but the frequency with which the surfer visits nodes before teleporting reveals interesting properties of this localized region of the network. Because of this power, teleportation behaviors are much more varied for these localized applications.

2. The mathematics of PageRank. There are many slight variations on the PageRank problem, yet there is a core definition that applies to the almost all of them. It arises from a generalization of the random surfer idea. Pages where the random surfer is likely to appear have large values in the stationary distribution of a Markov chain that, with probability α , randomly transitions according to the link structure of the web, and with probability $1 - \alpha$ teleports according to a *teleportation distribution vector* \mathbf{v} , where \mathbf{v} is usually a uniform distribution over all pages. In the generalization, we replace the notion of “transitioning according to the link structure of the web” with “transitioning according to a stochastic matrix \mathbf{P} .” This simple change divorces the mathematics of PageRank from the web and forms the basis for the applications we discuss. Thus, it abstracts the random surfer model from the introduction in a relatively seamless way. Furthermore, the vector \mathbf{v} is a critical modeling tool that distinguishes between the two typical uses of PageRank. For centrality uses, \mathbf{v} will resemble a uniform distribution over all possibilities; for localized uses, \mathbf{v} will focus the attention of the random surfer on a region of the graph.

Before stating the definition formally, let us fix some notation. Matrices and vectors are written in bold, Roman letters (\mathbf{A}, \mathbf{x}), scalars are Greek or indexed, unbold Roman ($\alpha, A_{i,j}$). The vector \mathbf{e} is the column vector of all ones, and all vectors are column vectors.

Let $P_{i,j}$ be the probability of transitioning from page j to page i . (Or more generally, from “thing j ” to “thing i ”.) The stationary distribution of the PageRank Markov chain is called the PageRank vector \mathbf{x} . It is the solution of the eigenvalue problem:

$$(\alpha\mathbf{P} + (1 - \alpha)\mathbf{v}\mathbf{e}^T)\mathbf{x} = \mathbf{x}. \quad (2.1)$$

Many take this eigensystem as the definition of PageRank [Langville and Meyer, 2006]. We prefer the following definition instead:

DEFINITION 2.1 (The PageRank Problem). *Let \mathbf{P} be a column-stochastic matrix where all entries are non-negative and the sum of entries in each column is 1. Let \mathbf{v} be a column stochastic vector ($\mathbf{e}^T\mathbf{v} = 1$), and let $0 < \alpha < 1$ be the teleportation parameter. Then the PageRank problem is to find the solution of the linear system*

$$(\mathbf{I} - \alpha\mathbf{P})\mathbf{x} = (1 - \alpha)\mathbf{v}, \quad (2.2)$$

where the solution \mathbf{x} is called the PageRank vector.

The eigenvector and linear system formulations are equivalent if we seek an eigenvector \mathbf{x} of (2.1) with $\mathbf{x} \geq 0$ and $\mathbf{e}^T\mathbf{x} = 1$, in which case:

$$\mathbf{x} = \alpha\mathbf{P}\mathbf{x} + (1 - \alpha)\mathbf{v}\mathbf{e}^T\mathbf{x} = \alpha\mathbf{P}\mathbf{x} + (1 - \alpha)\mathbf{v} \quad \Leftrightarrow \quad (\mathbf{I} - \alpha\mathbf{P})\mathbf{x} = (1 - \alpha)\mathbf{v}.$$

We prefer the linear system because of the following reasons. In the linear system setup, the existence and uniqueness of the solution is immediate: the matrix $\mathbf{I} - \alpha\mathbf{P}$ is a diagonally dominant M-matrix. The solution \mathbf{x} is non-negative for the same reason. Also, there is only one possible normalization of the solution: $\mathbf{x} \geq 0$ and $\mathbf{e}^T \mathbf{x} = 1$. Anecdotally, we note that, among the strategies to solve PageRank problems, those based on the linear system setup are both more straightforward and more effective than those based on the eigensystem approach. And in closing, Page et al. [1999] describe an iteration more akin to a linear system than an eigenvector.

Computing the PageRank vector \mathbf{x} is simple. The humble iteration

$$\mathbf{x}^{(k+1)} = \alpha\mathbf{P}\mathbf{x}^{(k)} + (1 - \alpha)\mathbf{v} \quad \text{where} \quad \mathbf{x}^{(0)} = \mathbf{v} \text{ or } \mathbf{x}^{(0)} = 0$$

is equivalent both to the power method on (2.1) and the Richardson method on (2.2), and more importantly, it has excellent convergence properties when α is not too close to 1. To see this fact, note that the true solution $\mathbf{x} = \alpha\mathbf{P}\mathbf{x} + (1 - \alpha)\mathbf{v}$ and consider the error after a single iteration:

$$\mathbf{x} - \mathbf{x}^{(k+1)} = \underbrace{[\alpha\mathbf{P}\mathbf{x} + (1 - \alpha)\mathbf{v}]}_{\text{the true solution } \mathbf{x}} - \underbrace{[\alpha\mathbf{P}\mathbf{x}^{(k)} + (1 - \alpha)\mathbf{v}]}_{\text{the updated iterate } \mathbf{x}^{(k+1)}} = \alpha\mathbf{P}(\mathbf{x} - \mathbf{x}^{(k)}).$$

Thus, the following theorem characterizes the error after k iterations from two different starting conditions:

THEOREM 2.2. *Let $\alpha, \mathbf{P}, \mathbf{v}$ be the data for a PageRank problem to compute a PageRank vector \mathbf{x} . Then the error after k iterations of the update $\mathbf{x}^{(k+1)} = \alpha\mathbf{P}\mathbf{x}^{(k)} + (1 - \alpha)\mathbf{v}$ is:*

1. *if $\mathbf{x}^{(0)} = \mathbf{v}$, then $\|\mathbf{x} - \mathbf{x}^{(k)}\|_1 \leq \|\mathbf{x} - \mathbf{v}\|_1 \alpha^k \leq 2\alpha^k$; or*
2. *if $\mathbf{x}^{(0)} = 0$, then the error vector $\mathbf{x} - \mathbf{x}^{(k)} \geq 0$ for all k and $\|\mathbf{x} - \mathbf{x}^{(k)}\|_1 = \mathbf{e}^T(\mathbf{x} - \mathbf{x}^{(k)}) = \alpha^k$.*

Common values of α range between 0.1 and 0.99; hence, in the worst case, this method needs at most 3656 iterations to converge to a global 1-norm error of $2^{-52} \approx 10^{-16}$ (because $\alpha^{3656} \leq 2^{-53}$ to account for the possible factor of 2 if starting from $\mathbf{x}^{(0)} = \mathbf{v}$). For the majority of applications we will see, the matrix \mathbf{P} is sparse with fewer than 10,000,000 non-zeros; and thus, these solutions can be computed efficiently on a modern laptop computer.

ASIDE 2.3. *Although this theorem seems to suggest that $\mathbf{x}^{(0)} = 0$ is a superior choice, practical experience suggests that starting with $\mathbf{x}^{(0)} = \mathbf{v}$ results in a faster method. This may be confirmed by using a computable bound on the error based on the residual. Let $\mathbf{r}^{(k)} = (1 - \alpha)\mathbf{v} - (\mathbf{I} - \alpha\mathbf{P})\mathbf{x}^{(k)} = \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}$ be the residual after k iterations. We can use $\|\mathbf{x} - \mathbf{x}^{(k)}\|_1 = \|(\mathbf{I} - \alpha\mathbf{P})^{-1}\mathbf{r}^{(k)}\|_1 \leq \frac{1}{1-\alpha}\|\mathbf{r}^{(k)}\|_1$ in order to check for early convergence.*

This setup for PageRank, where the choice of \mathbf{P} , \mathbf{v} , and α vary by application, applies broadly as the subsequent sections show. However, in many descriptions, authors are not always careful to describe their contributions in terms of a column stochastic matrix \mathbf{P} and distribution vector \mathbf{v} . Rather, they use the following pseudo-PageRank system instead:

DEFINITION 2.4 (The pseudo-PageRank problem). *Let $\bar{\mathbf{P}}$ be a column sub-stochastic matrix where $\bar{P}_{i,j} \geq 0$ and $\mathbf{e}^T \bar{\mathbf{P}} \leq \mathbf{e}^T$ element-wise. Let \mathbf{f} be a non-negative vector, and let $0 < \alpha < 1$ be a teleportation parameter. Then the pseudo-PageRank problem is to find the solution of the linear system*

$$(\mathbf{I} - \alpha\bar{\mathbf{P}})\mathbf{y} = \mathbf{f} \tag{2.3}$$

where the solution \mathbf{y} is called the *pseudo-PageRank vector*.

Again, the pseudo-PageRank vector always exists and is unique because $\mathbf{I} - \alpha\bar{\mathbf{P}}$ is also a diagonally dominant M-matrix. Boldi et al. [2007] was the first to formalize this definition and distinction between PageRank and pseudo-PageRank, although they used the term PseudoRank and the normalization $(\mathbf{I} - \alpha\bar{\mathbf{P}})\mathbf{y} = (1 - \alpha)\mathbf{f}$; some advantages of this alternative form are discussed in Section 5.2. The two problems are equivalent in the following formal sense (which has an intuitive understanding explained in Section 3.1, Strongly Preferential PageRank):

THEOREM 2.5. *Let \mathbf{y} be the solution of a pseudo-PageRank system with α , $\bar{\mathbf{P}}$ and \mathbf{f} . Let $\mathbf{v} = \mathbf{f}/(\mathbf{e}^T\mathbf{f})$. Then if \mathbf{y} is renormalized to sum to 1, that is $\mathbf{x} = \mathbf{y}/(\mathbf{e}^T\mathbf{y})$, then \mathbf{x} is the solution of a PageRank system with α , $\mathbf{P} = \bar{\mathbf{P}} + \mathbf{v}\mathbf{c}^T$, and \mathbf{v} , where $\mathbf{c}^T = \mathbf{e}^T - \mathbf{e}^T\bar{\mathbf{P}} \geq 0$ is a correction vector to make $\bar{\mathbf{P}}$ stochastic.*

Proof. First note that α , \mathbf{P} , and \mathbf{v} is a valid PageRank problem. This is because \mathbf{f} is non-negative and thus \mathbf{v} is column stochastic by definition, and also $\bar{\mathbf{P}}$ is column stochastic because $\mathbf{c} \geq 0$ (hence $\mathbf{P} \geq 0$) and $\mathbf{e}^T\mathbf{P} = \mathbf{e}^T\bar{\mathbf{P}} + \mathbf{c}^T = \mathbf{e}^T$. Next, note that the solution of the PageRank problem for \mathbf{x} satisfies:

$$\mathbf{x} = \alpha\bar{\mathbf{P}}\mathbf{x} + \alpha\mathbf{v}\mathbf{c}^T\mathbf{x} + (1 - \alpha)\mathbf{v} = \alpha\bar{\mathbf{P}}\mathbf{x} + \gamma\mathbf{f} \quad \text{where} \quad \gamma = \frac{\alpha\mathbf{c}^T\mathbf{x} + (1 - \alpha)}{\mathbf{e}^T\mathbf{f}}.$$

Hence $(\mathbf{I} - \alpha\bar{\mathbf{P}})\mathbf{x} = \gamma\mathbf{f}$ and so $\mathbf{x} = \gamma\mathbf{y}$. But, we know that $\mathbf{e}^T\mathbf{x} = 1$ because \mathbf{x} is a solution of a PageRank problem, and the theorem follows. \square

The importance of this theorem is it shows that underlying any pseudo-PageRank system is a true PageRank system in the sense of Definition 2.1. The difference is entirely in terms of the normalization of the solution – which was demonstrated by Del Corso et al. [2004]; Berkhin [2005]; Del Corso et al. [2005]. The result of Theorem 2.2 also applies to solving the pseudo-PageRank system, albeit with the following revisions:

THEOREM 2.6. *Let α , $\bar{\mathbf{P}}$, \mathbf{f} be the data for a pseudo-PageRank problem to compute a pseudo-PageRank vector \mathbf{y} . Then the error after k iterations of the update $\mathbf{y}^{(k+1)} = \alpha\bar{\mathbf{P}}\mathbf{y}^{(k)} + \mathbf{f}$ is:*

1. if $\mathbf{y}^{(0)} = \frac{1}{1-\alpha}\mathbf{f}$, then $\|\mathbf{y} - \mathbf{y}^{(k)}\|_1 \leq \|\mathbf{y} - \mathbf{f}\|_1\alpha^k \leq \frac{2\mathbf{e}^T\mathbf{f}}{1-\alpha}\alpha^k$; or
2. if $\mathbf{y}^{(0)} = 0$, then the error vector $\mathbf{y} - \mathbf{y}^{(k)} \geq 0$ for all k and $\|\mathbf{y} - \mathbf{y}^{(k)}\|_1 = \mathbf{e}^T(\mathbf{y} - \mathbf{y}^{(k)}) \leq \alpha^k$.

ASIDE 2.7. *The error progression proceeds at the same rate for both PageRank and pseudo-PageRank. This can be improved for pseudo-PageRank if the vector $\mathbf{c}^T = \mathbf{e}^T - \mathbf{e}^T\bar{\mathbf{P}} > 0$ (element-wise). In such cases, then we can derive an equivalent system with a smaller value of α and a suitably rescaled matrix $\bar{\mathbf{P}}$.*

These formal results represent the mathematical foundations of all of the PageRank systems that arise in the literature (with a few technical exceptions that we will study in Section 5). The results depend only on the construction of a stochastic matrix or sub-stochastic matrix, a teleportation distribution, and a parameter α . Thus, they apply generally and have no intrinsic relationship back to the original motivation of PageRank for the web. Each type of PageRank problem has a unique solution that always exists, and the two convergence theorems justify that simple algorithms for PageRank converge to the unique solutions quickly. These are two of the most attractive features of PageRank.

One final set of mathematical results is important to understand the behavior of localized PageRank; however, the precise statement of these results requires a lengthy and complicated diversion into graph partitioning, graph cuts, and spectral graph

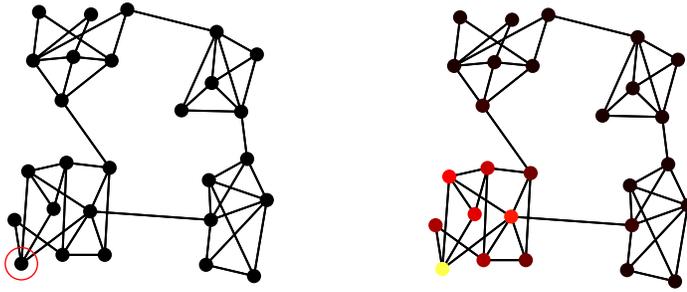


FIG. 2.1. An illustration of the empirical properties of localized PageRank vectors with teleportation to a single node in an isolated region. In the graph at left, the teleportation vector is the single circled node. The PageRank vector is shown as the node color in the right figure. PageRank values remain high within this region and are nearly zero in the rest of the graph. Theory from Andersen et al. [2006] explains when this property occurs.

theory. Instead, we'll state this a bit informally. Suppose that we solve a localized PageRank problem in a large graph, but the nodes we select for teleportation lie in a region that is somehow isolated, yet connected to the rest of the graph. Then the final PageRank vector is large only in this isolated region and has small values on the remainder of the graph. This behavior is exactly what most uses of localized PageRank want: they want to find out what is nearby the selected nodes and far from the rest of the graph. Proving this result involves spectral graph theory, Cheeger inequalities, and localized random walks – see Andersen et al. [2006] for more detail. Instead, we illustrate this theory with Figure 2.1.

Next, we will see some of the common constructions of the matrices \mathbf{P} and $\bar{\mathbf{P}}$ that arise when computing PageRank on a graph. These justify that PageRank is also a simple construction.

3. PageRank constructions. When a PageRank method is used within an application, there are two common motivations. In the centrality case, the input is a graph representing relationships or flows between a set of things – they may be documents, people, genes, proteins, roads, or pieces of software – and the goal is to determine the expected importance of each piece in light of the full set of relationships and the teleporting behavior. This motivation was Google's original goal in crafting PageRank. In the localized case, the input is also the same type of graph, but the goal is to determine the importance relative to a small subset of the objects. In either case, we need to build a stochastic or sub-stochastic matrix from a graph. In this section, we review some of the common constructions that produce a PageRank or pseudo-PageRank system. For a visual overview of some of the possibilities, see Figures 3.1 and 3.2.

Notation for graphs and matrices. Let \mathbf{A} be the adjacency matrix for a graph where we assume that the vertex set is $V = \{1, \dots, n\}$. The graph could be directed, in which case \mathbf{A} is non-symmetric, or undirected, in which case \mathbf{A} is symmetric. The graph could also be weighted, in which case $A_{i,j}$ gives the positive weight of edge (i, j) . Edges with zero weight are assumed to be irrelevant and equivalent to edges that are not present. For such a graph, let \mathbf{d} be the vector of node out-degrees, or equivalently, the vector of row-sums: $\mathbf{d} = \mathbf{A}\mathbf{e}$. The matrix \mathbf{D} is simply the diagonal matrix with \mathbf{d} on the diagonal. Weighted graphs are extremely common in applications when the

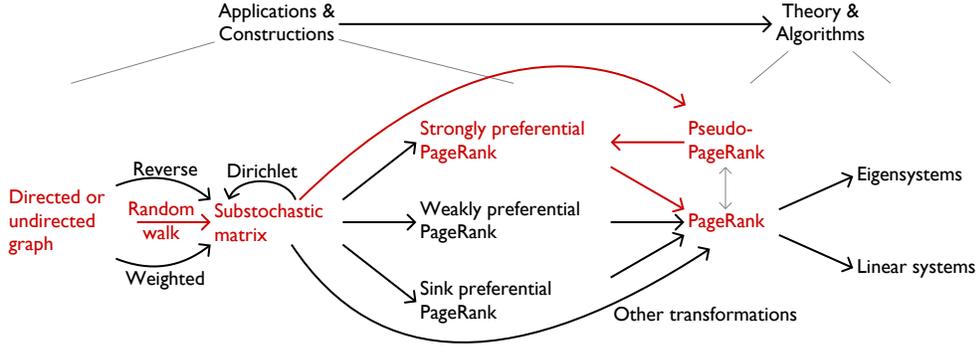


FIG. 3.1. An overview of PageRank constructions and how they relate. The vast majority of PageRank applications fall somewhere on the red path.

weights reflect a measure of the *strength* of the relationships between two nodes.

