

# A Singular Perturbation Approach for Choosing the PageRank Damping Factor

Konstantin Avrachenkov, Nelly Litvak, and Kim Son Pham

**Abstract.** We study the PageRank mass of principal components in a bow-tie web graph as a function of the damping factor  $c$ . It is known that the web graph can be divided into three principal components: SCC, IN, and OUT. The giant strongly connected component (SCC) contains a large group of pages having a hyperlink path connecting them. The pages in the IN (OUT) component have a path to (from) the SCC, but not back. Using a singular perturbation approach, we show that the PageRank share of the IN and SCC components remains high even for very large values of the damping factor, in spite of the fact that it drops to zero when  $c$  tends to one. However, a detailed study of the OUT component reveals the presence of “dead ends” (small groups of pages linking only to each other) that receive an unfairly high ranking when  $c$  is close to 1. We argue that this problem can be mitigated by choosing  $c$  as small as  $\frac{1}{2}$ .

## 1. Introduction

The link-based ranking schemes such as PageRank [Page et al. 98], HITS [Kleinberg 99], and SALSA [Lempel and Moran 00] have been successfully used in search engines to provide adequate importance measures for web pages. In the present work we restrict ourselves to the analysis of the PageRank criterion and use the following definition of PageRank from [Langville and Meyer 03]. Denote by  $n$  the total number of pages on the web and define the  $n \times n$  hyperlink matrix

$W$  as follows:

$$w_{ij} = \begin{cases} 1/d_i, & \text{if page } i \text{ links to } j, \\ 1/n, & \text{if page } i \text{ is dangling,} \\ 0, & \text{otherwise,} \end{cases} \quad (1.1)$$

for  $i, j = 1, \dots, n$ , where  $d_i$  is the number of outgoing links from page  $i$ . A page is called *dangling* if it does not have outgoing links. The PageRank is defined as a stationary distribution of a Markov chain whose state space is the set of all web pages, and the transition matrix is

$$G = cW + (1 - c)(1/n)\mathbf{1}\mathbf{1}^T. \quad (1.2)$$

Here and throughout this paper we use the symbol  $\mathbf{1}$  for a column vector of ones having by default an appropriate dimension. In (1.2),  $\mathbf{1}\mathbf{1}^T$  is a matrix all of whose entries are equal to one, and  $c \in (0, 1)$  is the parameter known as a *damping factor*. Let  $\pi$  be the PageRank vector. Then by definition,  $\pi G = \pi$ , and  $\|\pi\| = \pi\mathbf{1} = 1$ , where we write  $\|\mathbf{x}\|$  for the  $L_1$ -norm of the vector  $\mathbf{x}$ .

The damping factor  $c$  is a crucial parameter in the PageRank definition. It regulates the level of the uniform noise introduced to the system. Based on publicly available information, Google originally used  $c = 0.85$ , which appears to be a reasonable compromise between the true reflection of the web structure and numerical efficiency (see [Langville and Meyer 06] for more details). However, it was mentioned in [Boldi et al. 05] that a value of  $c$  too close to one results in distorted ranking of important pages. This phenomenon was also independently observed in [Avrachenkov and Litvak 06]. Moreover, with smaller  $c$ , PageRank is more robust, that is, one can bound the influence of outgoing links of a page (or a small group of pages) on the PageRank of other groups [Bianchini et al. 05] and on its own PageRank [Avrachenkov and Litvak 06].

In this paper we explore the idea of relating the choice of  $c$  to specific properties of the web structure. The authors of [Broder et al. 00, Kumar et al. 00] have shown that the web graph can be divided into three principal components. The giant strongly connected component (SCC) contains a large group of pages having a hyperlink path connecting them. The pages in the IN (OUT) component have a path to (from) the SCC, but not back. Furthermore, the SCC component is larger than the second-largest strongly connected component by several orders of magnitude.

In Section 3 we consider a Markov walk governed by the hyperlink matrix  $W$  and explicitly describe the limiting behavior of the PageRank vector as  $c \rightarrow 1$  with the help of singular perturbation theory [Avrachenkov 99, Korolyuk and Turbin 93, Pervozvanskii and Gaitsgori 88, Yin and Zhang 05]. We experimentally study the OUT component in more detail to discover a so-called pure OUT

component (the OUT component without dangling nodes and their predecessors) and show that pure OUT contains a number of small sub-SCCs, or dead ends, that absorb the total PageRank mass when  $c = 1$ . In Section 4 we analyze the shape of the PageRank of IN+SCC as a function of  $c$ . The dangling nodes turn out to play an unexpectedly important role in the qualitative behavior of this function.

Our analytical and experimental results suggest that the PageRank mass of IN+SCC is sustained on a high level for quite large values of  $c$ , in spite of the fact that it drops to zero as  $c \rightarrow 1$ . Furthermore, the PageRank mass of IN+SCC has a unique maximum. Then in Section 5 we show that the total PageRank mass of the pure OUT component increases with  $c$ . We argue that  $c = 0.85$  results in an inadequately high ranking for pure OUT pages, and we present an argument based on singular perturbation theory for choosing  $c$  as small as  $\frac{1}{2}$ . We confirm our theoretical argument by experiments with log files.

We would like to mention that the value  $c = \frac{1}{2}$  was also used in [Chen et al. 06] to find gems in scientific citations. This choice was justified intuitively by the observation that researchers may check references in cited papers, but on average they hardly go deeper than two levels. Nowadays, when search engines work really fast, this argument also applies to web search. Indeed, it is easier for the user to refine a query and receive a more relevant page in a fraction of a second than to look for this page by clicking on hyperlinks. Therefore, we may assume that a surfer searching for a page does not go deeper on average than two clicks.

## 2. Data Sets

We have collected two web graphs, which we denote by INRIA and FrMath-Info. The web graph INRIA was taken from the site of INRIA,<sup>1</sup> the French National Institute for Research in Computer Science and Control. The seed for the INRIA collection was the web page [www.inria.fr](http://www.inria.fr). It is a typical large web site with around 300,000 pages and two million hyperlinks. We have collected all pages belonging to INRIA. The web graph FrMathInfo was crawled with the initial seeds of 50 French mathematics and informatics laboratories, taken from Google Directory. The crawl was executed by a breadth-first search of depth 6. The FrMathInfo web graph contains around 700,000 pages and eight million hyperlinks. Because the web seems to have a fractal structure [Dill et al. 02], we expect our data sets to be sufficiently representative.

---

<sup>1</sup>*Institut National de Recherche en Informatique et en Automatique.*

#	INRIA	FrMathInfo
total nodes	318585	764119
nodes in SCC	154142	333175
nodes in IN	0	0
nodes in OUT	164443	430944
nodes in ESCC	300682	760016
nodes in Pure OUT	17903	4103
SCCs in OUT	1148	1382
SCCs in Pure Out	631	379

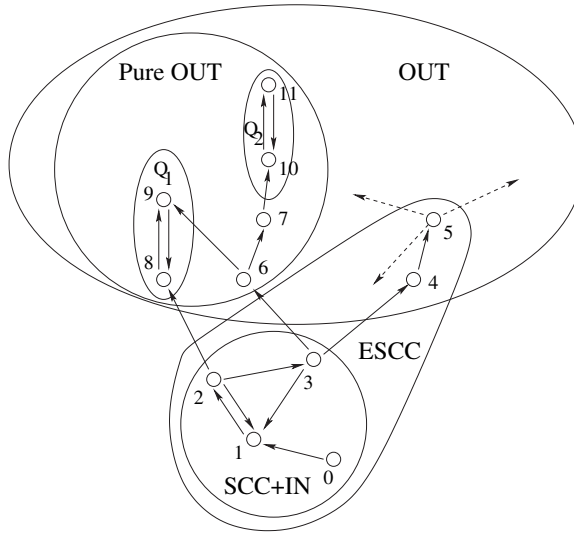
**Table 1.** Component sizes in the INRIA and FrMathInfo data sets.