3.1. The standard random walk. In the standard construction of PageRank, the matrix \mathbf{P} represents a uniform random walk operation on the graph \mathbf{A} . When the graph is weighted, the simple generalization is to model a non-uniform walk that chooses subsequent nodes with probability proportional to the connecting edge’s weight. The elements of $\bar{\mathbf{P}}$ are rather similar between the two cases:

$$\bar{P}_{j,i} = \frac{A_{i,j}}{\sum_k A_{i,k}} = \frac{A_{i,j}}{d_i} = \begin{array}{l} \text{probability of taking the transition} \\ \text{from } i \text{ to } j \text{ via a random walk step.} \end{array}$$

Notice two features of this construction. First, we transpose between j, i and i, j . This is because $A_{i,j}$ indicates an edge from node i to node j , whereas the probability transition matrix element i, j indicates that node i can be reached via node j . Second, we have written $\bar{\mathbf{P}}$ and $\bar{P}_{j,i}$ here because there may be nodes of the graph with *no outlinks*. These nodes are called *dangling nodes*. Dangling nodes complicate the construction of stochastic matrices \mathbf{P} in a few ways because we must specify a behavior for the random walk at these nodes in order to fully specify the stochastic matrix.

As a matrix formula, the standard random walk construction is:

$$\bar{\mathbf{P}} = \mathbf{A}^T \mathbf{D}^+.$$

Here, we have used the *pseudo-inverse* of the degree matrix to “invert” the diagonal matrix in light of the dangling nodes with 0 out-degrees. Let \mathbf{c}^T be the sub-stochastic correction vector. For the standard random walk construction, \mathbf{c}^T is just an indicator vector for the dangling nodes:

$$c_i = 1 - \sum_k \bar{P}_{k,i} = \begin{cases} 1 & \text{node } i \text{ is dangling} \\ 0 & \text{otherwise.} \end{cases}$$

We shall now see a few ideas that turn these sub-stochastic matrices into fully stochastic PageRank problems.

Strongly Preferential PageRank. Given a directed graph with dangling nodes, the standard random walk construction produces the sub-stochastic matrix $\bar{\mathbf{P}}$ described above. If we had just used this matrix to solve a pseudo-PageRank problem with

a stochastic teleportation vector $\mathbf{f} = (1 - \alpha)\mathbf{v}$, then, by Theorem 2.5, the result is equivalent up to normalization to computing PageRank on the matrix:

$$\mathbf{P} = \bar{\mathbf{P}} + \mathbf{c}\mathbf{v}^T.$$

This construction models a random walk that transitions according to the distribution \mathbf{v} when visiting a dangling node. This behavior reinforces the effect of the teleportation vector \mathbf{v} , or preference vector as it is sometimes called. Because of this reinforcement, Boldi et al. [2007] called the construction $\mathbf{P} = \bar{\mathbf{P}} + \mathbf{c}\mathbf{v}^T$ a *strongly preferential PageRank* problem. Again, many authors are not careful to *explicitly* choose a correction to turn the sub-stochastic matrix into a stochastic matrix. Their lack of choice, then, *implicitly* chooses the strongly preferential PageRank system.

Weakly Preferential PageRank & Sink Preferential PageRank. Boldi et al. [2007] also proposed the *weakly preferential PageRank* system. In this case, the behavior of the random walk at dangling nodes is adjusted *independently* of the choice of teleportation vector. For instance, Langville and Meyer [2004] advocates transitioning uniformly from dangling nodes. In such a case, let $\mathbf{u} = \mathbf{e}/n$ be the uniform distribution vector, then a weakly preferential PageRank system is:

$$\mathbf{P} = \bar{\mathbf{P}} + \mathbf{c}\mathbf{u}^T.$$

We note that another choice of behavior is for the random walk to remain at dangling nodes until it moves away via a teleportation step:

$$\mathbf{P} = \bar{\mathbf{P}} + \text{diag}(\mathbf{c}).$$

We call this final method *sink preferential PageRank*. These systems are less common. These choices should be used when the matrix \mathbf{P} models some type of information or material flow that must be decoupled from the teleporting behavior.

3.2. Reverse PageRank. In reverse PageRank, we compute PageRank on the transposed graph \mathbf{A}^T . This corresponds to reversing the direction of each edge (i, j) to be an edge (j, i) . Reverse PageRank is often used to determine *why* a particular node is important rather than *which* nodes are important [Fogaras, 2003; Gyöngyi et al., 2004; Bar-Yossef and Mashiach, 2008]. Intuitively speaking, in reverse PageRank, we model a random surfer that follows in-links instead of out-links. Thus, large reverse PageRank values suggest nodes that can reach many nodes in the graph. When these are localized, they then provide evidence for why a node has large PageRank.

3.3. Dirichlet PageRank. Consider a PageRank problem where we wish to fix the importance score of a subset of nodes [Chung et al., 2011]. Let S be a subset of nodes such that $i \in S$ implies that $v_i = 0$. A Dirichlet PageRank problem seeks a solution of PageRank where each node i in S is fixed to a *boundary* value b_i . Formally, the goal is to find \mathbf{x} :

$$(\mathbf{I} - \alpha\mathbf{P})\mathbf{x} = (1 - \alpha)\mathbf{v} \quad \text{where} \quad x_i = b_i \text{ for } i \in S.$$

These problems reduce to solving a pseudo-PageRank system. Consider a block partitioning of \mathbf{P} based on the set S and the complement set of vertices \bar{S} :

$$\mathbf{P} = \begin{bmatrix} \mathbf{P}_{S,S} & \mathbf{P}_{S,\bar{S}} \\ \mathbf{P}_{\bar{S},S} & \mathbf{P}_{\bar{S},\bar{S}} \end{bmatrix}.$$

Then the Dirichlet PageRank problem is

$$\begin{bmatrix} \mathbf{I} & 0 \\ -\alpha\mathbf{P}_{\bar{S},S} & \mathbf{I} - \alpha\mathbf{P}_{\bar{S},\bar{S}} \end{bmatrix} \begin{bmatrix} \mathbf{b} \\ \mathbf{x}_{\bar{S}} \end{bmatrix} = (1 - \alpha) \begin{bmatrix} 0 \\ \mathbf{v}_{\bar{S}} \end{bmatrix}.$$

This system is equivalent to a pseudo-PageRank problem with $\bar{\mathbf{P}} = \mathbf{P}_{\bar{S},\bar{S}}$ and $\mathbf{f} = (1 - \alpha)\mathbf{v}_{\bar{S}} + \alpha\mathbf{P}_{\bar{S},S}\mathbf{b}$.

3.4. Weighted PageRank. In the standard random walk construction for PageRank on an unweighted graph, the probability of transitioning from node i to any of its neighbors j is the same: $1/d_i$. Weighted PageRank [Xing and Ghorbani, 2004; Jiang, 2009] alters this assumption such that the walk preferentially visits high-degree nodes. Thus, the probability of transitioning from node i to node j depends on the degree of j relative to the total sum of degrees of all i 's neighbors. In our notation, if the input is adjacency matrix \mathbf{A} with degree matrix \mathbf{D} , then the sub-stochastic matrix $\bar{\mathbf{P}}$ is given by the non-uniform random walk construction on the weighted graph with adjacency matrix $\mathbf{W} = \mathbf{AD}$, that is, $\bar{\mathbf{P}} = \mathbf{DA}^T \text{diag}(\mathbf{ADe})^+$. More generally, let \mathbf{D}_W be a non-negative weighting matrix. It could be derived from the graph itself based on the out-degree, in-degree, or total-degree (the sum of in- and out-degree), or from some external source. Then $\bar{\mathbf{P}} = \mathbf{D}_W\mathbf{A}^T \text{diag}(\mathbf{AD}_W\mathbf{e})^{-1}$. *Let us note that weighted PageRank uses a specific choice of weights for the prior importance of each node; the setting here already adapts seamlessly to edge-weighted graphs.*

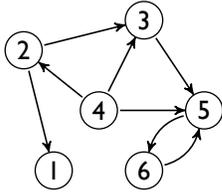
3.5. PageRank on an undirected graph. One final construction is to use PageRank on an undirected graph. Those familiar with Markov chain theory often find this idea puzzling at first. A uniform random walk on a connected, undirected graph has a well-known, unique stationary distribution [Stewart, 1994, is a good numerical treatment of such issues]:

$$\underbrace{\mathbf{A}^T\mathbf{D}^{-1}}_{\mathbf{P}}\mathbf{x} = \mathbf{x} \text{ is solved by } \mathbf{x} = \mathbf{De}/(\mathbf{e}^T\mathbf{d}).$$

This works because both the row and column sums of \mathbf{A} and \mathbf{A}^T are identical, and the resulting construction is a *reversible Markov chain* [Aldous and Fill, 2002, is a good reference on this topic]. If $\alpha < 1$, then the PageRank Markov chain *is not* a reversible Markov chain even on an undirected graph, and hence, has no simple stationary distribution. PageRank vectors of undirected graphs, when combined with carefully constructed teleportation vectors \mathbf{v} , yield important information about the presence of small isolated regions in the graph [Andersen et al., 2006; Gleich and Mahoney, 2014]; formally these results involve graph cuts and small conductance sets. These vectors are most useful when the teleportation vector is far away from the uniform distribution, such as the case in Figure 2.1 where the graph is undirected.

ASIDE 3.1. *Of course, if the teleportation distribution $\mathbf{v} = \mathbf{De}/(\mathbf{e}^T\mathbf{d})$, then the resulting chain is reversible. The PageRank vector is then equal to \mathbf{v} itself. There are also specialized PageRank-style constructions that preserve reversibility with more interesting stationary distributions [Avrachenkov et al., 2010].*

4. PageRank applications. When PageRank is used within applications, it tends to acquire a new name. We will see:



A directed graph

$$\mathbf{A} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix} \quad \mathbf{d} = \begin{bmatrix} 0 \\ 2 \\ 1 \\ 3 \\ 1 \\ 1 \end{bmatrix} \quad \mathbf{c} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

The adjacency matrix, degree vector, and correction vector

Random walk

$$\bar{\mathbf{P}} = \begin{bmatrix} 0 & 1/2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1/3 & 0 & 0 \\ 0 & 1/2 & 0 & 1/3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1/3 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

$$\bar{\mathbf{P}} = \mathbf{A}^T \mathbf{D}^+$$

Strongly preferential

$$\mathbf{P} = \begin{bmatrix} 0 & 1/2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1/3 & 0 & 0 \\ 1/3 & 1/2 & 0 & 1/3 & 0 & 0 \\ 1/3 & 0 & 0 & 0 & 0 & 0 \\ 1/3 & 0 & 1 & 1/3 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

$$\mathbf{P} = \bar{\mathbf{P}} + \mathbf{v}\mathbf{c}^T$$

Weakly preferential

$$\mathbf{P} = \begin{bmatrix} 1/6 & 1/2 & 0 & 0 & 0 & 0 \\ 1/6 & 0 & 0 & 1/3 & 0 & 0 \\ 1/6 & 1/2 & 0 & 1/3 & 0 & 0 \\ 1/6 & 0 & 0 & 0 & 0 & 0 \\ 1/6 & 0 & 1 & 1/3 & 0 & 1 \\ 1/6 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

$$\mathbf{P} = \bar{\mathbf{P}} + \mathbf{u}\mathbf{c}^T$$

$\mathbf{u} \neq \mathbf{v}$

Reverse

$$\bar{\mathbf{P}} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1/2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1/3 & 0 \\ 0 & 1 & 1/2 & 0 & 1/3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1/3 & 0 \end{bmatrix}$$

$$\bar{\mathbf{P}} = \mathbf{A} \text{diag}(\mathbf{A}^T \mathbf{e})^+$$

Dirichlet

$$\bar{\mathbf{P}} = \begin{bmatrix} 0 & 0 & 1/3 & 0 & 0 \\ 1/2 & 0 & 1/3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1/3 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

$$S = \{2, 3, 4, 5, 6\}$$

$$\bar{\mathbf{P}} = \bar{\mathbf{P}}_{\bar{S}, \bar{S}}$$

$S \subset V$

Weighted

$$\bar{\mathbf{P}} = \begin{bmatrix} 0 & 1/4 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 3/10 & 0 & 0 \\ 0 & 3/4 & 0 & 3/10 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 4/10 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

$$\bar{\mathbf{P}} = (\mathbf{D}_W \mathbf{A}^T) \text{diag}(\mathbf{A} \mathbf{D}_W \mathbf{e})^+$$

\mathbf{D}_W is a diagonal weighting matrix, e.g. total degree here

FIG. 3.2. A directed graph and some of the different PageRank constructions on that graph. For the stochastic constructions, we have $\mathbf{v}^T = [0 \ 0 \ \frac{1}{3} \ \frac{1}{3} \ \frac{1}{3} \ 0]$ and $\mathbf{u} = \mathbf{e}/n$. Note that node 4 is dangling in the reverse PageRank construction. For the weighted construction, the total degrees are [1 3 3 3 4 2].

GeneRank	TimedPageRank	ObjectRank	HostRank
ProteinRank	CiteRank	FolkRank	DirRank
IsoRank	AuthorRank	ItemRank	TrustRank
MonitorRank	PopRank	BuddyRank	BadRank
BookRank	FactRank	TwitterRank	VisualRank

The remainder of this section explores the uses of PageRank within different domains. It is devoted to the most interesting and diverse uses and should not, necessarily, be read linearly. Our intention is not to cover the full details, but to survey the diversity of applications of PageRank. We recommend returning to the primary sources for additional detail.

Chemistry · §4.1	Literature · §4.7
Biology · §4.2	Bibliometrics · §4.8
Neuroscience · §4.3	Databases & Knowledge systems · §4.9
Engineered systems · §4.4	Recommender systems · §4.10
Mathematical systems · §4.5	Social networks · §4.11
Sports · §4.6	The web, redux · §4.12

4.1. PageRank in chemistry. The term “graph” arose from “chemico-graph” or a picture of a chemical structure [Sylvester, 1878]. Much of this chemical terminology remains with us today. For instance, the valence of a molecule is the number of potential bonds it can make. The valence of a vertex is synonymous with its degree, or the number of connections it makes in the graph. It is fitting, then, that recent work by Mooney et al. [2012] uses PageRank to study molecules in chemistry. In particular, they use PageRank to assess the change in a network of molecules linked by hydrogen bonds among water molecules. Given the output of a molecular dynamics simulation that provides geometric locations for a solute in water, the graph contains edges between the water molecules if they have a potential hydrogen bond to a solute molecule. The goal is to assess the hydrogen bond potential of a solvent. The PageRank centrality scores using uniform teleportation with $\alpha = 0.85$ are strongly correlated with the degree of the node – which is expected – but the deviance of the PageRank score from the degree identifies important outlier molecules with smaller degree than many in their local regions. The authors compare the networks based the PageRank values with and without a solute to find structural differences.

4.2. PageRank in biology & bioinformatics: GeneRank, ProteinRank, IsoRank. Biology and bioinformatics are currently awash in network data. Some of the most interesting applications of PageRank arise when it is used to study these networks. Most of these applications use PageRank to reveal localized information about the graph based on some form of external data.

GeneRank. Microarray experiments are a measurement of whether or not a gene’s expression is promoted or repressed in an experimental condition. Microarrays estimate the outcomes for thousands of genes simultaneously in a few experimental conditions. The results are extremely noisy. GeneRank [Morrison et al., 2005] is a PageRank-inspired idea to help to denoise them. The essence of the idea is to use a graph of known relationships between genes to find genes that are highly related to those promoted or repressed in the experiment, but were not themselves promoted or repressed. Thus, they use the microarray expression results as the teleportation distribution vector for a PageRank problem on a network of known relationships between genes. The network of relationships between genes is undirected, unweighted with a few thousand nodes. This problem uses a localized teleportation behaviour and, experimentally, the best choice of α ranges between 0.75 and 0.85. Teleporting is used to focus the search.

Finding correlated genes. This same idea of using a network of known relationships in concert with an experiment encapsulates many of the other uses of PageRank in biology. Jiang et al. [2009] use a combination of PageRank and BlockRank [Kamvar et al., 2003; Kamvar, 2010] on tissue-specific protein-protein interaction networks in order to find genes related to type 2 diabetes. The teleportation is provided by 34 proteins known to be related to that disease with $\alpha = 0.92$.

Winter et al. [2012] use PageRank to study pancreatic ductal adenocarcinoma, a type of cancer responsible for 130,000 deaths each year, with a particularly poor prognosis (2% mortality after five years). They identified seven genes that better predicted patient survival than all existing tools, and validated this in a clinical trial. One curious feature is that their teleportation parameter was small, $\alpha = 0.3$. This was chosen based on a cross-validation strategy in a statistically rigorous way. The particular type of teleportation they used was based on the correlation between the expression level of a gene and the survival time of the patient.

ProteinRank. The goal of ProteinRank [Freschi, 2007] is similar, in spirit, to GeneRank. Given an undirected network of protein-protein interactions and human-

curated functional annotations about what these proteins do, the goal is to find proteins that may share a functional annotation. Thus, the PageRank problem is, again, a localized use. The teleportation distribution is given by a random choice of nodes with a specific functional annotation. The PageRank vector reveals proteins that are highly related to those with this function, but do not themselves have that function labeled.

Protein distance. Recall that the solution of a PageRank problem for a given teleportation vector \mathbf{v} involves solving $(\mathbf{I} - \alpha\mathbf{P})\mathbf{x} = (1 - \alpha)\mathbf{v}$. The resolvent matrix $\mathbf{X} = (1 - \alpha)(\mathbf{I} - \alpha\mathbf{P})^{-1}$ corresponds to computing PageRank vectors that teleport to every individual node. The entry $X_{i,j}$ is the value of the i th node when the PageRank problem is localized on node j . One interpretation for this score is the PageRank that node j contributes to node i , which has the flavor of a similarity score between node i and j . Voevodski et al. [2009] base an affinity measure between proteins on this idea. Formally, consider an undirected, unweighted protein-protein interaction network. Compute the matrix \mathbf{X} for $\alpha = 0.85$, and the affinity matrix $\mathbf{S} = \min(\mathbf{X}, \mathbf{X}^T)$. (For an undirected graph, a quick calculation shows that $\mathbf{X}^T = \mathbf{D}^{-1}\mathbf{X}\mathbf{D}$.) For each vertex i in the graph, form links to the k vertices with the largest values in row of i of \mathbf{S} . These PageRank affinity scores show a much larger correlation with known protein relationships than do other affinity or similarity metrics between vertices.

IsoRank. Consider the problem of deciding if the vertices of two networks can be mapped to each other. The relationship between this problem and PageRank is surprising and unexpected; although precursor literature existed [Jeh and Widom, 2002; Blondel et al., 2004]. Singh et al. [2007] proposes a PageRank problem to estimate how much of a match the two nodes are in a diffusion sense. They call it IsoRank based on the idea of ranking graph isomorphisms. Let \mathbf{P} be the Markov chain for one network and let \mathbf{Q} be the Markov chain for the second network. Then IsoRank solves a PageRank problem on $\mathbf{Q} \otimes \mathbf{P}$. The solution vector \mathbf{x} is a vectorized form of a matrix \mathbf{X} where X_{ij} indicates a likelihood that vertex i in the network underlying \mathbf{P} will match to vertex j in the network underlying \mathbf{Q} . See Figure 4.1 for an example. If we have an apriori measure of similarity between the vertices of the two networks, we can add this as a teleportation distribution term. IsoRank problems are some of the largest PageRank problems around due to the Kronecker product (e.g. Gleich et al. [2010] has a problem with 4 billion nodes and 100 billion edges). But there are quite a few good algorithmic approaches to tackle them by using properties of the Kronecker product [Bayati et al., 2013] and low-rank matrices [Kollias et al., 2011].

The IsoRank authors consider the problem of matching protein-protein interaction networks between distinct species. The goal is to leverage insight about the proteins from a species such as a mouse in concert with a matching between mouse proteins and human proteins, based on their interactions, in order to hypothesize about possible functions for proteins in a human. For these problems, each protein is coded by a gene sequence. The authors construct a teleportation distribution by comparing the gene sequences of each protein using a tool called BLAST. They found that using α around 0.9 gave the highest structural similarity between the two networks.

4.3. PageRank in neuroscience. The human brain connectome is one of the most important networks, about which we understand surprisingly little. Applied network theory is one of a variety of tools currently used to study it [Sporns, 2011]. Thus, it is likely not surprising that PageRank has been used to study the properties of networks related to the connectome [Zuo et al., 2011]. Most recently, PageRank helped evaluate the importance of brain regions given observed correlations of brain activity. In the resulting graph, two voxels of an MRI scan are connected if the correlation

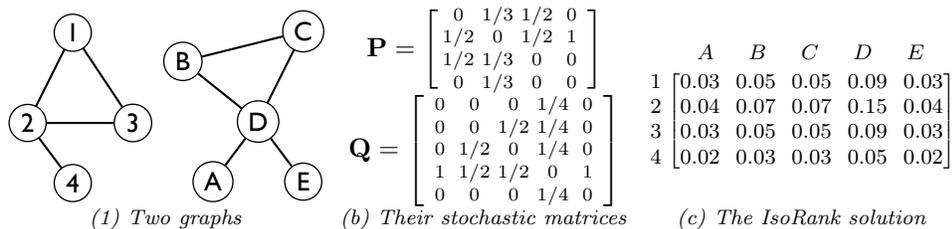


FIG. 4.1. An illustration of the IsoRank problem. The solution, written here as a matrix, gives the similarity between pairs of nodes of the graph. For instance, node 2 is most similar to node D. Removing this match, then nodes 1 and 3 are indistinguishable from B and C. Removing these leaves node 4 equally similar to A and E. In this example we solved $(\mathbf{I} - \alpha\mathbf{Q} \otimes \mathbf{P})\mathbf{x} = (1 - \alpha)\mathbf{e}/20$ with $\alpha = 0.85$.

between their functional MRI time-series is high. Edges with weak correlation are deleted and the remainder are retained with either binary weights or the correlation weights. The resulting graph is also undirected, and they use PageRank, combined with community detection and known brain regions, in order to understand changes in brain structure across a population of 1000 individuals that correlate with age.

Connectome networks are widely hypothesized to be hierarchically organized. Given a directed network that should express a hierarchical structure, how can we recover the order of the nodes that minimizes the discrepancy with a hierarchical hypothesis? Crofts and Higham [2011] consider PageRank for this application on networks of neural connections from *C. Elegans*. They find that this gives poor results compared with other network metrics such as the Katz score [Katz, 1953], and communicability [Estrada et al., 2008]. In their discussion, the authors note that this result may have been a mismatch of models, and conjecture that the flow of influence in PageRank was incorrect. Literature involving Reverse PageRank (Section 3.2) strengthens this conjecture. Let us reiterate that although PageRank models are easy to apply, they must be employed with some care in order to get the best results.

4.4. PageRank in complex engineered systems: MonitorRank. The applications of PageRank to networks in chemistry, biology, and neuroscience are part of the process of investigating and analyzing something we do not fully understand. PageRank methods are also used to study systems that we explicitly engineered. As these engineered systems grow, they become increasingly complex, with networks and submodules interacting in unpredictable, nonlinear ways. Network analysis methods like PageRank, then, help reveal these details. We’ll see two examples: software systems and city systems.

MonitorRank. Diagnosing root causes of issues in a modern distributed system is painstaking work. It involves repeatedly searching through error logs and tracing debugging information. MonitorRank [Kim et al., 2013] is a system to provide guidance to a systems administrator or developer as they perform these activities. It returns a ranked list of systems based on the likelihood that they contributed to, or participated in, an anomalous situation. Consider the systems underlying the LinkedIn website: each service provides one or more APIs that allow other services to utilize its resources. For instance, the web-page generator uses the database and photo store. The photo store in turn uses the database, and so on. Each combination of a service and a programming interface becomes a node in the MonitorRank graph. Edges are directed and indicate the direction of function calls – e.g. web-page to photo store. Given that an anomaly was detected in a system, MonitorRank solves a personalized PageRank

problem on a weighted, augmented version of the call graph, where the weights and augmentation depend on the anomaly detected. (The construction is interesting, albeit tangential, and we refer readers to that paper for the details.) The localized PageRank scores help determine the anomaly. The graphs involved are fairly small: a few hundred to a few thousand nodes.

PageRank of the Linux kernel. The Linux kernel is the foundation for an open source operating system. It has evolved over the past 20 years with contributions from nearly 2000 individuals in an effort with an estimated value of \$3 billion. As of July 2013, the Linux kernel comprised 15.8 million lines of code containing around 300,000 functions. The kernel call graph is a network that represents dependencies between functions and both PageRank and reverse PageRank, as centrality scores, produce an ordering of the most important functions in Linux [Chepelianskii, 2010]. The graphs were directed with a few million edges. Teleportation was typical: $\alpha = 0.85$ with a global, uniform $\mathbf{v} = \mathbf{e}/n$. They find that utility functions such as `printk`, which prints messages from the kernel, and `memset`, a routine that initializes a region of memory, have the highest PageRank, whereas routines that initialize the system such as `start_kernel` have the highest reverse PageRank. Chepelianskii [2010] further uses the distribution of PageRank and reverse PageRank scores to characterize the properties of a software system. (This same idea is later used for Wikipedia too, Zhurov et al. 2010, Section 4.12.)

Roads and Urban Spaces. Another surprising use of PageRank is with road and urban space networks. PageRank helps to predict both traffic flow and human movement in these systems. The natural road construction employed is an interesting graph. A natural road is more or less what it means: it's a continuous path, built from road segments by joining adjacent segments together if the angle is sufficiently small and there isn't a better alternative. (For help visualizing this idea, consider traffic directions that state: "Continue straight from High street onto Main street." This would mean that there is one natural road joining High street and Main street.) Using PageRank with $\alpha = 0.95$, Jiang et al. [2008] finds that PageRank is the best network measure in terms of predicting traffic on the individual roads. These graphs have around 15,000 nodes and around 50,000 edges. Another group used PageRank to study Markov chain models based on the line-graph of roads [Schlote et al., 2012]. That is, given a graph of intersections (nodes) and roads (edges), the line graph, or dual graph, changes the role of roads to the nodes and intersections to the edges. In this context, PageRank's teleportation mirrors the behavior of starting or ending a journey on each street. This produces a different value of α for each node that reflects the tendency of individuals to park, or end their journey, on each street. Note that this is slightly different setup where each node has a separate teleportation parameter α , rather than a different entry in the teleportation vector. Assuming that each street has some probability of a journey ending there, then this system is equivalent to a more general PageRank construction (Section 5.5). These Markov chains are used to study road planning and optimal routing in light of new constraints imposed by electric vehicles.

An urban space is the largest space of a city observable from a single vantage point. For instance, the Mission district of San Francisco is too large, but the area surrounding Dolores Park is sufficiently small to be appreciated as a whole. For the study by Jiang [2009], an urban space is best considered as a city neighborhood or block. The urban space network connects adjacent spaces, or blocks, if they are physically adjacent. The networks of urban spaces in London, for instance, have up to 20,000

nodes and 100,000 links. In these networks, weighted PageRank (Section 3.4) best predicts human mobility in a case study of movement within London. It outperforms PageRank, and in fact, they find that weighted PageRank with $\alpha = 1$ accounts for up to 60% of the observed movement. Both using weighted PageRank and $\alpha = 1$ make sense for these problems – individuals and businesses are likely to co-locate places with high connectivity, and individuals cannot teleport over the short time-frames used for the human mobility measurements. Based on the evidence here, we would hypothesize that using $\alpha < 1$ would better generalize over longer time-spans.

4.5. PageRank in mathematical systems. Graphs and networks arise in mathematics to abstract the properties of systems of equations and processes to relationships between simple sets. We present one example of what PageRank reveals about a dynamical system by abstracting the phase-space to a discrete set of points and modeling transitions among them. Curiously, PageRank and its localization properties has not yet been used to study properties of Cayley graphs from large, finite groups, although closely related structures have been examined [Frahm et al., 2012].

PageRank of symbolic images and Ulam networks. Let f be a discrete-time dynamical system on a compact state space M . For instance, M will be the subset of \mathbb{R}^2 formed by $[0, 2\pi] \times [0, 2\pi]$ for our example below. Consider a covering of M by cells C . In our forthcoming example, this covering will just be a set of non-overlapping cells that form a regular, discrete partition into cells of size $2\pi/N \times 2\pi/N$. The symbolic image [Osipenko, 2007] of f with respect to C is a graph where the vertices are the cells and $C_i \in C$ links to $C_j \in C$ if $x \in C_i$ and $f(x) \in C_j$. The Ulam network is a weighted approximation to this graph that is constructed by simulating s starting points within cell C_i and forming weighted links to their destinations C_j [Shepelyansky and Zhironov, 2010]. The example studied by those authors, and the example we will consider here, is the Chirikov typical map.

$$\begin{aligned} y_{t+1} &= \eta y_t + k \sin(x_t + \theta_t) \\ x_{t+1} &= x_t + y_{t+1}. \end{aligned}$$

It models a kicked oscillator. We generate T random phases θ_t and look at the map:

$$f(x, y) = (x_{T+1}, y_{T+1}) \bmod 2\pi \quad \text{where} \quad x_1 = x, y_1 = y.$$

That is, we iterate the map for T steps for each of the T random phase shifts $\theta_1, \dots, \theta_T$. Applying the construction above with $s = 1000$ random samples from each cell yields a directed weighted graph G with N^2 nodes and at most $N^2 s$ edges. PageRank on this graph, with uniform teleportation, yields beautiful pictures of the transient behaviors of this chaotic dynamical system; these are easy to highlight with modest teleportation parameters such as $\alpha = 0.85$ because this regime inhibits the dynamical system from converging to its stable attractors. This application is particularly useful for modeling the effects of different PageRank constructions as we illustrate in Figure 4.2. For that figure, the graph has 262,144 nodes and 4,106,079 edges, $\eta = 0.99$, $k = 0.22$, $T = 10$.

4.6. PageRank in sports. Stochastic matrices and eigenvector ranking methods are nothing new in the realm of sports ranking [Keener, 1993; Callaghan et al., 2007; Langville and Meyer, 2012]. One of the natural network constructions for sports is the winner network. Each team is a node in the network, and node i points to node j if j won in the match between i and j . These networks are often weighted by the score by which team j beat team i . Govan et al. [2008] used the centrality sense of PageRank with uniform teleportation and $\alpha = 0.85$ to rank football teams with these

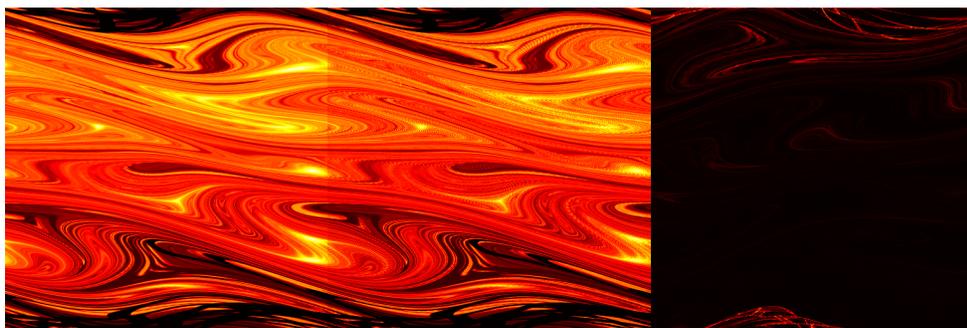


FIG. 4.2. PageRank vectors of the symbolic image, or Ulam network, of the Chirikov typical map with $\alpha = 0.9$ and uniform teleportation. From left to right, we show the standard PageRank vector, the weighted PageRank vector using the unweighted cell in-degree count as the weighting term, and the reverse PageRank vector. Each node in the graph is a point (x, y) , and it links to all other points (x, y) reachable via the map f (see the text). The graph is weighted by the likelihood of the transition. PageRank, itself, highlights both the attractors (the bright regions), and the contours of the transient manifold that leads to the attractor. The weighted vector looks almost identical, but it exhibits an interesting stippling effect. The reverse PageRank highlights regions of the phase-space that are exited quickly, and thus, these regions are dark or black in the PageRank vector. The solution vectors were scaled by the cube-root for visualization purposes. These figures are incredibly beautiful and show important transient regions of these dynamical systems.

winner networks. The intuitive idea underlying these rankings is that of a random fan that follows a team until another team beats them, at which point they pick up the new team, and periodically restarts with an arbitrary team. In the Govan et al. [2008] construction, they corrected dangling nodes using a strongly preferential modification, although, we note that a sink preferential modification may have been more appropriate given the intuitive idea of a random fan. Radicchi [2011] used PageRank on a network of tennis players with the same construction. Again, this was a weighted network. PageRank with $\alpha = 0.85$ and uniform teleportation on the tennis network placed Jimmy Conors in the best player position.

4.7. PageRank in literature: BookRank. PageRank methods help with three problems in literature. What are the most important books? Which story paths in hypertextual literature are most likely? And what should I read next?

For the first question, Jockers [2012] defines a complicated distance metric between books using topic modeling ideas from latent Dirichlet allocation [Blei et al., 2003]. Using PageRank as a centrality measure on this graph, in concert with other graph analytic tools, allows Jockers to argue that Jane Austin and Walter Scott are the most original authors of the 19th century.

Hypertextual literature contains multiple possible story paths for a single novel. Among American children of similar age to me, the most familiar would be the *Choose your own adventure* series. Each of these books consists of a set of storylets; at the conclusion of a storylet, the story either ends, or presents a set of possibilities for the next story. Kontopoulou et al. [2012] argue that the random surfer model for PageRank maps perfectly to how users read these books. Thus, they look for the most probable storylets in a book. For this problem, the graphs are directed and acyclic, the stochastic matrix is normalized by outdegree, and we have a standard PageRank problem. They are careful to model a weakly preferential PageRank system that deterministically transitions from a terminal (or dangling) storylet back to the start of

the book. Teleporting is uniform in their experiments. They find that both PageRank and a ranking system they derive give useful information about the properties of these stories.

Books & tags: BookRank. Traditional library catalogs use a carefully curated set of index terms to indicate the contents of books. These enabled content-based search prior to the existence of fast full-text search engines. Social cataloging sites such as LibraryThing and Shelfari allow their users to curate their own set of index terms for books that they read, and easily share this information among the user sites. The data on these websites consists of *books* and *tags* that indicate the topics of books. BookRank, which is localized PageRank on the bipartite book-tag graph [Meng, 2009], produces eerily accurate suggestions for what to read next. For instance, if we use teleportation to localize on Golub and van Loan’s text “Matrix Computations”, Boyd and Vandenberghe’s “Convex Optimization”, and Hastie, Tibshirani, and Friedman’s “Elements of Statistical Learning”, then the top suggestion is a book on Combinatorial Optimization by Papadimitriou and Steiglitz. A similar idea underlies the general FolkRank system [Hotho et al., 2006] that we’ll see shortly (Section 4.9).