The link structure of the two web graphs is stored in an Oracle database. We are able to store the adjacency lists in RAM to speed up the computation of PageRank and other quantities of interest. This enables us to make more iterations, which is extremely important when the damping factor  $c$  is close to one. Our PageRank computation program consumes about one hour to make 500 iterations for the FrMathInfo data set and about half an hour for the INRIA data set for the same number of iterations. Our algorithms for discovering the structures of the web graph are based on breadth-first search and depth-first search methods, which are linear in the sum of the number of nodes and links.

### 3. The Structure of the Hyperlink Transition Matrix

Let us refine the bow-tie structure of the web graph [Broder et al. 00, Kumar et al. 00]. We recall that the transition matrix  $W$  induces artificial links to all pages from dangling nodes. Obviously, the graph with many artificial links has a much higher connectivity than the original web graph. In particular, if the random walk can move from a dangling node to an arbitrary node with the uniform distribution, then the giant SCC component increases further in size. We refer to this new strongly connected component as the extended strongly connected component (ESCC). Due to the artificial links from the dangling nodes, the SCC component and IN component are now interconnected and are parts of the ESCC. Furthermore, if there are dangling nodes in the OUT component, then these nodes together with all their predecessors become a part of the ESCC.

In the mini-example in Figure 1, node 0 represents the IN component, nodes 1 to 3 form the SCC component, and the rest of the nodes, 4 to 11, are in the OUT component. Node 5 is a dangling node, and thus artificial links go from the dangling node 5 to all other nodes. After addition of the artificial links, all nodes from 0 to 5 form the ESCC.



**Figure 1.** Example of a graph.

The part of the OUT component without dangling nodes and their predecessors forms a block that we refer to as a pure OUT component. In Table 1 the pure OUT component consists of nodes 6 to 11. Typically, the pure OUT component is much smaller than the extended SCC.

The sizes of all components for our two data sets are given in Table 1. Here the size of the IN components is zero, because in the web crawl we used breadth-first search, and we started from important pages in the giant SCC. For the purposes of the present research this makes no difference, since we always consider IN and SCC together.

Let us now analyze the structure of the pure OUT component in more detail. It turns out that inside pure OUT there are many disjoint strongly connected components. We refer to these sub-SCCs as “dead ends,” since once the random walk induced by transition matrix  $W$  enters such a component, it will not be able to leave it. In Figure 1 there are two dead-end components,  $\{8, 9\}$  and  $\{10, 11\}$ . We have observed that in our two data sets the majority of dead ends are of size 2 or 3.

Let us now characterize the new refined structure of the web graph in terms of the ergodic structure of the Markov chain induced by the matrix  $W$ . First, we note that all states in the dead ends are *recurrent*, that is, the Markov chain started from any of these states always returns to it. In contrast, all the states

from ESCC are *transient*, that is, with probability 1, the Markov chain induced by  $W$  eventually leaves this set of states and never returns. The stationary probability of all these states is zero. We note that the pure OUT component also contains transient states that eventually bring the random walk into one of the dead ends. For simplicity, we add these states to the giant transient ESCC component.

By appropriate renumbering of the states, we can now refine the matrix  $W$  by subdividing all states into one giant transient block and a number of small recurrent blocks as follows:

$$W = \begin{bmatrix} Q_1 & & 0 & 0 \\ & \ddots & & \\ 0 & & Q_m & 0 \\ R_1 & \cdots & R_m & T \end{bmatrix} \begin{array}{l} \text{dead end (recurrent)} \\ \cdots \\ \text{dead end (recurrent)} \\ \text{ESCC} + [\text{transient states in pure OUT}] \text{ (transient)} \end{array}$$

Here for  $i = 1, \dots, m$ , a block  $Q_i$  corresponds to transitions inside the  $i$ th recurrent block, and a block  $R_i$  contains transition probabilities from transient states to the  $i$ th recurrent block. Block  $T$  corresponds to transitions between the transient states. For instance, in the example of the graph from Figure 1, nodes 8 and 9 correspond to block  $Q_1$ , nodes 10 and 11 correspond to block  $Q_2$ , and all other nodes belong to block  $T$ . Let us denote by  $\bar{\pi}_{\text{OUT},i}$  the stationary distribution corresponding to block  $Q_i$ .

We would like to emphasize that the recurrent blocks here are really small, constituting altogether about 5% for INRIA and about 0.5% for FrMathInfo. We believe that for larger data sets, this percentage will be even less. By far the most important portion of the pages is contained in the ESCC, which constitutes the major part of the giant transient block. However, if the random walk is governed by transition matrix  $W$ , it is absorbed with probability 1 into one of the recurrent blocks.

The use of the Google transition matrix  $G$  with  $c < 1$  (1.2) instead of  $W$  ensures that all the pages are recurrent states with positive stationary probabilities. However, if  $c = 1$ , the majority of pages turn into transient states with stationary probability zero. Hence, the random walk governed by the Google transition matrix  $G$  is in fact a singularly perturbed Markov chain. Informally, by singular perturbation we mean relatively small changes in elements of the matrix that lead to altered connectivity and stationary behavior of the chain. Using the results of singular perturbation theory (see, e.g., [Avrachenkov 99, Koryuk and Turbin 93, Pervozvanskii and Gaitsgori 88, Yin and Zhang 05]), in the next proposition we characterize explicitly the limiting PageRank vector as  $c \rightarrow 1$ .

**Proposition 3.1.** *Let  $\bar{\pi}_{\text{OUT},i}$  be a stationary distribution of the Markov chain governed by  $Q_i$ ,  $i = 1, \dots, m$ . Then we have*

$$\lim_{c \rightarrow 1} \pi(c) = [\pi_{\text{OUT},1} \ \cdots \ \pi_{\text{OUT},m} \ \mathbf{0}],$$

where

$$\pi_{\text{OUT},i} = \left( \frac{\# \text{ nodes in block } Q_i}{n} + \frac{1}{n} \mathbf{1}^T [I - T]^{-1} R_i \mathbf{1} \right) \bar{\pi}_{\text{OUT},i} \quad (3.1)$$

for  $i = 1, \dots, m$ ,  $I$  is the identity matrix, and  $\mathbf{0}$  is a row vector of zeros that correspond to stationary probabilities of the states in the transient block.