4.8. PageRank in bibliometrics: TimedPageRank, CiteRank, AuthorRank. The field of bibliometrics is another big producer and consumer of network ranking methods, starting with seminal work by Garfield on aggregating data into a citation network between journals [Garfield, 1955; Garfield and Sher, 1963] and proceeding through Pinski and Narin [1976], who defined a close analogue of PageRank. In almost all of these usages, PageRank is used as a centrality measure to reveal the most important journals, papers, and authors.

Citations among journals. The citation network Garfield originally collected and analyzed is the journal-journal citation network. It is a weighted network where each node is a journal and each edge is the number of citations between articles of the journals. ISI’s impact factor is a more refined analysis of these citation patterns. Bollen et al. [2006] takes ISI’s methods a step further and finds that a combination of the impact factor with the PageRank value in the journal citation produces a ranked list of journals that better correlates with experts’ judgements. PageRank is used as a centrality measure here with uniform teleportation and weights that correspond to the weighted citation network. The graph had around 6000 journals. The Eigenfactor system [West et al., 2010] uses a PageRank vector on the journal co-citation network with uniform teleportation and $\alpha = 0.85$ to measure the influence of a journals. It also shows these rankings on easy-to-browse website.

Citations among papers: TimedPageRank, CiteRank. Moving beyond individual journals, we can also study the citation network among individual papers using PageRank. In a paper citation network, each node is an individual article and the edges are directed based on the citation. Modern bibliographic and citation databases such as arXiv and DBLP make these networks easy to construct. They tend to have hundreds of thousands of nodes and a few million edges. TimedPageRank is an idea to weight the edges of the stochastic matrix in PageRank such that more recent citations are more important. Formally, it is the solution of

$$(\mathbf{I} - \alpha \mathbf{A}^T \mathbf{D}^{-1} \mathbf{W}) \mathbf{x} = (1 - \alpha) \mathbf{e}$$

where \mathbf{W} is a diagonal matrix with weights between 0 and 1 that reflects the age of the paper (1 is recent and 0 is old). The matrix $\mathbf{A}^T \mathbf{D}^{-1} \mathbf{W}$ is column sub-stochastic and so this is a pseudo-PageRank problem. CiteRank is a subsequent idea that uses the teleportation in PageRank to increase the rank of recent articles [Walker et al.,

2007]. Thus, v_i is smaller if paper i is older and v_i is larger if paper i is more recent. The goal of both methods is to produce temporally relevant orderings that remove the bias of older articles to acquire citations.

While the previous two papers focused on how to make article importance more accurate, Chen et al. [2007] attempts to use PageRank in concert with the number of citations to find *hidden gems*. One notable contribution is the study of α in citation analysis: based on a heuristic argument about how we build references for an article, they recommend $\alpha = 0.5$. Moreover, they find papers with higher PageRank scores than would be expected given their citation count. These are the hidden gems of the literature. Ma et al. [2008] uses the same idea in a larger study and find a similar effect.

Citations among authors: AuthorRank. Another type of bibliographic network is the co-authorship graph. For each paper, insert edges among all co-authors. Thus, each paper becomes a clique in the co-authorship network. The weights on each edge are either uniform (and set to 1), based on the number of papers co-authored, or based on another weighting construction defined in that paper. All of these constructions produce an undirected network. PageRank on this network gives a practical ranking of the most important authors [Liu et al., 2005]. The teleportation is uniform with $\alpha = 0.85$, or can be focused on a subset of authors to generate an area-specific ranking. Their data have a few thousand authors. These graphs are constructions based on an underlying bipartite matrix \mathbf{B} that relates authors and papers. More specifically, the weighted co-authorship network is the matrix $\mathbf{B}\mathbf{B}^T$. Many such constructions can be related back to the matrix $\begin{bmatrix} 0 & \mathbf{B} \\ \mathbf{B}^T & 0 \end{bmatrix}$ [Dhillon, 2001]. We are not aware of any analysis that makes a relationship between PageRank in the bipartite graph $\begin{bmatrix} 0 & \mathbf{B} \\ \mathbf{B}^T & 0 \end{bmatrix}$ and the weighted matrix $\mathbf{B}\mathbf{B}^T$.

Author, paper, citation networks. Citation analysis and co-authorship analysis can, of course, be combined, and that is exactly what Fiala et al. [2008] and Jezek et al. [2008] do. Whereas Liu et al. [2005] study the co-authorship network, here, they study a particular construction that joins the bipartite author-paper network to the citation network to produce an author-citation network. This is a network where author i links to author j if i has a paper that cites j where j is not a co-author on that particular paper. Using $\alpha = 0.9$ and uniform teleportation produces another helpful list of the most important authors. In the notation of the previous paragraph, a related construction is the network with adjacency matrix

$$\mathbf{A} = \begin{bmatrix} 0 & \mathbf{B} \\ \mathbf{B}^T & \mathbf{C} \end{bmatrix},$$

where \mathbf{B} is the bipartite author-paper matrix and \mathbf{C} is the citation matrix among papers. PageRank on these networks takes into account both the co-authorship and directed citation information, and it rewards authors that have many, highly cited papers. The graphs studied have a few hundred thousand authors and author-author citations.

4.9. PageRank in databases and knowledge information systems: PopRank, FactRank, ObjectRank, FolkRank. Knowledge information systems store codified forms of information, typically as a relational database. For instance, a knowledge system about movies consists of relationships between actors, characters, movies, directors, and so on. Contemporary information systems also often contain large bodies of user-generated content through tags, ratings, and such. Ratings are a sufficiently special case that we review them in a forthcoming section (Section 4.10),

but we will study PageRank and tags here. PageRank serves important roles as both a centrality measure and localized measure in networks derived from a knowledge system. We'll also present slightly more detail on four interesting applications.

Centrality scores: PopRank, FactRank. PageRank's role as a centrality measure in a knowledge information system is akin to its role on the web as an importance measure. For instance, the authors of PopRank [Nie et al., 2005] consider searching through large databases of objects – think of academic papers – that have their own internal set of relationships within the knowledge system – think of co-author relationships. But these papers are also linked to by websites. PopRank uses web-importance as a teleportation vector for a PageRank vector defined on the set of object relationships. The result is a measure of object popularity biased by its web popularity. One of the challenges in using such a system is that collecting good databases of relational information is hard. FactRank helps with this process [Jain and Pantel, 2010]. It is a measure designed to evaluate the importance and accuracy of a fact network. A fact is just a sentence that connects two objects, such as “David-Gleich wrote the-paper PageRank-Beyond-The-Web.” These sentences come from textual analysis of large web crawls. In a fact network, facts are connected if they involve the same set of objects. Variations on PageRank with uniform teleportation provide lists of important facts. The authors of FactRank found that weighting relationships between facts and using PageRank scores of this weighted network gave higher performance than both a baseline and standard PageRank method in the task of finding correct facts. The fact networks are undirected and have a few million nodes.

Localized scores: Random-walk with restart, Semi-Supervised Learning. Prediction tasks akin to the bioinformatics usages of PageRank are standard within knowledge information systems: networks contain noisy relationships, and the task is inferring, or predicting, missing data based on these relationships. Zhou et al. [2003] used a localized PageRank computation to infer the identity of handwritten digits from only a few examples. These problems were called *semi-supervised learning on graphs* because they model the case of finding a vector over vertices (or learning a function) based on a few values of the function (supervised). It differs from the standard supervised learning problem because the graph setup implies that only predictions on the vertices are required, instead of the general prediction problem with arbitrary future inputs. In the particular study, the graph among these images is based on a radial basis function construction. For this task $\alpha = 0.99$ in the pseudo-PageRank system $(\mathbf{I} - \alpha\mathbf{P})\tilde{\mathbf{Y}} = \mathbf{S}$, where \mathbf{S} is a binary matrix indicating *known samples* $S_{ij} = 1$ if image i is known to be digit j . The largest value in each row of $\mathbf{Y} = \mathbf{D}\tilde{\mathbf{Y}}$ gives the predicted digit for any unknown image. While these graphs were undirected, later work [Zhou et al., 2005] showed how to use PageRank with global teleportation, in concert with symmetric Laplacian structure defined on a directed graph [Chung, 2005], to enable the same methodology on a general directed graph.

Pan et al. [2004] define a random walk with restart, which is exactly a personalized PageRank system, to infer captions for a database of images. Given a set of images labeled by captions, define a graph where each image is connected to its regions, each region is connected to other regions via a similarity function, and each image is connected to the terms in its caption. A query image is a distribution over regions, and we find terms by solving a PageRank problem with this as the teleportation vector. These graphs are weighted, undirected graphs. Curiously, the authors chose α based on experimentation and found that $\alpha = 0.1$ or $\alpha = 0.2$ works best. They attribute the difference to the incredibly small diameter of their network. Subsequent work in the

same vein showed some of the relationships with the normalized Laplacian matrix of a graph [Tong et al., 2006] and returned to a larger value of α around 0.9.

Application 1 – Database queries: ObjectRank. ObjectRank is an interesting type of database query [Balmin et al., 2004]. A typical query to a database will retrieve all of the rows of a specified set of tables matching a precise criteria, such as, “find all students with a GPA of 3.5 that were born in Minnesota.” These tables often have internal relationships – the database schema – that would help determine which are the most important returned results. In the ObjectRank model, a user queries the database with a textual term. The authors describe a means to turn the database objects and schema into a sub-stochastic transition matrix and define ObjectRank as the query-dependent solution of the PageRank linear system where the teleportation vector reflects textual matches. They suggest a great deal of flexibility with defining the weights of this matrix. For instance, there may be no natural direction for many of these links and the authors suggest differently weighting forward edges and backward edges – their intuition is that a paper cited by many important papers is itself important, but that citing important papers does not transfer any importance. They use $\alpha = 0.85$ and the graphs have a few million edges.

Application 2 – Folksonomy search: FolkRank. A more specific situation is folksonomy search. A folksonomy is a collection of objects, users, and tags. Each entry is a triplet of these three items. A user such as myself may have tagged a picture on the flickr network with the term “sunset” if it contained a sunset, thus creating the triplet (picture,user,“sunset”). FolkRank scores [Hotho et al., 2006] are designed to measure the importance of an object, tag, or user with respect to a small set of objects, tags, or users that define a topic. (This idea is akin to topic-sensitive PageRank, Haveliwala 2002.) These scores then help reveal important objects *related* to a given search, as well as the tags that relate them. The scores are based on localized PageRank scores from an undirected, tripartite weighted network. There is a wrinkle, however. The FolkRank scores are taken as the *difference* between a PageRank vector computed with $\alpha = 1$ and $\alpha = 1/2$. The graph is undirected, so the solution with $\alpha = 1$ is just the weighted degree distribution. Thus, FolkRank downweights items that are important for *everyone*.

Application 3 – Semantic relatedness. The Open Directory Project, or ODP, is a hierarchical, categorical index of web-pages that organizes them into related groups. Bar-Yossef and Mashiach [2008] suggests a way of defining the relatedness of two categories on ODP using their localized PageRank scores. The goal is to generalize the idea of the least-common ancestor to random walks to give a different sense of the distance between categories. To do so, create a graph from the directed hierarchy in the ODP. Let \mathbf{x} be the reverse PageRank vector that teleports back to a single category, and let \mathbf{y} be the reverse PageRank vector that teleports back to another (single) category. Then the relatedness of these categories is the cosine of the angle between \mathbf{x} and \mathbf{y} . Let \mathbf{x} be the localized PageRank vector (Note the use of reverse PageRank here so that edges go from child to parent.) They show evidence that this is a useful measure of relationship in ODP.

Application 4 – Logic programming. A fundamental challenge with scaling logic programming systems like Prolog is that there is an exponential explosion of potential combinations and rules to evaluate and, unless the system is extremely well-designed, these cannot be pruned away. This limits applications to almost trivial problems. Internally, Prolog-type systems resolve, or prove, logical statements using a search procedure over an implicitly defined graph that may be infinite. At each node of the

graph, the proof system generates all potential neighbors of the node by applying a rule set given by the logic system. Thus, given one node in the graph, the search procedure eventually visits all nodes. Localized PageRank provides a natural way to restrict the search space to only “short” and “likely” proofs [Wang et al., 2013]. Formally, they use PageRank’s random teleportation to control the expansion of the search procedure. However, there is an intuitive explanation for the random restarts in such a problem: periodically we all abandon our current line of attack in a proof and start out fresh. Their system with localized PageRank allows them to realize this behavior in a rigorous way.

4.10. PageRank in recommender systems: ItemRank. A recommender system attempts to predict what its users will do based on their past behavior. Netflix and Amazon have some of the most famous recommendation systems that predict movies and products, respectively, their users will enjoy. Localized PageRank helps to score potential predictions in many research studies on recommender systems.

Query reformulation. A key component of modern web-search systems is predicting future queries. Boldi et al. [2008] run localized PageRank on a query reformulation-graph that describes how users rewrite queries with $\alpha = 0.85$. Two queries, q_1 and q_2 , are connected in this graph if a user searched for q_1 before q_2 within a close time-frame and both q_1 and q_2 have some non-trivial textual relationships. This graph is directed and weighted. The teleportation vector is localized on the current query, or a small set of previously used terms. PageRank has since had great success for many tasks related to query suggestion and often performs among the best methods evaluated [Song et al., 2012].

Item recommendation: ItemRank. Both Netflix and Amazon’s recommender systems are called *item recommendation* problems. Users rate items – typically with a 5-star scale – and we wish to recommend items that a user will rate highly. The ratings matrix is an items-by-users matrix where R_{ij} is the numeric rating given to item i by user j . These ratings form a bipartite network between the two groups and we collapse this to a graph over items as follows. Let G be a weighted graph where the weights on an edge (i, j) are the number of users that rated both items i and j . (These weights are equivalent to the number of paths of length 2 between each pair of items in terms of the bipartite graph.) Let \mathbf{P} be the standard weighted random walk construction on G . Then the ItemRank scores [Gori and Pucci, 2007] are the solutions of:

$$(\mathbf{I} - \alpha\mathbf{P})\mathbf{S} = (1 - \alpha)\mathbf{R}\mathbf{D}_{\mathbf{R}}^{-1}$$

where $\mathbf{D}_{\mathbf{R}}$ are column sums of the rating matrix. Each column of \mathbf{S} is a set of recommendations for user j , and S_{ij} is a proxy for the interest of user j in item i . Note that any construction of the transition matrix \mathbf{P} based on correlations between items based on user ratings would work in this application as well.

Link prediction. Given the current state of a network, link prediction tries to predict which edges will come into existence in the future. Liben-Nowell and Kleinberg [2006] evaluated the localized PageRank score of an unknown edge in terms of its predictive power. These PageRank values were entries in the matrix $(\mathbf{I} - \alpha\mathbf{P})^{-1}$ for edges that currently do not exist in the graph. PageRank with α between 0.5 and 0.99 was not one of their best predictors, but the Katz matrix $(\mathbf{I} - \alpha\mathbf{A})^{-1}$ was one of the best with $\alpha = 0.0005$. Note that Katz’s matrix is, implicitly, a pseudo-PageRank problem if $\alpha < \frac{1}{d_{\max}}$ where d_{\max} is the largest degree in the graph. The co-authorship graphs tested seem to have had degrees less than 2000, making this hidden pseudo-PageRank

problem one of the best predictors of future co-authorship. More recent work using PageRank for predicting links on the Facebook social network includes a training phase to estimate weights of the matrix \mathbf{P} to achieve higher prediction [Backstrom and Leskovec, 2011]. Localized PageRank is believed to be part of Twitter’s follower suggestion scheme too [Bahmani et al., 2010].

4.11. PageRank in social networks: BuddyRank, TwitterRank. PageRank serves three purposes in a social network, where the nodes are people and the edges are some type of social relationship. First, as we discussed in the previous section, it can help solve link prediction problems to find individuals that will become friends soon. Second, it serves a classic role in evaluating the centrality of the people involved to estimate their social status and power. Third, it helps evaluate the potential influence of a node on the opinions of the network.

Centrality: BuddyRank. Centrality methods have a long history in social networks – see Katz [1953] and Vigna [2009] for a good discussion. The following claim is difficult to verify, but we suspect that the first use of PageRank in a large-scale social network was the BuddyRank measure employed by BuddyZoo in 2003.¹ BuddyZoo collected contact lists from users of the AOL Instant Messenger service and assembled them into one of the first large-scale social networks studied via graph theoretic methods. Since then, PageRank has been used to rank individuals in the Twitter network by their importance [Java, 2007] and to help characterize properties of the Twitter social network by the PageRank values of their users [Kwak et al., 2010]. These are standard applications of PageRank with global teleportation and $\alpha \approx 0.85$.

Influence. Finding influential individuals is one of the important questions in social network analysis. This amounts to finding nodes that can spread their influence widely. More formalizations of this question result in NP-hard optimization problems [Kempe et al., 2003] and thus, heuristics and approximation algorithms abound [Kempe et al., 2003, 2005]. Using Reverse PageRank with global teleportation as a heuristic outperforms out-degree for this task, as shown by Java et al. [2006] for web-blog influence and Bar-Yossef and Mashiach [2008] for the social network LiveJournal. Reverse PageRank, instead of traditional PageRank, is the correct model to understand the *origins* of influence – the distinction is much like the treatment of hubs and authorities in other ranking models on networks [Kleinberg, 1999; Blondel et al., 2004]. These ideas also extend to finding topical authorities in social networks by using the teleportation vector and topic-specific transition probabilities to localize the PageRank vector in TwitterRank [Weng et al., 2010].

4.12. PageRank in the web, redux: HostRank, DirRank, TrustRank, BadRank, VisualRank. At the conclusion of our survey of applications, we return to uses of PageRank on the web itself. Before we begin, let us address the elephant in the room, so to speak. Does Google still use PageRank? Google reportedly uses a basket of ranking metrics to determine the final order that results are returned. These evolve continuously and vary depending on where and when you are searching. It is unclear to what extent PageRank, or more generally, link analysis measures play a role in Google’s search ordering, and this is a closely guarded secret unlikely to be known outside of an inner-circle at Google. One the one hand, in perhaps the only large-scale published study on PageRank’s effectiveness in a search engine, Najork et al. [2007] found that it underperformed in-degree. On the other hand, PageRank is still widely believed to still play some role based on statements from Google. For instance,

¹<http://web.archive.org/web/20050724231459/http://buddyzoo.com/>

Matt Cutts, a Google engineer, wrote about how Google uses PageRank to determine crawling behavior [Cutts, 2006], and later wrote about how Google moved to a full substochastic matrix in terms of their PageRank vector [Cutts, 2009]. The latter case was designed to handle a new class of link on the web called `rel=nofollow`. This was an optional HTML parameter that would tell a crawler that the following link is not useful for relevance judgements. All the major web companies created this parameter to combat links created in the comment sections of extremely high quality pages such as the Washington Post. These links are created by users of the Washington Post, not the staff themselves, and shouldn't constitute an endorsement on a page. Cutts described how Google's new PageRank equation would count these `rel=nofollow` links in the *degree* of a node when it was computing a stochastic normalization, but would remove the links when computing relevance. For instance, if my page had three true links and two `rel=nofollow` links, then my true links would have probabilities $1/5$ instead of $1/3$, and the sum of my outgoing probability would be $3/5$ instead of 1. Thus, Google's PageRank computation is a pseudo-PageRank problem now.

Outside of Google's usage, PageRank is also used to evaluate the web at coarser levels of granularity through HostRank and DirRank. Reverse PageRank provides a good measure of a page's similarity to a hub, according to both Fogaras [2003] and Bar-Yossef and Mashiach [2008]. PageRank and reverse PageRank also provide information on the "spaminess" of particular pages through metrics such as TrustRank and BadRank. PageRank-based information also helped to identify spam directly in a study by Becchetti et al. [2008]. Finally, PageRank helps identify canonical images to place on a web-search result (VisualRank).

Coarse PageRank: HostRank, DirRank. Arasu et al. [2002] was an important early paper that defined HostRank, where the web is aggregated at the level of hostnames. In this case, all links to and from a hostname, such as `www.cs.purdue.edu`, become equivalent. This particular construction models a random surfer that, when visiting a page, makes a censored, or silent, transition within all pages on the same host, and then follows a random link. The HostRank scores are the sums of these modified PageRank scores on the pages within each host [Gleich and Polito, 2007]. Later work included BlockRank [Kamvar et al., 2003], which used HostRank to initialize PageRank, and DirRank [Eiron et al., 2004], which forms an aggregation at the level of *directories* of websites.

Trust, Reputation, & Spam: TrustRank, BadRank. PageRank typically provides authority scores to estimate the importance of a page on the web. As the commercial value of websites grew, it became highly profitable to create *spam* sites that contain no new information content but attempt to capture Google search results by appearing to contain information. BadRank [Sobek, 2003] and TrustRank [Gyöngyi et al., 2004] emerged as new, link analysis tools to combat the problem. Essentially, these ideas solve localized, reverse PageRank problems. The results are either used directly, or as a "safe teleportation" vector for PageRank, as in TrustRank, or in concert with other techniques, as likely done in BadRank. Kolda and Procopio [2009] generalizes these models and includes the idea of adding self-links to fix the dangling nodes, like in sink preferential PageRank, but they add them everywhere, not just at dangling nodes. For spam-link applications, this way of handling dangling nodes is superior – in a modeling sense – to the alternatives.

Wikipedia. Wikipedia is often used as a subset of the web for studying ranking. It is easy to download the data for the entire website, which makes building the web-graph convenient. (A crawl from a few years ago is in the sparse matrix repository, Davis and

Hu 2010, as the matrix `Gleich/wikipedia-20070206`.) Current graphs of the English language pages have around 100,000,000 links and 10,000,000 articles. The nature of the pages on Wikipedia also makes it easy to evaluate results anecdotally. For instance, we would all raise an eyebrow and demand explanation if “Gene Golub” was the page with highest global PageRank in Wikipedia. On the other hand, this result might be expected if we solve a localized PageRank problem around the Wikipedia article for “numerical linear algebra.” Wissner-Gross [2006] used Wikipedia as a test set to build reading lists using a combination of localized and global PageRank scores. Later, Zhironov et al. [2010] computed a 2d ranking on Wikipedia by combining global PageRank and reverse PageRank. Finally, this 2d ranking showed that Frank Sinatra was one of the most important people [Eom et al., 2014].

Image search: VisualRank. PageRank also helps to identify “canonical” images to display as a visual summary of a larger set of images returned from an image search engine. In the VisualRank system, Jing and Baluja [2008] compute PageRank of an image similarity graph generated from an image search result. The graphs are small – around 1000 nodes – which reflects the standard textual query results, and they are also symmetric and weighted. They solve a global PageRank problem with uniform teleportation or high-result biased teleportation. The highest ranked images are canonical images of Mona Lisa amid a diverse collection of views.

5. PageRank generalizations. Beyond the applications discussed so far, there is an extremely wide set of PageRank-like models that do not fit into the canonical definition and constructions from Section 3. These support a wide range of additional applications with mathematics that differs slightly, and some of them are formal mathematical generalizations of the PageRank vectors. For instance, in prior work, we studied PageRank with a random teleportation parameter [Constantine and Gleich, 2010]. The standard deviation of these vectors resulted in increased accuracy in detecting spam pages on the web. We now survey some of these formal generalizations.

5.1. Diffusions, damped sums, \mathcal{E} heat kernels. Recall that the pseudo-PageRank vector is the solution of (2.3),

$$(\mathbf{I} - \alpha \bar{\mathbf{P}})\mathbf{y} = \mathbf{f}.$$

Since all of the eigenvalues of $\bar{\mathbf{P}}$ are bounded by 1 in magnitude, the solution \mathbf{y} has an expansion in terms of the Neumann series:

$$\mathbf{y} = \sum_{k=0}^{\infty} \alpha^k \bar{\mathbf{P}}^k \mathbf{f}.$$

This expressions gives the pseudo-PageRank vector as a damped sum of powers of $\bar{\mathbf{P}}$ where each power, $\bar{\mathbf{P}}^k$, has the geometrically decaying weight α^k . These are often called *damped diffusions* because this equation models how the quantities in \mathbf{f} probabilistically diffuse through the graph where the probability of a path of length k is damped by α^k . Many other sequences serve the same purpose as pointed out by a variety of authors.

Generalized damping. Perhaps the most general setting for these ideas is the generalized damped PageRank vector:

$$\mathbf{z} = \sum_{k=0}^{\infty} \gamma_k \bar{\mathbf{P}}^k \mathbf{f} \tag{5.1}$$

where γ_k is a non-negative ℓ_1 -sequence (that is, $\sum_k \gamma_k < \infty$ and $\gamma_k \geq 0$). This reduces to PageRank if $\gamma_k = \alpha^k$. Huberman et al. [1998] suggested using such a construction where γ_k arises from real-world path following behaviors on the web, which they found to resemble inverse Gaussian functions. Later results from Baeza-Yates et al. [2006] proposed essentially the same formula in (5.1). They suggested a variety of interesting functions γ_k , including some with only a finite number of non-zero terms. These authors drew their motivation from the earlier work of TotalRank [Boldi, 2005], which suggested $\gamma_k = \frac{1}{k+1} - \frac{1}{k+2}$ in order to evaluate the TotalRank vector:

$$\mathbf{z} = \int_0^1 (\mathbf{I} - \alpha \bar{\mathbf{P}})^{-1} (1 - \alpha) \mathbf{v} d\alpha.$$

This integrates over all possible values of α . (As an aside, this integral is well defined because a unique limiting PageRank value exists at $\alpha = 1$, see Section 5.2. This sidesteps a technical issue with the singular matrix at $\alpha = 1$.) Our work with making the value of α in PageRank a random variable is really a further generalization [Constantine and Gleich, 2010]. Let $\mathbf{x}(\alpha)$ be a parameterized form for the PageRank vector for a fixed graph and teleportation vector. Let A be a random variable supported on $[0, 1]$ with an infinite number of finite moments, that is, $E[A^k] < \infty$ for all k . Intuitively, A is the probability that a random user of the web follows a link. Our idea was to use the expected value of PageRank $E[\mathbf{x}(A)]$ to produce a ranking that reflected the distribution of path-following behaviors in the random surfers. We showed:

$$E[\mathbf{x}(A)] = \sum_{k=0}^{\infty} (E[A^k] - E[A^{k+1}]) \mathbf{P}^k \mathbf{v}.$$

This results in a family of sequences of γ_k that depend on the random variable A . Recent work by Kollias et al. [2013] shows how to evaluate these generalized damped vectors as a polynomial combination of PageRank vectors in the sense of (2.2).

Heat kernels & matrix exponentials. Another specific case of generalized damping arises from the matrix exponential, or heat kernel:

$$\mathbf{z} = e^{\beta \bar{\mathbf{P}}} \mathbf{f} = \sum_{k=0}^{\infty} \frac{\beta^k}{k!} \bar{\mathbf{P}}^k \mathbf{f}.$$

Such functions arose in a wide variety of domains that would be tangential to review here [Estrada, 2000; Miller et al., 2001; Kondor and Lafferty, 2002; Farahat et al., 2006; Chung, 2007; Kunegis and Lommatzsch, 2009; Estrada and Higham, 2010]. In terms of a specific relationship with PageRank, Yang et al. [2007] noted that the pseudo-PageRank vector itself was a single-term approximation to these heat kernel diffusions. Consider

$$\mathbf{z} = e^{\beta \bar{\mathbf{P}}} \mathbf{f} \quad \Leftrightarrow \quad e^{-\beta \bar{\mathbf{P}}} \mathbf{z} = \mathbf{f} \quad \Leftrightarrow \quad (\mathbf{I} - \beta \bar{\mathbf{P}} + \dots) \mathbf{z} = \mathbf{f}.$$

If we truncate the heat kernel expansion after just the first two terms ($\mathbf{I} - \beta \bar{\mathbf{P}}$), then we arise at the pseudo-PageRank vector. (A similar result holds for the formal PageRank vector too.)

5.2. PageRank limits & eigenvector centrality. In the definition of PageRank used in this paper, we assume that $\alpha < 1$. PageRank, however, has a unique well-defined limit as $\alpha \rightarrow 1$ [Serra-Capizzano, 2005; Boldi et al., 2005, 2009b]. This

is easy to prove using the Jordan canonical form for the case of PageRank (2.2), but extensions to pseudo-PageRank are slightly more nuanced. As in the previous section, let $\mathbf{x}(\alpha)$ be the PageRank vector as a function of α for a fixed stochastic \mathbf{P} : $(\mathbf{I} - \alpha\mathbf{P})\mathbf{x}(\alpha) = (1 - \alpha)\mathbf{v}$. Let $\mathbf{X}\mathbf{J}\mathbf{X}^{-1}$ be the Jordan canonical form of \mathbf{P} . Because \mathbf{P} is stochastic, it's eigenvalues on the unit circle are all semi-simple [Meyer, 2000, page 696]. Thus:

$$\mathbf{J} = \begin{bmatrix} \mathbf{I} & & \\ & \mathbf{D}_1 & \\ & & \mathbf{J}_2 \end{bmatrix},$$

where \mathbf{D}_1 is a diagonal matrix of the eigenvalues on the unit circle and \mathbf{J}_2 is a Jordan block for all eigenvalues with $|\lambda| < 1$. We now substitute this into the PageRank equation:

$$(\mathbf{I} - \alpha\mathbf{P})\mathbf{x}(\alpha) = (1 - \alpha)\mathbf{v} \Leftrightarrow (\mathbf{I} - \alpha\mathbf{J})^{-1} \underbrace{\hat{\mathbf{x}}(\alpha)}_{\mathbf{X}^{-1}\mathbf{x}(\alpha)} = (1 - \alpha) \underbrace{\hat{\mathbf{v}}}_{\mathbf{X}^{-1}\mathbf{v}}.$$

Using the structure of \mathbf{J} decouples these equations:

$$\left(\begin{bmatrix} \mathbf{I} & & \\ & \mathbf{I} & \\ & & \mathbf{I} \end{bmatrix} - \alpha \begin{bmatrix} \mathbf{I} & & \\ & \mathbf{D}_1 & \\ & & \mathbf{J}_2 \end{bmatrix} \right) \begin{bmatrix} \hat{\mathbf{x}}(\alpha)_0 \\ \hat{\mathbf{x}}(\alpha)_1 \\ \hat{\mathbf{x}}(\alpha)_2 \end{bmatrix} = (1 - \alpha) \begin{bmatrix} \hat{\mathbf{v}}_0 \\ \hat{\mathbf{v}}_1 \\ \hat{\mathbf{v}}_2 \end{bmatrix}.$$

As $\alpha \rightarrow 1$, both $\hat{\mathbf{x}}(\alpha)_1$ and $\hat{\mathbf{x}}(\alpha)_2$ go to 0 because these linear systems remain non-singular. Also, note that $\hat{\mathbf{x}}(\alpha)_0 = \hat{\mathbf{v}}_1$ for all $\alpha \neq 1$, so this point is a removable singularity. Thus, $\hat{\mathbf{x}}$ can be uniquely defined at $\alpha = 1$, and hence, so can \mathbf{x} . Vigna [2005] uses the structure of this limit to argue that taking $\alpha \rightarrow 1$ in practical applications is not useful unless the underlying graph is strongly connected, and they propose a new PageRank construction to ensure this property. Subsequent work by Vigna [2009] does a nice job of showing how limiting cases of PageRank vectors converge to traditional eigenvector centrality measures from bibliometrics [Pinski and Narin, 1976] and social network analysis [Katz, 1953].

The pseudo-PageRank problem does not have nice limiting properties in our formulation. Let $\mathbf{y}(\alpha)$ be a parametric form for the solution of the pseudo-PageRank system $(\mathbf{I} - \alpha\bar{\mathbf{P}})\mathbf{y} = \mathbf{f}$. As $\alpha \rightarrow 1$, then $\mathbf{y} \rightarrow \infty$, unless the non-zero support of \mathbf{f} lies outside of a recurrent class, in which case $\mathbf{y} \rightarrow 0$. Boldi et al. [2005] defines the PseudoRank system as:

$$(\mathbf{I} - \alpha\bar{\mathbf{P}})\mathbf{y} = (1 - \alpha)\mathbf{f}$$

instead. This system always has a non-infinite limit as $\alpha \rightarrow 1$. It could, however, have zero as a limit if $\bar{\mathbf{P}}$ has all eigenvalues less than 1.

5.3. Over-teleportation, negative teleportation, \mathcal{E} the Fiedler vector.

The next generalization of PageRank is to values of $\alpha > 1$. These arose in our prior work to understand the convergence of quadrature formulas for approximating the expected value of PageRank with random teleportation parameters [Constantine and Gleich, 2010]. Mahoney et al. [2012] subsequently showed an amazing relationship among (i) the Fiedler vector of a graph [Fiedler, 1973; Anderson and Morley, 1985; Pothén et al., 1990], (ii) a particular generalization of the PageRank vector, which we call MOV, and (iii) values of $\alpha > 1$.

The Fiedler vector. In contrast to the remainder of this paper, the constructions and statements in this section are specific to connected, undirected graphs with symmetric adjacency matrices. The conductance of a set of vertices in a graph is defined as the number of edges leaving that set, divided by the sum of the degrees of the vertices within the set. Conductance and its relatives are often used as numeric quality scores for graph partitioning in parallel processing [Pothen et al., 1990] and for community detection in graphs [Schaeffer, 2007]. It is NP-hard to find the set of smallest conductance, but Fiedler’s vector reveals information about the presence of small conductance sets in a graph through the Cheeger inequality [Chung, 1992]. Let G be a connected, undirected graph with symmetric adjacency matrix \mathbf{A} and diagonal degree matrix \mathbf{D} . The Fiedler vector is the generalized eigenvector of $(\mathbf{D}-\mathbf{A})\mathbf{q} = \lambda_*\mathbf{D}\mathbf{q}$, with the smallest positive eigenvalue $\lambda_* > 0$. All of the generalized eigenvalues are non-negative, the smallest is 0, and the largest is bounded above by 1. Cheeger’s inequality bounds the relationship between λ_* and the set of smallest conductance in the graph.

MOV. The MOV vector is defined as the pseudo-inverse solution \mathbf{r} in the consistent linear system of equations:

$$[(\mathbf{D} - \mathbf{A}) - \gamma\mathbf{D}]\mathbf{r} = \rho(\gamma)\mathbf{D}\mathbf{s}, \quad (5.2)$$

where $\gamma < \lambda_*$, \mathbf{s} is a “seed” vector such that $\mathbf{s}^T\mathbf{D}\mathbf{e} = 0$, and $\rho(\gamma)$ is a scaling constant such that \mathbf{r} has a fixed norm. When $\gamma = 0$, this system is singular but consistent, and thus, we take the pseudo-inverse solution. Note that this is equivalent to the pseudo-PageRank problem:

$$(\mathbf{I} - \alpha\mathbf{P})\mathbf{z} = \alpha\rho(\gamma)\hat{\mathbf{f}}$$

where $\alpha = \frac{1}{1-\gamma}$, $\mathbf{z} = \mathbf{D}\mathbf{r}$, and $\hat{\mathbf{f}} = \mathbf{D}\mathbf{s}$. The properties of \mathbf{s} in MOV imply that $\hat{\mathbf{f}}^T\mathbf{e} = 0$, and thus, $\hat{\mathbf{f}}$ must have negative elements, which generalizes the standard pseudo-PageRank.