**Proof.** First, we note that if we make a change of variables  $\varepsilon = 1 - c$ , the Google matrix becomes a transition matrix of a singularly perturbed Markov chain as in Lemma 6.1 (see the appendix, Section 6) with  $A = W$  and  $C = \frac{1}{n} \mathbf{1} \mathbf{1}^T - W$ . Specifically,  $A_i = Q_i$ ,  $L_i = R_i$ ,  $E = T$ , and  $\mu_i = \bar{\pi}_{\text{OUT},i}$ . Next, define the aggregated generator matrix  $D$  as follows:

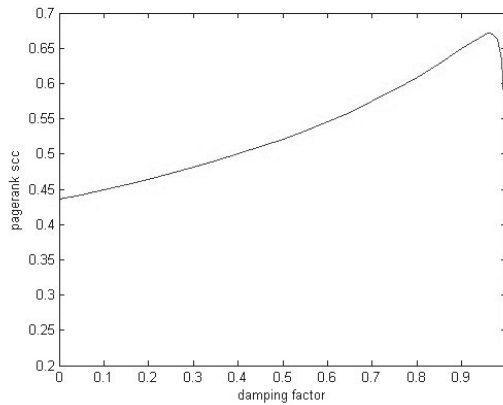
$$D = \frac{1}{n} \mathbf{1} \mathbf{1}^T B - I = \frac{1}{n} \mathbf{1} [n_1 + \mathbf{1} [I - T]^{-1} R_1 \mathbf{1}, \dots, n_m + \mathbf{1} [I - T]^{-1} R_m \mathbf{1}] - I. \quad (3.2)$$

Using the definition of  $C$  together with identities  $\bar{\pi}_{\text{OUT},i} (1/n) \mathbf{1} \mathbf{1}^T = (1/n) \mathbf{1} \mathbf{1}^T$  and  $\bar{\pi}_{\text{OUT},i} Q_i = \bar{\pi}_{\text{OUT},i}$ , it is easy to verify that the matrix  $D$  in (3.2) has been computed in exactly the same way as the matrix  $D$  in Lemma 6.1. Furthermore, since the aggregated transition matrix  $D + I$  has identical rows, its stationary distribution  $\nu$  is simply equal to each of these rows. Thus, invoking Lemma 6.1, we obtain (3.1).  $\square$

The second term inside the parentheses in formula (3.1) corresponds to the PageRank mass received by a dead end from the extended SCC. If  $c$  is close to one, then this contribution can outweigh by far the fair share of the PageRank, whereas the PageRank mass of the giant transient block decreases to zero. How large is the neighborhood of one where the ranking is skewed toward the pure OUT? Is the value  $c = 0.85$  already too large? We will address these questions in the remainder of the paper. In the next section we analyze the PageRank mass IN+SCC component, which is an important part of the transient block.

## 4. PageRank Mass of IN+SCC

In Figure 2 we depict the PageRank mass of the giant component IN+SCC for FrMathInfo as a function of the damping factor.



**Figure 2.** The PageRank mass of IN+SCC as a function of  $c$ .

Here we see a typical behavior also observed for several pages in the mini-web from [Boldi et al. 05]: the PageRank first grows with  $c$  and then decreases to zero. In our case, the PageRank mass of IN+SCC drops drastically starting from some value  $c$  close to one. Our goal now is to explain this behavior. Clearly, since IN+SCC is a part of the transient block, we do expect that the corresponding PageRank mass drops to zero when  $c$  goes to one. Thus, the two phenomena that remain to be justified are the growth of the PageRank mass when  $c$  is not too large, and the abrupt drop to zero after reaching a (unique) extreme point.

The plan of the analysis in this section is as follows. First, we write the expression for  $\|\pi_{\text{IN+SCC}}\|$ , the PageRank mass of IN+SCC, as a function of  $c$ . Then we consider the derivative of  $\|\pi_{\text{IN+SCC}}(c)\|$  at  $c = 0$  and prove that surprisingly, this derivative is always positive in a graph with a sufficiently large fraction of dangling nodes. This explains the fact that  $\|\pi_{\text{IN+SCC}}(c)\|$  is initially increasing. Further, we use singular perturbation theory to show that the derivative of  $\|\pi_{\text{IN+SCC}}(c)\|$  at  $c = 1$  is a large negative number, and that  $\|\pi_{\text{IN+SCC}}(c)\|$  can have only one extreme point in  $(0, 1)$ .

We base our analysis on the model in which the web graph sample is subdivided into three subsets of nodes: IN+SCC, OUT, and the set of dangling nodes DN. We assume that all links to dangling nodes come from IN+SCC. This simplifies the derivation but does not alter our conclusions. Then the web hyperlink matrix  $W$  in (1.1) can be written in the form

$$W = \begin{bmatrix} Q & 0 & 0 \\ R & P & S \\ \frac{1}{n}\mathbf{1}\mathbf{1}^T & \frac{1}{n}\mathbf{1}\mathbf{1}^T & \frac{1}{n}\mathbf{1}\mathbf{1}^T \end{bmatrix} \begin{array}{l} \text{OUT} \\ \text{IN+SCC} \\ \text{DN} \end{array}$$



where the block  $Q$  corresponds to the hyperlinks inside the OUT component, the block  $R$  corresponds to the hyperlinks from IN+SCC to OUT, the block  $P$  corresponds to the hyperlinks inside the IN+SCC component, and the block  $S$  corresponds to the hyperlinks from SCC to dangling nodes. In the above,  $n$  is the total number of pages in the web graph sample, and the blocks  $\mathbf{1}\mathbf{1}^T$  are the matrices of ones adjusted to appropriate dimensions.

Let us derive the expression for the PageRank mass of IN+SCC. Dividing the PageRank vector into segments corresponding to the blocks OUT, IN+SCC, and DN, namely  $\pi = [\pi_{\text{OUT}}, \pi_{\text{IN+SCC}}, \pi_{\text{DN}}]$ , we can rewrite the well-known formula (see, e.g., [Moler and Moler 03])

$$\pi = \frac{1-c}{n} \mathbf{1}^T [I - cW]^{-1} \quad (4.1)$$

as a system of three linear equations:

$$\pi_{\text{OUT}}[I - cQ] - \pi_{\text{IN+SCC}}cR - \frac{c}{n}\pi_{\text{DN}}\mathbf{1}\mathbf{1}^T = \frac{1-c}{n}\mathbf{1}^T, \quad (4.2)$$

$$\pi_{\text{IN+SCC}}[I - cP] - \frac{c}{n}\pi_{\text{DN}}\mathbf{1}\mathbf{1}^T = \frac{1-c}{n}\mathbf{1}^T, \quad (4.3)$$

$$-\pi_{\text{IN+SCC}}cS + \pi_{\text{DN}} - \frac{c}{n}\pi_{\text{DN}}\mathbf{1}\mathbf{1}^T = \frac{1-c}{n}\mathbf{1}^T. \quad (4.4)$$

Now we would like to solve (4.2)–(4.4) for  $\pi_{\text{IN+SCC}}$ . To this end, we first observe that if  $\pi_{\text{IN+SCC}}$  and  $\pi_{\text{DN}}\mathbf{1}$  are known, then from (4.2) it is straightforward to obtain  $\pi_{\text{OUT}}$ :

$$\pi_{\text{OUT}} = \pi_{\text{IN+SCC}}cR[I - cQ]^{-1} + \left( \frac{1-c}{n} + \pi_{\text{DN}}\mathbf{1}\frac{c}{n} \right) \mathbf{1}^T [I - cQ]^{-1}.$$