In a small surprise, allowing \mathbf{f} to take on negative values results in no additional modeling power in the case of symmetric \mathbf{A} . To establish this result, we first observe that:

$$(\mathbf{I} - \alpha\mathbf{A}\mathbf{D}^{-1})\frac{\sigma}{1-\alpha}\mathbf{d} = \sigma\mathbf{d}.$$

This preliminary fact shows that the pseudo-PageRank vector of an undirected graph with teleportation according to the degree vector \mathbf{d} simply results in a rescaling. We can use this property to shift any \mathbf{f} with negative values in a controlled manner:

$$(\mathbf{I} - \alpha\mathbf{P})\underbrace{\left(\mathbf{z} + \frac{\sigma}{1-\alpha}\mathbf{d}\right)}_{\mathbf{y}} = \underbrace{\alpha\hat{\mathbf{f}} + \sigma\mathbf{d}}_{\mathbf{f}},$$

where σ is chosen such that $\mathbf{f} \geq 0$ element-wise. Solving these shifted pseudo-PageRank systems, then, effectively computes the solution \mathbf{z} with a well-understood bias term $\theta\mathbf{d}$. This is easy to remove afterwards: $\mathbf{z} = \mathbf{y} - \frac{\sigma}{1-\alpha}\mathbf{d}$, at which point we can normalize \mathbf{z} to account for $\rho(\gamma)$ if desired.

Values of $\alpha > 1$. While this generalization with negative entries in \mathbf{f} gives no additional mathematical power, it does permit a seamless limit from PageRank vectors to the Fiedler vector. Let $\alpha_* = \frac{1}{1-\lambda_*} > 1$. The formal result is that the limit $\lim_{\alpha \rightarrow \alpha_*} \frac{1}{\rho(\alpha)}\mathbf{z}(\alpha) = \mathbf{q}$, the Fiedler vector. Note that for the construction of

$\mathbf{P} = \mathbf{A}\mathbf{D}^{-1}$ on an undirected, connected graph, we have that $\mathbf{P}^k \rightarrow \frac{1}{e^T \mathbf{d}} \mathbf{d} \mathbf{e}^T$ as $k \rightarrow \infty$. Thus, when $\alpha = 1$, the MOV solution \mathbf{z} is equivalent to the solution of $(\mathbf{I} - (\mathbf{P} - \frac{1}{e^T \mathbf{d}} \mathbf{d} \mathbf{e}^T)) \mathbf{z} = \mathbf{f}$ because the right hand side of \mathbf{f} is orthogonal to the left eigenvector \mathbf{e}^T . As all of the eigenvalues of $(\mathbf{P} - \frac{1}{e^T \mathbf{d}} \mathbf{d} \mathbf{e}^T)$ are distinct from 1, this is a non-singular system. And this fact allows the limit construction to pass through $\alpha = 1$ seamlessly. If we additionally assume that $\mathbf{f}^T \mathbf{q} \neq 0$, then

$$\lim_{\alpha \rightarrow \alpha_*} \frac{1}{\rho(\alpha)} \mathbf{z}(\alpha) = \mathbf{q},$$

and the limiting value of PageRank with over-teleportation is the Fiedler vector. The analysis in Mahoney et al. [2012], then, interpolates many of the arguments in Vigna [2009] beyond $\alpha = 1$ to yield important relationships between spectral graph theory and PageRank vectors.

5.4. Complex-valued teleportation parameters and a time-dependent generalization. Again, let $\mathbf{x}(\alpha)$ be the PageRank vector (in the sense of (2.2)) as a function of α for a fixed graph and teleportation vector. Mathematically, the PageRank vector is a rational function of α . This simple insight produces a host of possibilities, one of which is evaluating the derivative of the PageRank vector [Boldi et al., 2005; Golub and Greif, 2006; Gleich et al., 2007]. Another is that PageRank with complex-valued α is a reasonable mathematical generalization [Horn and Serra-Capizzano, 2007]. Let $\alpha \in \mathbb{C}$ with $|\alpha| < 1$, then $\mathbf{x}(\alpha)$ has some interesting properties and usages. In Constantine and Gleich [2010], we needed to bound $\|\mathbf{x}(\alpha)\|_1$ when α was complex. If α is real and $0 < \alpha < 1$, then $\|\mathbf{x}(\alpha)\|_1 = 1$ independent of the choice of α . However, if α is complex we have: $\|\mathbf{x}\|_1 \leq \frac{|1-\alpha|}{1-|\alpha|}$. Later, in Gleich and Rossi [2014], we found that complex values of α arise in computing closed form solutions to PageRank dynamical systems where the teleportation vector is a function of time, but the graph remains fixed. Specifically, the PageRank vector with complex teleportation arises in the steady-state time-dependent solution of

$$\mathbf{x}'(t) = (1 - \alpha)\mathbf{v}(t) - (\mathbf{I} - \alpha\mathbf{P})\mathbf{x}(t),$$

when $\mathbf{v}(t)$ oscillates between a fixed set of vectors. Thus, PageRank with complex teleportation is both an interesting mathematical problem and has practical applications in a time-dependent generalization of PageRank.

5.5. Censored node constructions. The final generalized PageRank construction we wish to discuss is, in fact, a PageRank system hiding inside a Markov chain construction with a different type of teleportation. In order to motivate the particular form of this construction, we first review an alternative derivation of the PageRank vector.

A censored node in a Markov chain is one that exhibits a virtual influence on the chain in the sense that walks proceed through it as if it were not present. Let us illustrate this idea by crafting teleportation behavior into a Markov chain in a different way and computing the PageRank vector itself by censoring that Markov chain. Suppose that we want to find the stationary distribution of a walk where, if a surfer wants to teleport, they first transition to a teleport state, and then move from the teleport state according to the teleportation distribution. The transition matrix of the Markov chain is:

$$\mathbf{P}' = \begin{bmatrix} \alpha\mathbf{P} & \mathbf{v} \\ (1 - \alpha)\mathbf{e}^T & 0 \end{bmatrix}.$$

And the stationary distribution of this Markov chain is:

$$\begin{bmatrix} \alpha \mathbf{P} & \mathbf{v} \\ (1-\alpha)\mathbf{e}^T & 0 \end{bmatrix} \begin{bmatrix} \mathbf{x}' \\ \gamma \end{bmatrix} = \begin{bmatrix} \mathbf{x}' \\ \gamma \end{bmatrix}, \quad \mathbf{e}^T \mathbf{x}' + \gamma = 1.$$

Censoring the final teleportation state amounts to modeling its influence on the stationary distribution, but leaving it with no final contribution. Put more formally, the stationary distribution of the censored chain is just \mathbf{x}' renormalized to be a probability distribution: $\mathbf{x} = \mathbf{x}' / \mathbf{e}^T \mathbf{x}'$. In other words, censoring that state models pretending that it wasn't there when determining the stationary distribution, but the transitions through it still took place; this is equivalent to the standard teleporting behavior. The vector \mathbf{x} is also the PageRank vector of $\alpha, \mathbf{P}, \mathbf{v}$, which follows from

$$\mathbf{x} = \frac{1-\alpha}{\gamma} \mathbf{x}' = \frac{1-\alpha}{\gamma} [\alpha \mathbf{P} \mathbf{x}' + \gamma \mathbf{v}] = \alpha \mathbf{P} \mathbf{x} + (1-\alpha) \mathbf{v}.$$

Tomlin [2003], Eiron et al. [2004] and Lee et al. [2007, written in 2003] were some of the first to observe this property in the context of PageRank; although censoring Markov chains goes back much further.

There is a more general class of PageRank-style methods that craft transitions akin to non-uniform teleportation through a censored node construction. Consider, for example, adding a teleportation node c that connects to all nodes of a network as in Figure 5.1. This construction gives rise to an implicit PageRank problem with $\alpha = \frac{d_{\max}}{d_{\max}+1}$ as we now show. Let

$$\mathbf{A}' = \begin{bmatrix} \mathbf{A} & \mathbf{e} \\ \mathbf{v}^T & 0 \end{bmatrix}$$

be the adjacency matrix for the modified graph, where \mathbf{v} is the teleportation destination vector. A uniform random walk on this adjacency structure has a single recurrent class, and thus, a unique stationary distribution [Berman and Plemmons, 1994, Theorem 3.23]. The stationary distribution satisfies:

$$\mathbf{P}' \mathbf{x} = \mathbf{x} \Leftrightarrow \begin{bmatrix} \mathbf{A}^T (\mathbf{D} + \mathbf{I})^{-1} & \mathbf{v} / \mathbf{e}^T \mathbf{v} \\ \mathbf{e} (\mathbf{D} + \mathbf{I})^{-1} & 0 \end{bmatrix} \begin{bmatrix} \mathbf{x}' \\ \gamma \end{bmatrix} = \begin{bmatrix} \mathbf{x}' \\ \gamma \end{bmatrix}.$$

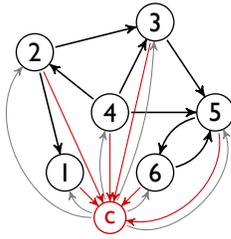
Let $\bar{\mathbf{P}}' = \mathbf{A}^T (\mathbf{D} + \mathbf{I})^{-1}$. The censored distribution $\mathbf{x} = \mathbf{x}' / \mathbf{e}^T \mathbf{x}'$ is a normalized solution of the linear system:

$$(\mathbf{I} - \bar{\mathbf{P}}') \mathbf{x} = \mathbf{v}. \quad (5.3)$$

Note that $\mathbf{c}^T = \mathbf{e}^T - \mathbf{e}^T \bar{\mathbf{P}}' > 0$, and so all columns are substochastic. This means that all of the nodes “leak probability” in a semi-formal sense. Scaling $\bar{\mathbf{P}}'$ by $\frac{1}{1-c_{\max}} > 1$ adjusts the probabilities such that there is at least one column that is stochastic. Consequently, we can write $\bar{\mathbf{P}}' = \alpha \bar{\mathbf{P}}$ where $\alpha = (1 - c_{\max})$ and $\bar{\mathbf{P}} = \frac{1}{1-c_{\max}} \bar{\mathbf{P}}'$. By substituting this form into (5.3), we have that \mathbf{x} is the normalized solution of a pseudo-PageRank problem where $\alpha = \frac{1}{1-c_{\max}}$. Assuming that \mathbf{A} is an unweighted graph, then $\alpha = \frac{d_{\max}}{d_{\max}+1}$.

This idea frequently reappears; for instance, Bini et al. [2010], Lü et al. [2011] and Schlote et al. [2012] all use it in different contexts.

In a different context, this same type of analysis shows that the Colley matrix for ranking sports teams is a diagonally perturbed, generalized pseudo-PageRank



$$\bar{\mathbf{P}} = \begin{bmatrix} 0 & 1/3 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1/4 & 0 & 0 & 0 \\ 0 & 1/3 & 0 & 1/4 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1/2 & 1/4 & 0 & 1/2 & 0 \\ 0 & 0 & 0 & 0 & 1/2 & 0 & 0 \end{bmatrix} \quad \mathbf{c} = \begin{bmatrix} 1 \\ 1/3 \\ 1/2 \\ 1/4 \\ 1/2 \\ 1/2 \\ 1/2 \end{bmatrix}$$

(a) A directed graph with a censored node c (b) The substochastic matrix and correction vector for the Markov chain construction after node c is censored.

FIG. 5.1. In this teleportation construction we add a node c to the original graph as in subfigure (a). The probability of transitioning to c , or teleporting after we censor node c , then depends on the degree of each node. A random surfer teleports from node 2 with probability $1/3$ and from node 4 with probability $1/4$. This construction yields a substochastic matrix $\bar{\mathbf{P}}$ where all the elements of the correction vector \mathbf{c} are positive. This means it's equivalent to a PageRank construction with $\alpha = 1 - \min \mathbf{c}$, or $\alpha = 3/4$ for this problem.

system [Colley, 2002; Langville and Meyer, 2012]. Let the symmetric, weighted graph G represent the network of times team i played team j . And let \mathbf{f} be a vector of the accumulated scores differences over all of those games. It could have negative entries, rendering it outside of our traditional framework, however, as we saw in Section 5.3, this is a technical detail that is avoidable. The vector of Colley scores \mathbf{r} is the solution of:

$$(\mathbf{D} + 2\mathbf{I} - \mathbf{A})\mathbf{r} = \mathbf{f}.$$

Let $\mathbf{y} = (\mathbf{D} + 2\mathbf{I})^{-1}\mathbf{r}$. Then,

$$(\mathbf{I} - \alpha\bar{\mathbf{P}})\mathbf{y} = \mathbf{f}$$

where $\alpha = \frac{d_{\max}}{d_{\max} + 2}$. This analysis establishes a formal relationship between Markov style ranking metrics [Langville and Meyer, 2012] and the least-squares style ranking metrics employed by Colley. It also enables us to use fast PageRank solvers for these Colley systems.

6. Discussion & a positive outlook on PageRank's wide usage. PageRank has gone from being used to evaluate the importance of web pages to a much broader set of applications. The method is easy to understand, is robust and reliable, and converges quickly. Most applications solve PageRank problems of only a modest size, with fewer than 100,000,000 vertices; this regime permits a much wider variety of algorithmic possibilities than those that must only work on the web.

We have avoided the discussion of PageRank algorithms entirely in this manuscript because, by and large, simple iterations suffice for fast convergence in this regime. Values of α tend to be less than 0.99, which requires fewer than 2000 iterations to converge to machine precision. Nevertheless, there is ample opportunity to accelerate PageRank computations in this regime as there are ideas that involve computing *multiple* PageRank vectors for a single task. One example is PerturbationRank [Du et al., 2008], which uses the perturbation induced in a PageRank vector by removing a node to compute a new ranking of all nodes. Thus, innovations in PageRank algorithms are still relevant, but must be made within the context of these small-scale uses.

There are also a great number of PageRank-like ideas outside of our specific canon. For instance, none of the following models fit our PageRank framework:

- BrowseRank** Liu et al. [2008] define a continuous time Markov chain to model a random surfer that *remains* on a specified node for some period of time before transitioning away. This model handles sites like Facebook, where users spend significant time interacting within a single page.
- Voting** Boldi et al. [2009a] and Boldi et al. [2011] define a voting model on a social network inspired by computing Katz or PageRank on a random network where each node picks a *single outlink*.
- SimRank** This problem is another way to use PageRank-like ideas to evaluate similarity between the nodes of a graph (like the IsoRank problem) [Jeh and Widom, 2002]. SimRank, however, involves solving a linear system on a *row sub-stochastic matrix*.
- Food webs** The food web is a network where species are linked by the feeding relationships. Allesina and Pascual [2009] point out a few modifications to PageRank to make it more appropriate. First, they use teleportation to model a constant loss of nutrients from higher-level species and reinject these nutrients through primary producers (such as bacteria). Second, they note that the flow of importance ought to be reversed so that species i points to species j if i is important for j 's survival. The result is an eigenvector computation on a fully stochastic matrix.
- Opinion dynamics** Models of opinion formation on social network posit strikingly similar dynamics to a PageRank iteration [Friedkin and Johnsen, 1990, 1999]. The essential difference is that a node's opinion is the average of its in-links, instead of propagating its value to its out-links. Like SimRank, this results in a row sub-stochastic iteration.

The details and implications of these models are fascinating, and this manuscript would double in size if we were to treat them.

In most successful studies, PageRank is used as a type of baseline measure. Its widespread success above extremely simple baselines suggests that its modified random walk is a generally useful alternative worth investigating. In this sense, it resembles a form of regularization. And this is how we feel that PageRank should be used. Note that studies must use care when determining the type of PageRank construction – weighted, reverse, Dirichlet, etc. – as this can make a large difference in the quality of the results. Consider, for instance, the use of weighted PageRank in Jiang [2009]. In their application, they wanted to model where people move, and it makes good sense that businesses would locate in places with many connections and therefore, that people would preferentially move to these same locations. Given the generality of the idea and its intuitive appeal, we anticipate continued widespread use of the PageRank idea over the next 20 years in new and exciting applications as network data continues to proliferate.

REFERENCES

- D. ALDOUS and J. A. FILL. *Reversible Markov chains and random walks on graphs*. 2002. Unfinished monograph, recompiled 2014, available at <http://www.stat.berkeley.edu/~aldous/RWG/book.html>. Cited on page 8.
- S. ALLESINA and M. PASCUAL. *Googling food webs: Can an eigenvector measure species' importance for coextinctions?* PLoS Comput Biol, 5 (9), p. e1000494, 2009. doi:10.1371/journal.pcbi.1000494. Cited on page 30.
- R. ANDERSEN, F. CHUNG, and K. LANG. *Local graph partitioning using PageRank vectors*. In *Proceedings of the 47th Annual IEEE Symposium on Foundations of Computer Science*. 2006. Cited on pages 5 and 8.