Therefore, let us solve equations (4.3) and (4.4). We sum the elements of the vector equation (4.4), which corresponds to the postmultiplication of equation (4.4) by the vector  $\mathbf{1}$ :

$$-\pi_{\text{IN+SCC}}cS\mathbf{1} + \pi_{\text{DN}}\mathbf{1} - \frac{c}{n}\pi_{\text{DN}}\mathbf{1}\mathbf{1}^T\mathbf{1} = \frac{1-c}{n}\mathbf{1}^T\mathbf{1}.$$

Now denote by  $n_{\text{IN}}$ ,  $n_{\text{OUT}}$ ,  $n_{\text{IN+SCC}}$ , and  $n_{\text{DN}}$  the number of pages in the IN component, the OUT component, the SCC component, and the number of dangling nodes. Since  $\mathbf{1}^T\mathbf{1} = n_{\text{DN}}$ , we have

$$\pi_{\text{DN}}\mathbf{1} = \frac{n}{n - cn_{\text{DN}}} \left( \pi_{\text{IN+SCC}}cS\mathbf{1} + \frac{1-c}{n}n_{\text{DN}} \right).$$

Substituting the above expression for  $\pi_{\text{DN}}\mathbf{1}$  into (4.3), we get

$$\pi_{\text{IN+SCC}} \left[ I - cP - \frac{c^2}{n - cn_{\text{DN}}}S\mathbf{1}\mathbf{1}^T \right] = \frac{c}{n - cn_{\text{DN}}} \frac{1-c}{n} n_{\text{DN}} \mathbf{1}^T + \frac{1-c}{n} \mathbf{1}^T.$$

Denote by  $\alpha = (n_{\text{IN}} + n_{\text{IN+SCC}})/n$  and  $\beta = n_{\text{DN}}/n$  the fractions of nodes in IN+SCC and DN, respectively, and let  $\mathbf{u}_{\text{IN+SCC}} = (n_{\text{IN}} + n_{\text{IN+SCC}})^{-1} \mathbf{1}^T$  be a uniform probability row vector of dimension  $n_{\text{IN}} + n_{\text{IN+SCC}}$ . Then from the last equation we directly obtain

$$\pi_{\text{IN+SCC}}(c) = \frac{(1-c)\alpha}{1-c\beta} \mathbf{u}_{\text{IN+SCC}} \left[ I - cP - \frac{c^2\alpha}{1-c\beta} S \mathbf{1} \mathbf{u}_{\text{IN+SCC}} \right]^{-1}. \quad (4.5)$$

Equation (4.5) gives the desired expression for the PageRank mass of IN+SCC as a function of  $c$ , and we can analyze the behavior of this function by looking at its derivatives. Define

$$k(c) = \frac{(1-c)\alpha}{1-c\beta} \quad \text{and} \quad U(c) = P + \frac{c\alpha}{1-c\beta} S \mathbf{1} \mathbf{u}_{\text{IN+SCC}}. \quad (4.6)$$

Then the derivative of  $\pi_{\text{IN+SCC}}(c)$  with respect to  $c$  is given by

$$\pi'_{\text{IN+SCC}}(c) = \mathbf{u}_{\text{IN+SCC}} \{k'(c)I + k(c)[I - cU(c)]^{-1}(cU(c))'\} [I - cU(c)]^{-1}, \quad (4.7)$$

where from (4.6) after simple calculations we get

$$\begin{aligned} k'(c) &= -(1-\beta)\alpha/(1-c\beta)^2, \\ (cU(c))' &= U(c) + c\alpha(1-c\beta)^{-2} S \mathbf{1} \mathbf{u}_{\text{IN+SCC}}. \end{aligned}$$

Now we are ready to explain the fact that  $\|\pi_{\text{IN+SCC}}(c)\|$  is increasing when  $c$  is small. Consider the point  $c = 0$ . Using (4.7), we get

$$\pi'_{\text{IN+SCC}}(0) = -\alpha(1-\beta)\mathbf{u}_{\text{IN+SCC}} + \alpha\mathbf{u}_{\text{IN+SCC}}P. \quad (4.8)$$

One can see from the above equation that the PageRank of pages in IN+SCC with many incoming links will increase as  $c$  increases from zero, which explains the graphs presented in [Boldi et al. 05]. Next, for the total mass of the IN+SCC component, from (4.8) we obtain

$$\|\pi'_{\text{IN+SCC}}(0)\| = -\alpha(1-\beta)\mathbf{u}_{\text{IN+SCC}} + \alpha\mathbf{u}_{\text{IN+SCC}}P\mathbf{1} = \alpha(-1 + \beta + p_1),$$

where  $p_1 = \mathbf{u}_{\text{IN+SCC}}P\mathbf{1}$  is the probability that a random walk on the hyperlink matrix stays in IN+SCC for one step if the initial distribution is uniform over IN+SCC. If  $1-\beta < p_1$ , then the derivative at 0 is positive. Since dangling nodes typically constitute more than 25% of the graph [Eiron et al. 04], and  $p_1$  is usually close to one, the condition  $1-\beta < p_1$  seems to be comfortably satisfied in web samples. Thus, the total PageRank of IN+SCC increases in  $c$  when  $c$  is small. Note, by the way, that if  $\beta = 0$ , then  $\|\pi_{\text{IN+SCC}}(c)\|$  is strictly decreasing

in  $c$ . Hence, surprisingly, the presence of dangling nodes qualitatively changes the behavior of the IN+SCC PageRank mass.

Now let us consider the point  $c = 1$ . Again using (4.7), we get

$$\pi'_{\text{IN}+\text{SCC}}(1) = -\frac{\alpha}{1-\beta} \mathbf{u}_{\text{IN}+\text{SCC}} \left[ I - P - \frac{\alpha}{1-\beta} S \mathbf{1} \mathbf{u}_{\text{IN}+\text{SCC}} \right]^{-1}. \quad (4.9)$$

We will show that the derivative above is a negative number with a large absolute value. Note that the matrix in the square brackets is close to singular. Denote by  $\bar{P}$  the hyperlink matrix of IN+SCC when the outer links are neglected. Then  $\bar{P}$  is an irreducible stochastic matrix. Denote its stationary distribution by  $\bar{\pi}_{\text{IN}+\text{SCC}}$ . Then we can apply Lemma 6.2 from singular perturbation theory to (4.9) by taking

$$A = \bar{P} \quad \text{and} \quad \varepsilon C = \bar{P} - P - \alpha(1-\beta)^{-1} S \mathbf{1} \mathbf{u}_{\text{IN}+\text{SCC}},$$

and noting that

$$\varepsilon C \mathbf{1} = R \mathbf{1} + (1-\alpha-\beta)(1-\beta)^{-1} S \mathbf{1}.$$

Combining all terms and using

$$\bar{\pi}_{\text{IN}+\text{SCC}} \mathbf{1} = \|\bar{\pi}_{\text{IN}+\text{SCC}}\| = 1 \quad \text{and} \quad \mathbf{u}_{\text{IN}+\text{SCC}} \mathbf{1} = \|\mathbf{u}_{\text{IN}+\text{SCC}}\| = 1,$$

by Lemma 6.2 we obtain

$$\|\pi'_{\text{IN}+\text{SCC}}(1)\| \approx -\frac{\alpha}{1-\beta} \frac{1}{\bar{\pi}_{\text{IN}+\text{SCC}} R \mathbf{1} + \frac{1-\beta-\alpha}{1-\beta} \bar{\pi}_{\text{IN}+\text{SCC}} S \mathbf{1}}.$$

It is expected that the value in the denominator of the second fraction is typically small (indeed, in our data set INRIA, the value is 0.022), and hence the mass  $\|\pi_{\text{IN}+\text{SCC}}(c)\|$  decreases very fast as  $c$  approaches one.

Having described the behavior of the PageRank mass  $\|\pi_{\text{IN}+\text{SCC}}(c)\|$  at the boundary points  $c = 0$  and  $c = 1$ , now we would like to show that there is at most one extremum in  $(0, 1)$ . It is sufficient to prove that if  $\|\pi'_{\text{IN}+\text{SCC}}(c_0)\| \leq 0$  for some  $c_0 \in (0, 1)$  then  $\|\pi'_{\text{IN}+\text{SCC}}(c)\| \leq 0$  for all  $c > c_0$ . To this end, we apply the Sherman–Morrison formula to (4.5), which yields

$$\pi_{\text{IN}+\text{SCC}}(c) = \tilde{\pi}_{\text{IN}+\text{SCC}}(c) + \frac{\frac{c^2\alpha}{1-c\beta} \mathbf{u}_{\text{IN}+\text{SCC}} [I - cP]^{-1} S \mathbf{1}}{1 + \frac{c^2\alpha}{1-c\beta} \mathbf{u}_{\text{IN}+\text{SCC}} [I - cP]^{-1} S \mathbf{1}} \tilde{\pi}_{\text{IN}+\text{SCC}}(c), \quad (4.10)$$

where

$$\tilde{\pi}_{\text{IN}+\text{SCC}}(c) = \frac{(1-c)\alpha}{1-c\beta} \mathbf{u}_{\text{IN}+\text{SCC}} [I - cP]^{-1} \quad (4.11)$$

represents the main term on the right-hand side of (4.10). (The second summand in (4.10) is about 10% of the total sum for the INRIA data set for  $c = 0.85$ .) Now the behavior of  $\pi_{\text{IN}+\text{SCC}}(c)$  in Figure 2 can be explained by means of the following proposition.

**Proposition 4.1.** *The term  $\|\tilde{\pi}_{\text{IN}+\text{SCC}}(c)\|$  given by (4.11) has exactly one local maximum at some  $c_0 \in [0, 1]$ . Moreover,  $\|\tilde{\pi}_{\text{IN}+\text{SCC}}''(c)\| < 0$  for  $c \in (c_0, 1]$ .*

**Proof.** Multiplying both sides of (4.11) by  $\mathbf{1}$  and taking the derivatives, after some tedious algebra we obtain

$$\|\tilde{\pi}_{\text{IN}+\text{SCC}}'(c)\| = -a(c) + \frac{\beta}{1 - c\beta} \|\tilde{\pi}_{\text{IN}+\text{SCC}}(c)\|, \quad (4.12)$$

where the real-valued function  $a(c)$  is given by

$$a(c) = \frac{\alpha}{1 - c\beta} \mathbf{u}_{\text{IN}+\text{SCC}}[I - cP]^{-1}[I - P][I - cP]^{-1}\mathbf{1}.$$

Differentiating (4.12) and substituting  $\frac{\beta}{1 - c\beta} \|\tilde{\pi}_{\text{IN}+\text{SCC}}(c)\|$  from (4.12) into the resulting expression, we get

$$\|\tilde{\pi}_{\text{IN}+\text{SCC}}''(c)\| = \left\{ -a'(c) + \frac{\beta}{1 - c\beta} a(c) \right\} + \frac{2\beta}{1 - c\beta} \|\tilde{\pi}_{\text{IN}+\text{SCC}}'(c)\|.$$

Note that the term in the curly braces is negative by the definition of  $a(c)$ . Hence, if  $\|\tilde{\pi}_{\text{IN}+\text{SCC}}'(c)\| \leq 0$  for some  $c \in [0, 1]$ , then  $\|\tilde{\pi}_{\text{IN}+\text{SCC}}''(c)\| < 0$  for this value of  $c$ .  $\square$

We conclude that  $\|\tilde{\pi}_{\text{IN}+\text{SCC}}(c)\|$  is decreasing and concave for  $c \in [c_0, 1]$ , where  $\|\tilde{\pi}_{\text{IN}+\text{SCC}}'(c_0)\| = 0$ . This is exactly the behavior we observe in our experiments. The analysis and experiments suggest that  $c_0$  is definitely larger than 0.85 and actually is quite close to one. Thus, one may want to choose a large value for  $c$  in order to maximize the PageRank mass of IN+SCC. However, in the next section we will indicate important drawbacks of this choice.

## 5. PageRank Mass of ESCC

Let us now consider the PageRank mass of the extended SCC component (ESCC) described in Section 3 as a function of  $c \in [0, 1]$ . Subdividing the PageRank vector in the blocks  $\pi = [\pi_{\text{PureOUT}}, \pi_{\text{ESCC}}]$ , from (4.1) we obtain

$$\pi_{\text{ESCC}}(c) = (1 - c)\gamma \mathbf{u}_{\text{ESCC}}[I - cT]^{-1} = (1 - c)\gamma \mathbf{u}_{\text{ESCC}} \sum_{k=1}^{\infty} c^k T^k, \quad (5.1)$$

where  $T$  represents the transition probabilities inside the ESCC block,  $\gamma = |\text{ESCC}|/n$  is the fraction of pages contained in the ESCC, and  $\mathbf{u}_{\text{ESCC}}$  is a uniform-probability row vector over ESCC. Clearly, we have that  $\|\pi_{\text{ESCC}}(0)\| = \gamma$  and  $\|\pi_{\text{ESCC}}(1)\| = 0$ . Furthermore, it is easy to see that  $\|\pi_{\text{ESCC}}(c)\|$  is a concave decreasing function, since

$$\frac{d}{dc} \|\pi_{\text{ESCC}}(c)\| = -\gamma \mathbf{u}_{\text{ESCC}} [I - cT]^{-2} [I - T] \mathbf{1} < 0$$

and

$$\frac{d^2}{dc^2} \|\pi_{\text{ESCC}}(c)\| = -2\gamma \mathbf{u}_{\text{ESCC}} [I - cT]^{-3} T [I - T] \mathbf{1} < 0.$$

The next proposition establishes upper and lower bounds for  $\|\pi_{\text{ESCC}}(c)\|$ .