- W. N. J. ANDERSON and T. D. MORLEY. *Eigenvalues of the Laplacian of a graph*. Linear and Multilinear Algebra, 18 (2), pp. 141–145, 1985. doi:10.1080/03081088508817681. Cited on page 25.
- A. ARASU, J. NOVAK, A. TOMKINS, and J. TOMLIN. *PageRank computation and the structure of the web: Experiments and algorithms*. In *Proceedings of the 11th international conference on the World Wide Web*. 2002. Poster session. Cited on page 22.
- K. AVRACHENKOV, B. RIBEIRO, and D. TOWSLEY. *Improving random walk estimation accuracy with uniform restarts*. In *Algorithms and Models for the Web-Graph*, pp. 98–109. Springer Berlin Heidelberg, 2010. doi:10.1007/978-3-642-18009-5_10. Cited on page 8.
- L. BACKSTROM and J. LESKOVEC. *Supervised random walks: Predicting and recommending links in social networks*. In *Proceedings of the Fourth ACM International Conference on Web Search and Data Mining*, pp. 635–644. 2011. doi:10.1145/1935826.1935914. Cited on page 21.
- R. BAEZA-YATES, P. BOLDI, and C. CASTILLO. *Generalizing PageRank: Damping functions for link-based ranking algorithms*. In *Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval (SIGIR2006)*, pp. 308–315. 2006. doi:10.1145/1148170.1148225. Cited on page 24.
- B. BAHMANI, A. CHOWDHURY, and A. GOEL. *Fast incremental and personalized pagerank*. Proc. VLDB Endow., 4 (3), pp. 173–184, 2010. Cited on page 21.
- A. BALMIN, V. HRISTIDIS, and Y. PAPAOKONSTANTINOY. *ObjectRank: authority-based keyword search in databases*. In *Proceedings of the Thirtieth international conference on Very large data bases - Volume 30*, pp. 564–575. 2004. Cited on page 19.
- Z. BAR-YOSSEF and L.-T. MASHIACH. *Local approximation of PageRank and reverse PageRank*. In *CIKM '08: Proceeding of the 17th ACM conference on Information and knowledge management*, pp. 279–288. 2008. doi:10.1145/1458082.1458122. Cited on pages 7, 19, 21, and 22.
- M. BAYATI, D. F. GLEICH, A. SABERI, and Y. WANG. *Message-passing algorithms for sparse network alignment*. ACM Trans. Knowl. Discov. Data, 7 (1), pp. 3:1–3:31, 2013. doi:10.1145/2435209.2435212. Cited on page 11.
- L. BECCHETTI, C. CASTILLO, D. DONATO, R. BAEZA-YATES, and S. LEONARDI. *Link analysis for web spam detection*. ACM Trans. Web, 2 (1), pp. 1–42, 2008. doi:10.1145/1326561.1326563. Cited on page 22.
- P. BERKHIN. *A survey on PageRank computing*. Internet Mathematics, 2 (1), pp. 73–120, 2005. Cited on page 4.
- A. BERMAN and R. J. PLEMMONS. *Nonnegative Matrices in the Mathematical Sciences*, SIAM, 1994. Cited on page 28.
- D. A. BINI, G. M. D. CORSO, and F. ROMANI. *A combined approach for evaluating papers, authors and scientific journals*. Journal of Computational and Applied Mathematics, 234 (11), pp. 3104 – 3121, 2010. Numerical Linear Algebra, Internet and Large Scale Applications. doi:10.1016/j.cam.2010.02.003. Cited on page 28.
- D. M. BLEI, A. Y. NG, and M. I. JORDAN. *Latent Dirichlet allocation*. Journal of Machine Learning Research, 3, pp. 993–1022, 2003. Cited on page 15.
- V. D. BLONDEL, A. GAJARDO, M. HEYMANS, P. SENELLART, and P. V. DOOREN. *A measure of similarity between graph vertices: Applications to synonym extraction and web searching*. SIAM Review, 46 (4), pp. 647–666, 2004. doi:10.1137/S0036144502415960. Cited on pages 11 and 21.
- P. BOLDI. *TotalRank: Ranking without damping*. In *Poster Proceedings of the 14th international conference on the World Wide Web (WWW2005)*, pp. 898–899. 2005. Cited on page 24.
- P. BOLDI, F. BONCHI, C. CASTILLO, D. DONATO, A. GIONIS, and S. VIGNA. *The query-flow graph: Model and applications*. In *Proceedings of the 17th ACM Conference on Information and Knowledge Management*, pp. 609–618. 2008. doi:10.1145/1458082.1458163. Cited on page 20.
- P. BOLDI, F. BONCHI, C. CASTILLO, and S. VIGNA. *Voting in social networks*. In *CIKM '09: Proceeding of the 18th ACM conference on Information and knowledge management*, pp. 777–786. 2009a. doi:10.1145/1645953.1646052. Cited on page 30.
- . *Viscous democracy for social networks*. Commun. ACM, 54 (6), pp. 129–137, 2011. doi:10.1145/1953122.1953154. Cited on page 30.
- P. BOLDI, R. POSENATO, M. SANTINI, and S. VIGNA. *Traps and pitfalls of topic-biased PageRank*. In *WAW2006, Fourth International Workshop on Algorithms and Models for the Web-Graph*, pp. 107–116. 2007. doi:10.1007/978-3-540-78808-9_10. Cited on pages 4 and 7.
- P. BOLDI, M. SANTINI, and S. VIGNA. *PageRank as a function of the damping factor*. In *Proceedings of the 14th international conference on the World Wide Web (WWW2005)*. 2005. doi:10.1145/1060745.1060827. Cited on pages 24, 25, and 27.
- . *PageRank: Functional dependencies*. ACM Trans. Inf. Syst., 27 (4), pp. 1–23, 2009b. doi:10.1145/1629096.1629097. Cited on page 24.
- J. BOLLEN, M. A. RODRIQUEZ, and H. VAN DE SOMPEL. *Journal status*. Scientometrics, 69 (3), pp.

- 669–687, 2006. doi:10.1007/s11192-006-0176-z. Cited on page 16.
- S. BRIN and L. PAGE. *The anatomy of a large-scale hypertextual web search engine*. Comput. Netw. ISDN Syst., 30 (1-7), pp. 107–117, 1998. doi:10.1016/S0169-7552(98)00110-X. Cited on page 1.
- T. CALLAGHAN, P. J. MUCHA, and M. A. PORTER. *Random walker ranking for NCAA division I-A football*. The American Mathematical Monthly, 114 (9), pp. 761–777, 2007. Cited on page 14.
- P. CHEN, H. XIE, S. MASLOV, and S. REDNER. *Finding scientific gems with Google’s PageRank algorithm*. Journal of Informetrics, 1 (1), pp. 8–15, 2007. doi:10.1016/j.joi.2006.06.001. Cited on page 17.
- A. CHEPELIANSKIĬ. *Towards physical laws for software architecture*. arXiv, cs.SE, p. 1003.5455, 2010. Cited on page 13.
- F. CHUNG. *Laplacians and the Cheeger inequality for directed graphs*. Annals of Combinatorics, 9 (1), pp. 1–19, 2005. 10.1007/s00026-005-0237-z. doi:10.1007/s00026-005-0237-z. Cited on page 18.
- . *The heat kernel as the PageRank of a graph*. Proceedings of the National Academy of Sciences, 104 (50), pp. 19735–19740, 2007. doi:10.1073/pnas.0708838104. Cited on page 24.
- F. CHUNG, A. TSIATAS, and W. XU. *Dirichlet pagerank and trust-based ranking algorithms*. In *Algorithms and Models for the Web Graph*, pp. 103–114. Springer Berlin Heidelberg, 2011. doi:10.1007/978-3-642-21286-4_9. Cited on page 7.
- F. R. L. CHUNG. *Spectral Graph Theory*, American Mathematical Society, 1992. Cited on page 26.
- W. N. COLLEY. *Colley’s bias free college football ranking method: The Colley matrix explained*. Technical report, Princeton University, 2002. Cited on page 29.
- P. G. CONSTANTINE and D. F. GLEICH. *Random alpha PageRank*. Internet Mathematics, 6 (2), pp. 189–236, 2010. doi:10.1080/15427951.2009.10129185. Cited on pages 23, 24, 25, and 27.
- J. J. CROFTS and D. J. HIGHAM. *Googling the brain: Discovering hierarchical and asymmetric network structures, with applications in neuroscience*. Internet Mathematics, 7, pp. 233–254, 2011. doi:10.1080/15427951.2011.604284. Cited on page 12.
- M. CUTTS. *Q&A march 27, 2006*. Matt Cutts: Gadgets, Google, and SEO. Available online <http://www.matcutts.com/blog/q-a-thread-march-27-2006/>, 2006. Cited on page 22.
- . *PageRank sculpting*. Matt Cutts: Gadgets, Google, and SEO blog. Available online <http://www.matcutts.com/blog/pagerank-sculpting/>, 2009. Cited on page 22.
- T. A. DAVIS and Y. HU. *The University of Florida sparse matrix collection*. ACM Transactions on Mathematical Software, 2010. To appear. Cited on page 22.
- G. DEL CORSO, A. GULLÍ, and F. ROMANI. *Fast PageRank computation via a sparse linear system (extended abstract)*. In *Algorithms and Models for the Web-Graph*, pp. 118–130. Springer Berlin Heidelberg, 2004. doi:10.1007/978-3-540-30216-2_10. Cited on page 4.
- G. M. DEL CORSO, A. GULLÍ, and F. ROMANI. *Fast PageRank computation via a sparse linear system*. Internet Mathematics, 2 (3), pp. 251–273, 2005. Cited on page 4.
- I. S. DHILLON. *Co-clustering documents and words using bipartite spectral graph partitioning*. In *Proceedings of the seventh ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 269–274. 2001. doi:10.1145/502512.502550. Cited on page 17.
- Y. DU, J. LEUNG, and Y. SHI. *PerturbationRank: A non-monotone ranking algorithm*. Technical report, University of Michigan, 2008. Cited on page 29.
- N. EIRON, K. S. MCCURLEY, and J. A. TOMLIN. *Ranking the web frontier*. In *Proceedings of the 13th international conference on the World Wide Web (WWW2004)*, pp. 309–318. 2004. doi:10.1145/988672.988714. Cited on pages 22 and 28.
- Y.-H. EOM, P. ARAGÓN, D. LANIADO, A. KALTENBRUNNER, S. VIGNA, and D. L. SHEPELYANSKY. *Interactions of cultures and top people of Wikipedia from ranking of 24 language editions*. arXiv, cs.SI, p. 1405.7183, 2014. Cited on page 23.
- E. ESTRADA. *Characterization of 3d molecular structure*. Chemical Physics Letters, 319 (5-6), pp. 713–718, 2000. doi:10.1016/S0009-2614(00)00158-5. Cited on page 24.
- E. ESTRADA and D. J. HIGHAM. *Network properties revealed through matrix functions*. SIAM Review, 52 (4), pp. 696–714, 2010. doi:10.1137/090761070. Cited on page 24.
- E. ESTRADA, D. J. HIGHAM, and N. HATANO. *Communicability and multipartite structures in complex networks at negative absolute temperatures*. Phys. Rev. E, 78, p. 026102, 2008. doi:10.1103/PhysRevE.78.026102. Cited on page 12.
- A. FARAHAT, T. LOFARO, J. C. MILLER, G. RAE, and L. A. WARD. *Authority rankings from HITS, PageRank, and SALSA: Existence, uniqueness, and effect of initialization*. SIAM Journal on Scientific Computing, 27 (4), pp. 1181–1201, 2006. doi:10.1137/S1064827502412875. Cited on page 24.
- D. FIALA, F. ROUSSELOT, and K. JEEK. *PageRank for bibliographic networks*. Scientometrics, 76 (1), pp. 135–158, 2008. doi:10.1007/s11192-007-1908-4. Cited on page 17.

- M. FIEDLER. *Algebraic connectivity of graphs*. Czechoslovak Mathematical Journal, 23 (98), pp. 298–305, 1973. Cited on page 25.
- D. FOGARAS. *Where to start browsing the web?* In *Innovative Internet Community Systems*, pp. 65–79. Springer Berlin Heidelberg, 2003. doi:10.1007/978-3-540-39884-4_6. Cited on pages 7 and 22.
- K. M. FRAHM, A. D. CHEPELIANSKII, and D. L. SHEPELYANSKY. *PageRank of integers*. Journal of Physics A: Mathematical and Theoretical, 45 (40), p. 405101, 2012. doi:10.1088/1751-8113/45/40/405101. Cited on page 14.
- V. FRESCHI. *Protein function prediction from interaction networks using a random walk ranking algorithm*. In *Proceedings of the 7th IEEE International Conference on Bioinformatics and Bioengineering (BIBE 2007)*, pp. 42–48. 2007. doi:10.1109/BIBE.2007.4375543. Cited on page 10.
- N. E. FRIEDKIN and E. C. JOHNSEN. *Social influence and opinions*. The Journal of Mathematical Sociology, 15 (3-4), pp. 193–206, 1990. doi:10.1080/0022250X.1990.9990069. Cited on page 30.
- . *Social influence networks and opinion change*. Advances in Group Processes, 16 (1), pp. 1–29, 1999. Cited on page 30.
- E. GARFIELD. *Citation indexes for science: A new dimension in documentation through association of ideas*. Science, 122 (3159), pp. 108–111, 1955. doi:10.1126/science.122.3159.108. Cited on page 16.
- E. GARFIELD and I. H. SHER. *New factors in the evaluation of scientific literature through citation indexing*. American Documentation, 14 (3), pp. 195–201, 1963. doi:10.1002/asi.5090140304. Cited on page 16.
- D. F. GLEICH, P. GLYNN, G. H. GOLUB, and C. GREIF. *Three results on the PageRank vector: eigenstructure, sensitivity, and the derivative*. In *Web Information Retrieval and Linear Algebra Algorithms*. 2007. Cited on page 27.
- D. F. GLEICH, A. P. GRAY, C. GREIF, and T. LAU. *An inner-outer iteration for PageRank*. SIAM Journal of Scientific Computing, 32 (1), pp. 349–371, 2010. doi:10.1137/080727397. Cited on page 11.
- D. F. GLEICH and M. M. MAHONEY. *Algorithmic anti-differentiation: A case study with min-cuts, spectral, and flow*. In *Proceedings of the International Conference on Machine Learning (ICML)*. 2014. Cited on page 8.
- D. F. GLEICH and M. POLITO. *Approximating personalized PageRank with minimal use of webgraph data*. Internet Mathematics, 3 (3), pp. 257–294, 2007. doi:10.1080/15427951.2006.10129128. Cited on page 22.
- D. F. GLEICH and R. A. ROSSI. *A dynamical system for PageRank with time-dependent teleportation*. Internet Mathematics, 10 (1–2), pp. 188–217, 2014. doi:10.1080/15427951.2013.814092. Cited on page 27.
- G. GOLUB and C. GREIF. *An Arnoldi-type algorithm for computing PageRank*. BIT Numerical Mathematics, 46 (4), pp. 759–771, 2006. doi:10.1007/s10543-006-0091-y. Cited on page 27.
- M. GORI and A. PUCCI. *ItemRank: a random-walk based scoring algorithm for recommender engines*. In *IJCAI'07: Proceedings of the 20th international joint conference on Artificial intelligence*, pp. 2766–2771. 2007. Cited on page 20.
- A. Y. GOVAN, C. D. MEYER, and R. ALBRIGHT. *Generalizing Google’s PageRank to rank national football league teams*. In *SAS Global Forum 2008*. 2008. Cited on pages 14 and 15.
- Z. GYÖNGYI, H. GARCIA-MOLINA, and J. PEDERSEN. *Combating web spam with TrustRank*. In *Proceedings of the 30th International Very Large Database Conference*. 2004. Cited on pages 7 and 22.
- T. H. HAVELIWALA. *Topic-sensitive PageRank*. In *Proceedings of the 11th international conference on the World Wide Web*. 2002. Cited on page 19.
- D. J. HIGHAM. *Google PageRank as mean playing time for pinball on the reverse web*. Applied Mathematics Letters, 18 (12), pp. 1359 – 1362, 2005. doi:10.1016/j.aml.2005.02.020. Cited on page 1.
- R. A. HORN and S. SERRA-CAPIZZANO. *A general setting for the parametric Google matrix*. Internet Mathematics, 3 (4), pp. 385–411, 2007. Cited on page 27.
- A. HOTH, R. JÄSCHKE, C. SCHMITZ, and G. STUMME. *Information retrieval in folksonomies: Search and ranking*. In *Proceedings of the 3rd European Semantic Web Conference*, pp. 411–426. 2006. doi:10.1007/11762256_31. Cited on pages 16 and 19.
- B. A. HUBERMAN, P. L. T. PIROLI, J. E. PITKOW, and R. M. LUKOSE. *Strong regularities in World Wide Web surfing*. Science, 280 (5360), pp. 95–97, 1998. doi:10.1126/science.280.5360.95. Cited on page 24.
- A. JAIN and P. PANTEL. *FactRank: Random walks on a web of facts*. In *Proceedings of the 23rd*