**Proposition 5.1.** *Let  $\lambda_1$  be the Perron–Frobenius eigenvalue of  $T$ , and let  $p_1 = \mathbf{u}_{\text{ESCC}} T \mathbf{1}$  be the probability that the random walk started from a randomly chosen state in ESCC stays in ESCC for one step. If  $p_1 \leq \lambda_1$  and*

$$p_1 \leq \frac{\mathbf{u}_{\text{ESCC}} T^k \mathbf{1}}{\mathbf{u}_{\text{ESCC}} T^{k-1} \mathbf{1}} \leq \lambda_1 \quad \text{for all } k \geq 1, \quad (5.2)$$

*then*

$$\frac{\gamma(1-c)}{1-cp_1} < \|\pi_{\text{ESCC}}(c)\| < \frac{\gamma(1-c)}{1-c\lambda_1}, \quad c \in (0, 1). \quad (5.3)$$

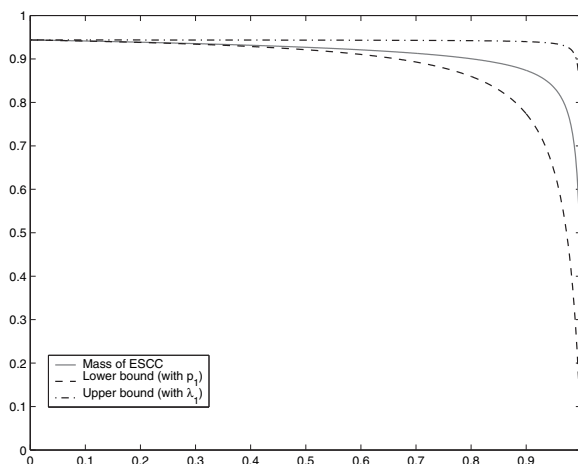
**Proof.** From condition (5.2) it follows by induction that

$$p_1^k \leq \mathbf{u}_{\text{ESCC}} T^k \mathbf{1} \leq \lambda_1^k, \quad k \geq 1,$$

and thus the statement of the proposition is obtained directly from the series expansion of  $\pi_{\text{ESCC}}(c)$  in (5.1).  $\square$

The conditions of Proposition 5.1 have a natural probabilistic interpretation. The value  $p_1$  is the probability that the Markov random walk on the web sample stays in the block  $T$  for one step, starting from the uniform distribution over  $T$ . Furthermore,  $p_k = \mathbf{u}_{\text{ESCC}} T^k \mathbf{1} / (\mathbf{u}_{\text{ESCC}} T^{k-1} \mathbf{1})$  is the probability that the random walk stays in  $T$  for one step provided that it has stayed there for the first  $k-1$  steps.

It is a well-known fact that as  $k \rightarrow \infty$ ,  $p_k$  converges to  $\lambda_1$ , the Perron–Frobenius eigenvalue of  $T$ . Let  $\hat{\pi}_{\text{ESCC}}$  be the probability-normed left Perron–Frobenius eigenvector of  $T$ . Then  $\hat{\pi}_{\text{ESCC}}$ , also known as a *quasistationary* distribution of  $T$ , is the limiting probability distribution of the Markov chain given that the random walk never leaves the block  $T$  (see, e.g., [Seneta 06]). Since



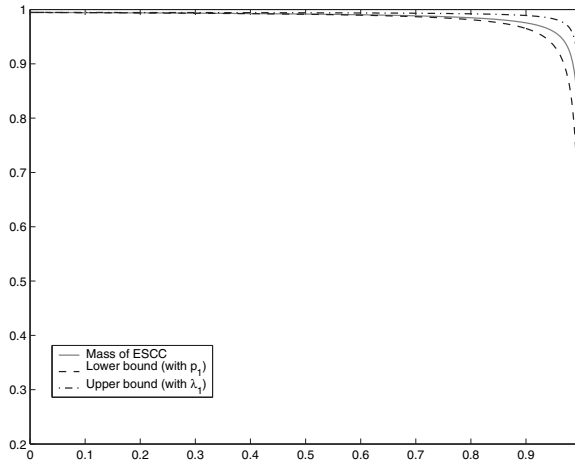
**Figure 3.** PageRank mass of ESCC and bounds for INRIA .

$\hat{\pi}_{\text{ESCC}}T = \lambda_1 \hat{\pi}_{\text{ESCC}}$ , the condition  $p_1 < \lambda_1$  means that the chance of staying in ESCC for one step in the quasistationary regime is higher than that in starting from the uniform distribution  $\mathbf{u}_{\text{ESCC}}$ . This is quite natural, since the quasistationary distribution tends to avoid the states from which the random walk is likely to leave the block  $T$ .

Furthermore, the condition in (5.2) says that if the random walk is about to make its  $k$ th step in  $T$ , then it leaves  $T$  most easily at step  $k = 1$ , and it is most difficult to leave  $T$  after an infinite number of steps. Both conditions of Proposition 5.1 are satisfied in our experiments on both data sets. Moreover, we noticed that the sequence  $(p_k)$ ,  $k \geq 1$ , was increasing from  $p_1$  to  $\lambda_1$ .

With the help of the derived bounds we conclude that  $\|\pi_{\text{ESCC}}(c)\|$  decreases very slowly for small and moderate values of  $c$ , and it decreases extremely fast when  $c$  becomes close to 1. This typical behavior is clearly seen in Figures 3 and 4, where  $\|\pi_{\text{ESCC}}(c)\|$  is plotted with a solid line. The bounds are plotted with dashed lines. For the INRIA data set we have  $p_1 = 0.97557$ ,  $\lambda_1 = 0.99954$ , and for the FrMathInfo data set we have  $p_1 = 0.99659$ ,  $\lambda_1 = 0.99937$ .

From the above we conclude that the PageRank mass of ESCC is smaller than  $\gamma$  for any value  $c > 0$ . In contrast, the PageRank mass of pure OUT increases in  $c$  beyond its “fair share”  $\delta = |\text{pure OUT}|/n$ . With  $c = 0.85$ , the PageRank mass of the pure OUT component in the INRIA data set is equal to  $1.95\delta$ . In the FrMathInfo data set, the unfairness is even more pronounced: the PageRank mass of the pure OUT component is equal to  $3.44\delta$ . This gives users an incentive to create dead ends: groups of pages that link only to each other. Clearly, this



**Figure 4.** PageRank mass of ESCC and bounds for FrMathInfo.

can be mitigated by choosing a smaller damping factor. Below we propose one way to determine an “optimal” value of  $c$ .

Since the PageRank mass of ESCC is always smaller than  $\gamma$ , we would like to choose the damping factor in such a way that the ESCC receives a “fair” fraction of  $\gamma$ . Formally, we would like to define a number  $\rho \in (0, 1)$  such that a desirable PageRank mass of ESCC could be written as  $\rho\gamma$ , and then find the value  $c^*$  that satisfies

$$\|\pi_{\text{ESCC}}(c^*)\| = \rho\gamma. \quad (5.4)$$

Then  $c \leq c^*$  will ensure that  $\|\pi_{\text{ESCC}}(c)\| \geq \rho\gamma$ . Naturally,  $\rho$  should somehow reflect the properties of the substochastic block  $T$ . For instance, as  $T$  becomes closer to a stochastic matrix,  $\rho$  should also increase. One possibility is to define

$$\rho = \mathbf{v}T\mathbf{1},$$

where  $\mathbf{v}$  is a row vector representing some probability distribution on ESCC. Then the damping factor  $c$  should satisfy

$$c \leq c^*,$$

where  $c^*$  is given by

$$\|\pi_{\text{ESCC}}(c^*)\| = \gamma\mathbf{v}T\mathbf{1}. \quad (5.5)$$

In this setting,  $\rho$  is the probability of staying in ESCC for one step if the initial distribution is  $\mathbf{v}$ . For a given  $\mathbf{v}$ , this number increases as  $T$  becomes closer to

$\mathbf{v}$	$\mathbf{c}$	INRIA	FrMathInfo
$\hat{\pi}_{\text{ESCC}}$	$c_1$	0.0184	0.1956
	$c_2$	0.5001	0.5002
	$c^*$	.02	.16
$\mathbf{u}_{\text{ESCC}}$	$c_1$	0.5062	0.5009
	$c_2$	0.9820	0.8051
	$c^*$	.604	.535
$\pi_{\text{ESCC}}/\ \pi_{\text{ESCC}}\ $	$1/(1 + \lambda_1)$	0.5001	0.5002
	$1/(1 + p_1)$	0.5062	0.5009

**Table 2.** Values of  $c^*$  with bounds.

a stochastic matrix. The problem of choosing  $\rho$  comes down to the problem of choosing  $\mathbf{v}$ . The advantage of this approach is twofold. First, we have lost no flexibility, because depending on  $\mathbf{v}$ , the value of  $\rho$  may vary considerably, except it cannot become too small if  $T$  is really close to a stochastic matrix. Second, we can use a probabilistic interpretation of  $\mathbf{v}$  to make a reasonable choice.

One can think, for instance, of the following three intuitive choices of  $\mathbf{v}$ : (1)  $\hat{\pi}_{\text{ESCC}}$ , the quasistationary distribution of  $T$ , (2) the uniform vector  $\mathbf{u}_{\text{ESCC}}$ , and (3) the normalized PageRank vector  $\pi_{\text{ESCC}}(c)/\|\pi_{\text{ESCC}}(c)\|$ . The first choice reflects the proximity of  $T$  to a stochastic matrix. The second choice is inspired by the definition of PageRank (restart from the uniform distribution), and the third choice combines both these features.

If the conditions of Proposition 5.1 are satisfied, then (5.3) holds, and thus the value of  $c^*$  satisfying (5.5) must be in the interval  $(c_1, c_2)$ , where

$$\frac{1 - c_1}{1 - p_1 c_1} = \|\mathbf{v}T\|, \quad \frac{1 - c_2}{1 - \lambda_1 c_2} = \|\mathbf{v}T\|.$$

Numerical results for all three choices of  $\mathbf{v}$  are presented in Table 2.

If  $\mathbf{v} = \hat{\pi}_{\text{ESCC}}$  then we have  $\|\mathbf{v}T\| = \lambda_1$ , which implies  $c_1 = (1 - \lambda_1)/(1 - \lambda_1 p_1)$  and  $c_2 = 1/(\lambda_1 + 1)$ . In this case, the upper bound  $c_2$  is only slightly larger than  $\frac{1}{2}$ , and  $c^*$  is close to zero in our data sets (see Table 2). Such small  $c$ , however, leads to ranking that takes into account only local information about the web graph (see, e.g., [Fortunato and Flammini 06]). The choice  $\mathbf{v} = \hat{\pi}_{\text{ESCC}}$  does not seem to represent the dynamics of the system, probably because the “easily bored surfer” random walk that is used in PageRank computations never follows a quasistationary distribution, since it often restarts itself from the uniform probability vector.

For the uniform vector  $\mathbf{v} = \mathbf{u}_{\text{ESCC}}$ , we have  $\|\mathbf{v}T\| = p_1$ , which gives  $c_1, c_2, c^*$ , presented in Table 2. We have obtained a higher upper bound, but the values of  $c^*$  are still much smaller than 0.85.



Finally, consider the normalized PageRank vector  $\mathbf{v}(c) = \pi_{\text{ESCC}}(c) / \|\pi_{\text{ESCC}}(c)\|$ . This choice of  $\mathbf{v}$  can also be justified as follows. Consider the derivative of the total PageRank mass of ESCC. Since  $[I - cT]^{-1}$  and  $[I - T]$  commute, we can write

$$\frac{d}{dc} \|\pi_{\text{ESCC}}(c)\| = -\gamma \mathbf{u}_{\text{ESCC}} [I - cT]^{-1} [I - T] [I - cT]^{-1} \mathbf{1},$$

or equivalently,

$$\begin{aligned} \frac{d}{dc} \|\pi_{\text{ESCC}}(c)\| &= -\frac{1}{1-c} \pi_{\text{ESCC}} [I - T] [I - cT]^{-1} \mathbf{1} \\ &= -\frac{1}{1-c} \left( \pi_{\text{ESCC}} - \|\pi_{\text{ESCC}}\| \frac{\pi_{\text{ESCC}}}{\|\pi_{\text{ESCC}}\|} T \right) [I - cT]^{-1} \mathbf{1} \\ &= -\frac{1}{1-c} (\pi_{\text{ESCC}} - \|\pi_{\text{ESCC}}\| \mathbf{v}(c) T) [I - cT]^{-1} \mathbf{1}, \end{aligned}$$

with  $\mathbf{v}(c) = \pi_{\text{ESCC}} / \|\pi_{\text{ESCC}}\|$ . It is easy to see that

$$\|\pi_{\text{ESCC}}(c)\| = \gamma - \gamma(1 - \mathbf{u}_{\text{ESCC}} T \mathbf{1})c + o(c).$$

Consequently, we obtain

$$\begin{aligned} \frac{d}{dc} \|\pi_{\text{ESCC}}(c)\| &= -\frac{1}{1-c} (\pi_{\text{ESCC}} - \gamma \mathbf{v}(c) T + \gamma(1 - \mathbf{u}_{\text{ESCC}} T \mathbf{1}) c \mathbf{v}(c) T + o(c)) [I - cT]^{-1} \mathbf{1}. \end{aligned}$$

Since in practice  $T$  is very close to stochastic, we have

$$1 - \mathbf{u}_{\text{ESCC}} T \mathbf{1} \approx 0 \quad \text{and} \quad [I - cT]^{-1} \mathbf{1} \approx \frac{1}{1-c} \mathbf{1}.$$

The latter approximation follows from Lemma 6.2. Thus, satisfying condition (5.5) means keeping the value of the derivative small.