- International Conference on Computational Linguistics*, pp. 501–509. 2010. Cited on page 18.
- A. JAVA. *Twitter social network analysis*. UMBC ebquity blog, <http://ebiquity.umbc.edu/blogger/2007/04/19/twitter-social-network-analysis/>, 2007. Cited on page 21.
- A. JAVA, P. KOLARI, T. FININ, , and T. OATES. *Modeling the spread of influence on the blogosphere*. Technical Report UMBC TR-CS-06-03, University of Maryland, Baltimore, 2006. Cited on page 21.
- G. JEH and J. WIDOM. *SimRank: a measure of structural-context similarity*. In *KDD '02: Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 538–543. 2002. doi:10.1145/775047.775126. Cited on pages 11 and 30.
- K. JEZEK, D. FIALA, and J. STEINBERGER. *Exploration and evaluation of citation networks*. In *Proceedings of the 12th International Conference on Electronic Publishing*, pp. 351–362. 2008. Cited on page 17.
- B. JIANG. *Ranking spaces for predicting human movement in an urban environment*. *Int. J. Geogr. Inf. Sci.*, 23 (7), pp. 823–837, 2009. doi:10.1080/13658810802022822. Cited on pages 8, 13, and 30.
- B. JIANG, S. ZHAO, and J. YIN. *Self-organized natural roads for predicting traffic flow: a sensitivity study*. *Journal of Statistical Mechanics: Theory and Experiment*, 2008 (07), p. P07008, 2008. doi:10.1088/1742-5468/2008/07/P07008. Cited on page 13.
- B.-B. JIANG, J.-G. WANG, J.-F. XIAO, and Y. WANG. *Gene prioritization for type 2 diabetes in tissue-specic protein interaction networks*. In *Third International Symposium on Optimization and Systems Biology*. 2009. Cited on page 10.
- Y. JING and S. BALUJA. *VisualRank: Applying pagerank to large-scale image search*. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 30 (11), pp. 1877–1890, 2008. doi:10.1109/TPAMI.2008.121. Cited on page 23.
- M. JOCKERS. *Computing and visualizing the 19th-century literary genome*. In *Digital Humanities*. 2012. Cited on page 15.
- S. KAMVAR. *Numerical Algorithms for Personalized Search in Self-organizing Information Networks*, Princeton University Press, 2010. Cited on page 10.
- S. KAMVAR, T. HAVELIWALA, C. MANNING, and G. GOLUB. *Exploiting the block structure of the web for computing PageRank*. Technical Report 2003-17, Stanford InfoLab, 2003. Cited on pages 10 and 22.
- L. KATZ. *A new status index derived from sociometric analysis*. *Psychometrika*, 18 (1), pp. 39–43, 1953. doi:10.1007/BF02289026. Cited on pages 12, 21, and 25.
- J. P. KEENER. *The Perron-Frobenius theorem and the ranking of football teams*. *SIAM Review*, 35 (1), pp. 80–93, 1993. doi:10.1137/1035004. Cited on page 14.
- D. KEMPE, J. KLEINBERG, and E. TARDOS. *Maximizing the spread of influence through a social network*. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 137–146. 2003. doi:10.1145/956750.956769. Cited on page 21.
- . *Influential nodes in a diffusion model for social networks*. In *Proceedings of the 32nd International Conference on Automata, Languages and Programming*, pp. 1127–1138. 2005. doi:10.1007/11523468_91. Cited on page 21.
- M. KIM, R. SUMBALY, and S. SHAH. *Root cause detection in a service-oriented architecture*. In *Proceedings of the ACM SIGMETRICS/international conference on Measurement and modeling of computer systems*, pp. 93–104. 2013. doi:10.1145/2465529.2465753. Cited on page 12.
- J. M. KLEINBERG. *Authoritative sources in a hyperlinked environment*. *Journal of the ACM*, 46 (5), pp. 604–632, 1999. Cited on page 21.
- T. G. KOLDA and M. J. PROCOPIO. *Generalized BadRank with graduated trust*. Technical Report SAND2009-6670, Sandia National Laboratories, 2009. Cited on page 22.
- G. KOLLIAS, E. GALLOPOULOS, and A. GRAMA. *Surfing the network for ranking by multidamping*. *Knowledge and Data Engineering, IEEE Transactions on*, PP (99), pp. 1–1, 2013. doi:10.1109/TKDE.2013.15. Cited on page 24.
- G. KOLLIAS, S. MOHAMMADI, and A. GRAMA. *Network similarity decomposition (NSD): A fast and scalable approach to network alignment*. *Knowledge and Data Engineering, IEEE Transactions on*, PP (99), p. 1, 2011. doi:10.1109/TKDE.2011.174. Cited on page 11.
- R. I. KONDOR and J. D. LAFFERTY. *Diffusion kernels on graphs and other discrete input spaces*. In *Proceedings of the Nineteenth International Conference on Machine Learning*, pp. 315–322. 2002. Cited on page 24.
- E.-M. KONTOPOULOU, M. PREDARI, T. KOSTAKIS, and E. GALLOPOULOS. *Graph and matrix metrics to analyze ergodic literature for children*. In *Proceedings of the 23rd ACM conference on Hypertext and social media*, pp. 133–142. 2012. doi:10.1145/2309996.2310018. Cited on page 15.
- D. KOSCHÜTZKI, K. A. LEHMANN, L. PEETERS, S. RICHTER, D. TENFELDE-PODEHL, , and O. ZLO-

- TOWSKI. *Centrality indices*. In *Network Analysis: Methodological Foundations*, chapter 3, pp. 16–61. Springer, 2005. doi:10.1007/b106453. Cited on page 1.
- J. KUNEGIS and A. LOMMATZSCH. *Learning spectral graph transformations for link prediction*. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pp. 561–568. 2009. doi:10.1145/1553374.1553447. Cited on page 24.
- H. KWAK, C. LEE, H. PARK, and S. MOON. *What is Twitter, a social network or a news media?* In *WWW '10: Proceedings of the 19th international conference on World wide web*, pp. 591–600. 2010. doi:10.1145/1772690.1772751. Cited on page 21.
- A. N. LANGVILLE and C. D. MEYER. *Deeper inside PageRank*. *Internet Mathematics*, 1 (3), pp. 335–380, 2004. Cited on page 7.
- . *Google's PageRank and Beyond: The Science of Search Engine Rankings*, Princeton University Press, 2006. Cited on pages 1 and 2.
- . *Who's #1? The Science of Rating and Ranking*, Princeton University Press, 2012. Cited on pages 14 and 29.
- C. P. LEE, G. H. GOLUB, and S. A. ZENIOS. *A two-stage algorithm for computing PageRank and multistage generalizations*. *Internet Mathematics*, 4 (4), pp. 299–327, 2007. Cited on page 28.
- D. LIBEN-NOWELL and J. KLEINBERG. *The link-prediction problem for social networks*. *Journal of the American Society of Information Science and Technology*, 58 (7), pp. 1019–1031, 2006. doi:10.1002/asi.20591. Cited on page 20.
- X. LIU, J. BOLLEN, M. L. NELSON, and H. VAN DE SOMPEL. *Co-authorship networks in the digital library research community*. *Inf. Process. Manage.*, 41 (6), pp. 1462–1480, 2005. doi:10.1016/j.ipm.2005.03.012. Cited on page 17.
- Y. LIU, B. GAO, T.-Y. LIU, Y. ZHANG, Z. MA, S. HE, and H. LI. *BrowseRank: letting web users vote for page importance*. In *SIGIR '08: Proceedings of the 31st annual international ACM SIGIR conference on Research and development in information retrieval*, pp. 451–458. 2008. doi:10.1145/1390334.1390412. Cited on page 30.
- L. LÜ, Y.-C. ZHANG, C. H. YEUNG, and T. ZHOU. *Leaders in social networks, the Delicious case*. *PLoS ONE*, 6 (6), p. e21202, 2011. doi:10.1371/journal.pone.0021202. Cited on page 28.
- N. MA, J. GUAN, and Y. ZHAO. *Bringing PageRank to the citation analysis*. *Information Processing & Management*, 44 (2), pp. 800–810, 2008. jce:title;Evaluating Exploratory Search Systems;jce:title;jce:title;Digital Libraries in the Context of Users Broader Activities;jce:title;. doi:http://dx.doi.org/10.1016/j.ipm.2007.06.006. Cited on page 17.
- M. W. MAHONEY, L. ORECCHIA, and N. K. VISHNOI. *A local spectral method for graphs: With applications to improving graph partitions and exploring data graphs locally*. *Journal of Machine Learning Research*, 13, p. 23392365, 2012. Cited on pages 25 and 27.
- X. MENG. *Computing BookRank via social cataloging*. 2009. <http://cads.stanford.edu/projects/presentations/2009visit/bookrank.pdf>. Cited on page 16.
- C. D. MEYER. *Matrix analysis and applied linear algebra*, Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2000. Cited on page 25.
- J. C. MILLER, G. RAE, F. SCHAEFER, L. A. WARD, T. LOFARO, and A. FARAHAT. *Modifications of Kleinberg's HITS algorithm using matrix exponentiation and web log records*. In *SIGIR '01: Proceedings of the 24th annual international ACM SIGIR conference on Research and development in information retrieval*, pp. 444–445. 2001. doi:10.1145/383952.384086. Cited on page 24.
- B. L. MOONEY, L. R. CORRALES, and A. E. CLARK. *Molecularnetworks: An integrated graph theoretic and data mining tool to explore solvent organization in molecular simulation*. *Journal of Computational Chemistry*, 33 (8), pp. 853–860, 2012. doi:10.1002/jcc.22917. Cited on page 10.
- J. L. MORRISON, R. BREITLING, D. J. HIGHAM, and D. R. GILBERT. *GeneRank: using search engine technology for the analysis of microarray experiments*. *BMC Bioinformatics*, 6 (1), p. 233, 2005. doi:10.1186/1471-2105-6-233. Cited on pages 1 and 10.
- M. A. NAJORK, H. ZARAGOZA, and M. J. TAYLOR. *HITS on the web: How does it compare?* In *Proceedings of the 30th annual international ACM SIGIR conference on Research and Development in information retrieval (SIGIR2007)*, pp. 471–478. 2007. doi:10.1145/1277741.1277823. Cited on page 21.
- Z. NIE, Y. ZHANG, J.-R. WEN, and W.-Y. MA. *Object-level ranking: Bringing order to web objects*. In *Proceedings of the 14th International Conference on World Wide Web*, pp. 567–574. 2005. doi:10.1145/1060745.1060828. Cited on page 18.
- G. OSIPENKO. *Dynamical systems, graphs, and algorithms*, Springer, 2007. doi:10.1007/3-540-35593-6. Cited on page 14.
- L. PAGE, S. BRIN, R. MOTWANI, and T. WINOGRAD. *The PageRank citation ranking: Bringing order to the web*. Technical Report 1999-66, Stanford University, 1999. Cited on pages 1 and 3.

- J.-Y. PAN, H.-J. YANG, C. FALOUTSOS, and P. DUYGULU. *Automatic multimedia cross-modal correlation discovery*. In *KDD '04: Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 653–658. 2004. doi:10.1145/1014052.1014135. Cited on pages 1 and 18.
- G. PINSKI and F. NARIN. *Citation influence for journal aggregates of scientific publications: Theory, with application to the literature of physics*. *Information Processing & Management*, 12 (5), pp. 297–312, 1976. doi:10.1016/0306-4573(76)90048-0. Cited on pages 16 and 25.
- A. POTHEN, H. D. SIMON, and K.-P. LIOU. *Partitioning sparse matrices with eigenvectors of graphs*. *SIAM J. Matrix Anal. Appl.*, 11, pp. 430–452, 1990. doi:10.1137/0611030. Cited on pages 25 and 26.
- F. RADICCHI. *Who is the best player ever? a complex network analysis of the history of professional tennis*. *PLoS ONE*, 6 (2), p. e17249, 2011. doi:10.1371/journal.pone.0017249. Cited on page 15.
- S. E. SCHAEFFER. *Graph clustering*. *Computer Science Review*, 1 (1), pp. 27–64, 2007. doi:10.1016/j.cosrev.2007.05.001. Cited on page 26.
- A. SCHLOTE, E. CRISOSTOMI, S. KIRKLAND, and R. SHORTEN. *Traffic modelling framework for electric vehicles*. *International Journal of Control*, 85 (7), pp. 880–897, 2012. doi:10.1080/00207179.2012.668716. Cited on pages 13 and 28.
- S. SERRA-CAPIZZANO. *Jordan canonical form of the Google matrix: A potential contribution to the PageRank computation*. *SIAM Journal on Matrix Analysis and Applications*, 27 (2), pp. 305–312, 2005. doi:10.1137/S0895479804441407. Cited on page 24.
- D. L. SHEPELYANSKY and O. V. ZHIROV. *Google matrix, dynamical attractors, and ulam networks*. *Phys. Rev. E*, 81 (3), p. 036213, 2010. doi:10.1103/PhysRevE.81.036213. Cited on page 14.
- R. SINGH, J. XU, and B. BERGER. *Pairwise global alignment of protein interaction networks by matching neighborhood topology*. In *Proceedings of the 11th Annual International Conference on Research in Computational Molecular Biology (RECOMB)*, pp. 16–31. 2007. doi:10.1007/978-3-540-71681-5_2. Cited on page 11.
- M. SOBEK. *PR0 - Google's PageRank 0 penalty*. Online., 2003. <http://pr.efactory.de/e-pr0.shtml>. Accessed 2013-09-19. Cited on page 22.
- Y. SONG, D. ZHOU, and L.-w. HE. *Query suggestion by constructing term-transition graphs*. In *Proceedings of the Fifth ACM International Conference on Web Search and Data Mining*, pp. 353–362. 2012. doi:10.1145/2124295.2124339. Cited on page 20.
- O. SPORNS. *Networks of the Brain*, The MIT Press, 2011. Cited on page 11.
- W. J. STEWART. *Introduction to the Numerical Solution of Markov Chains*, Princeton University Press, 1994. Cited on page 8.
- J. J. SYLVESTER. *Chemistry and algebra*. *Nature*, 17, p. 284, 1878. doi:10.1038/017284a0. Cited on page 10.
- J. A. TOMLIN. *A new paradigm for ranking pages on the world wide web*. In *Proceedings of the 12th International Conference on World Wide Web*, pp. 350–355. 2003. doi:10.1145/775152.775202. Cited on page 28.
- H. TONG, C. FALOUTSOS, and J.-Y. PAN. *Fast random walk with restart and its applications*. In *ICDM '06: Proceedings of the Sixth International Conference on Data Mining*, pp. 613–622. 2006. doi:10.1109/ICDM.2006.70. Cited on page 19.
- S. VIGNA. *TruRank: Taking PageRank to the limit*. In *WWW '05: Special interest tracks and posters of the 14th international conference on World Wide Web*, pp. 976–977. 2005. doi:10.1145/1062745.1062826. Cited on page 25.
- . *Spectral ranking*. arXiv, cs.IR, p. 0912.0238, 2009. Cited on pages 1, 21, 25, and 27.
- K. VOEVODSKI, S.-H. TENG, and Y. XIA. *Spectral affinity in protein networks*. *BMC Systems Biology*, 3 (1), p. 112, 2009. doi:10.1186/1752-0509-3-112. Cited on page 11.
- D. WALKER, H. XIE, K.-K. YAN, and S. MASLOV. *Ranking scientific publications using a model of network traffic*. *Journal of Statistical Mechanics: Theory and Experiment*, 2007 (06), p. P06010, 2007. doi:10.1088/1742-5468/2007/06/P06010. Cited on page 16.
- W. Y. WANG, K. MAZAITIS, and W. W. COHEN. *Programming with personalized PageRank: A locally groundable first-order probabilistic logic*. In *Proceedings of the 22Nd ACM International Conference on Conference on Information and Knowledge Management*, pp. 2129–2138. 2013. doi:10.1145/2505515.2505573. Cited on page 20.
- J. WENG, E.-P. LIM, J. JIANG, and Q. HE. *TwitterRank: finding topic-sensitive influential twitterers*. In *WSDM '10: Proceedings of the third ACM international conference on Web search and data mining*, pp. 261–270. 2010. doi:10.1145/1718487.1718520. Cited on page 21.
- J. D. WEST, T. C. BERGSTROM, and C. T. BERGSTROM. *The Eigenfactor metrics: A network approach to assessing scholarly journals*. *College & Research Libraries*, 71 (3), pp. 236–244, 2010. arXiv:<http://cr1.acrl.org/content/71/3/236.full.pdf+html>. Cited on page 16.

- C. WINTER, G. KRISTIANSEN, S. KERSTING, J. ROY, D. AUST, T. KNÖSEL, P. RMMELE, B. JAHNKE, V. HENTRICH, F. RÜCKERT, M. NIEDERGETHMANN, W. WEICHERT, M. BAHRA, H. J. SCHLITT, U. SETTMACHER, H. FRIESS, M. BÜCHLER, H.-D. SAEGER, M. SCHROEDER, C. PILARSKY, and R. GRITZMANN. *Google goes cancer: Improving outcome prediction for cancer patients by network-based ranking of marker genes*. PLoS Comput Biol, 8 (5), p. e1002511, 2012. doi:10.1371/journal.pcbi.1002511. Cited on page 10.
- A. D. WISSNER-GROSS. *Preparation of topical reading lists from the link structure of Wikipedia*. In *ICALT '06: Proceedings of the Sixth IEEE International Conference on Advanced Learning Technologies*, pp. 825–829. 2006. Cited on page 23.
- W. XING and A. GHORBANI. *Weighted PageRank algorithm*. In *Communication Networks and Services Research, 2004. Proceedings. Second Annual Conference on*, pp. 305–314. 2004. doi:10.1109/DNSR.2004.1344743. Cited on page 8.
- H. YANG, I. KING, and M. R. LYU. *DiffusionRank: a possible penicillin for web spamming*. In *Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval*, pp. 431–438. 2007. doi:10.1145/1277741.1277815. Cited on page 24.
- A. O. ZHIROV, O. V. ZHIROV, and D. L. SHEPELYANSKY. *Two-dimensional ranking of Wikipedia articles*. The European Physical Journal B, 77 (4), pp. 523–531, 2010. doi:10.1140/epjb/e2010-10500-7. Cited on pages 13 and 23.
- D. ZHOU, O. BOUSQUET, T. N. LAL, J. WESTON, and B. SCHÖLKOPF. *Learning with local and global consistency*. In *NIPS*. 2003. Cited on page 18.
- D. ZHOU, J. HUANG, and B. SCHÖLKOPF. *Learning from labeled and unlabeled data on a directed graph*. In *ICML '05: Proceedings of the 22nd International Conference on Machine Learning*, pp. 1036–1043. 2005. doi:10.1145/1102351.1102482. Cited on page 18.
- X.-N. ZUO, R. EHMKE, M. MENNES, D. IMPERATI, F. X. CASTELLANOS, O. SPORNS, and M. P. MILHAM. *Network centrality in the human functional connectome*. Cerebral Cortex, 2011. doi:10.1093/cercor/bhr269. Cited on page 11.

Acknowledgments. We acknowledge the following individuals for their discussions about this manuscript: Sebastiano Vigna, Amy Langville, Michael Saunders, Chen Greif, Des Higham, and Stratis Gallopoulos, as well as Kyle Kloster for carefully reading several early drafts. This work was supported in part by NSF CAREER award CCF-1149756.