Let us now solve (5.5) for  $\mathbf{v}(c) = \pi_{\text{ESCC}}(c) / \|\pi_{\text{ESCC}}(c)\|$ . Using (5.1), we rewrite (5.5) as

$$\|\pi_{\text{ESCC}}(c)\| = \frac{\gamma}{\|\pi_{\text{ESCC}}(c)\|} \pi_{\text{ESCC}}(c) T \mathbf{1} = \frac{\gamma^2(1-c)}{\|\pi_{\text{ESCC}}(c)\|} \mathbf{u}_{\text{IN+SCC}} [I - cT]^{-1} T \mathbf{1}.$$

Multiplying by  $\|\pi_{\text{ESCC}}(c)\|$ , after some algebra we obtain

$$\|\pi_{\text{ESCC}}(c)\|^2 = \frac{\gamma}{c} \|\pi_{\text{ESCC}}(c)\| - \frac{(1-c)\gamma^2}{c}.$$

Solving the quadratic equation for  $\|\pi_{\text{ESCC}}(c)\|$ , we get

$$\|\pi_{\text{ESCC}}(c)\| = r(c) = \begin{cases} \gamma & \text{if } c \leq \frac{1}{2}, \\ \frac{\gamma(1-c)}{c} & \text{if } c > \frac{1}{2}. \end{cases}$$

Hence, the value  $c^*$  solving (5.5) corresponds to the point where the graphs of  $\|\pi_{\text{ESCC}}(c)\|$  and  $r(c)$  cross each other. There is only one such point in  $(0, 1)$ , and since  $\|\pi_{\text{ESCC}}(c)\|$  decreases very slowly unless  $c$  is close to one, whereas  $r(c)$  decreases relatively fast for  $c > \frac{1}{2}$ , we expect that  $c^*$  is only slightly larger than  $\frac{1}{2}$ . Under the conditions of Proposition 5.1,  $r(c)$  first crosses the line  $\gamma(1 - c) \div (1 - \lambda_1 c)$ , then  $\|\pi_{\text{ESCC}}(c)\|_1$ , and then  $\gamma(1 - c)/(1 - p_1 c)$ . Thus, we obtain  $(1 + \lambda_1)^{-1} < c^* < (1 + p_1)^{-1}$ . Since both  $\lambda_1$  and  $p_1$  are very close to 1, this suggests that  $c$  should be chosen around  $\frac{1}{2}$ . This is also reflected in Table 2.

Last but not least, to support our theoretical argument about the undeserved high ranking of pages from pure OUT, we carry out the following experiment. In the INRIA data set we have chosen an absorbing component in pure OUT consisting just of two nodes. We have added an artificial link from one of these nodes to a node in the giant SCC and recomputed the PageRank.

In Table 3 in the column “PR rank w/o link” we give a ranking of a page according to the PageRank value computed before the addition of the artificial link, and in the column “PR rank with link” we give a ranking of a page according to the PageRank value computed after the addition of the artificial link. We have also analyzed the log file of the site INRIA Sophia Antipolis ([www-sop.inria.fr](http://www-sop.inria.fr)) and ranked the pages according to the number of clicks for the period of one year up to May 2007. We note that since we have access only to the log file of the INRIA Sophia Antipolis site, we also use the PageRank ranking only for the pages from the INRIA Sophia Antipolis site. For instance, for  $c = 0.85$ , the ranking of page A without an artificial link is 731 (this means that 730 pages are ranked higher than page A among the pages of INRIA Sophia Antipolis). However, its ranking according to the number of clicks is much lower, 2588.

This confirms our conjecture that the nodes in pure OUT obtain unjustifiably high ranking. Next, we note that the addition of an artificial link significantly diminishes the ranking. In fact, it brings it close to the ranking provided by the number of clicks. Finally, we draw the reader’s attention to the fact that choosing  $c = \frac{1}{2}$  also significantly reduces the gap between the ranking by PageRank and the ranking by the number of clicks.

To summarize, our results indicate that with  $c = 0.85$ , the pure OUT component receives an unfairly large share of the PageRank mass. Remarkably, in order to satisfy any of the three intuitive criteria of fairness presented above, the value of  $c$  should be drastically reduced. The experiment with the log files

$c$	PR rank w/o link	PR rank with link	rank by no. of clicks
Node A			
0.5	1648	2307	2588
0.85	731	2101	2588
0.95	226	2116	2588
Node B			
0.5	1648	4009	3649
0.85	731	3279	3649
0.95	226	3563	3649

**Table 3.** Comparison between PR- and click-based rankings.

confirms the same. Of course, a drastic reduction of  $c$  also considerably accelerates the computation of PageRank by numerical methods [Avrachenkov et al. 07, Langville and Meyer 06, Berkhin 05].

## 6. Appendix: Results from Singular Perturbation Theory

**Lemma 6.1.** *Let  $A(\varepsilon) = A + \varepsilon C$  be a transition matrix of a perturbed Markov chain. The perturbed Markov chain is assumed to be ergodic for sufficiently small  $\varepsilon$  different from zero. Let the unperturbed Markov chain ( $\varepsilon = 0$ ) have  $m$  ergodic classes. Namely, the transition matrix  $A$  can be written in the form*

$$A = \begin{bmatrix} A_1 & & 0 & 0 \\ & \ddots & & \\ 0 & & A_m & 0 \\ L_1 & \cdots & L_m & E \end{bmatrix} \in \mathbb{R}^{n \times n}.$$

*Then the stationary distribution of the perturbed Markov chain has a limit*

$$\lim_{\varepsilon \rightarrow 0} \pi(\varepsilon) = [\nu_1 \mu_1 \ \cdots \ \nu_m \mu_m \ 0],$$

*where zeros correspond to the set of transient states in the unperturbed Markov chain,  $\mu_i$  is a stationary distribution of the unperturbed Markov chain corresponding to the  $i$ th ergodic set, and  $\nu_i$  is the  $i$ th element of the aggregated stationary distribution vector that can be found by solving*

$$\nu D = \nu, \quad \nu \mathbf{1} = 1,$$

where  $D = MCB$  is the generator of the aggregated Markov chain and

$$M = \begin{bmatrix} \mu_1 & & 0 & 0 \\ & \ddots & & \\ 0 & & \mu_m & 0 \end{bmatrix} \in \mathbb{R}^{m \times n}, \quad B = \begin{bmatrix} \mathbf{1} & & 0 \\ & \ddots & \\ 0 & & \mathbf{1} \\ \phi_1 & \cdots & \phi_m \end{bmatrix} \in \mathbb{R}^{n \times m},$$

with  $\phi_i = [I - E]^{-1} L_i \mathbf{1}$ .

The proof of this lemma can be found in [Avrachenkov 99, Korolyuk and Turbin 93, Yin and Zhang 05].

**Lemma 6.2.** *Let  $A(\varepsilon) = A - \varepsilon C$  be a perturbation of an irreducible stochastic matrix  $A$  such that  $A(\varepsilon)$  is substochastic. Then for sufficiently small  $\varepsilon$  the following Laurent series expansion holds:*

$$[I - A(\varepsilon)]^{-1} = \frac{1}{\varepsilon} X_{-1} + X_0 + \varepsilon X_1 + \cdots, \quad (6.1)$$

with

$$X_{-1} = \frac{1}{\mu C \mathbf{1}} \mathbf{1} \mu, \quad (6.2)$$

where  $\mu$  is the stationary distribution of  $A$ . It follows that

$$[I - A(\varepsilon)]^{-1} = \frac{1}{\mu \varepsilon C \mathbf{1}} \mathbf{1} \mu + O(1) \quad \text{as } \varepsilon \rightarrow 0. \quad (6.3)$$

**Proof.** The proof of this result is based on the approach developed in [Avrachenkov 99, Avrachenkov et al. 01]. The existence of the Laurent series (6.1) is a particular case of more-general results on the inversion of analytic matrix functions [Avrachenkov et al. 01]. To calculate the terms of the Laurent series, let us equate the terms with the same powers of  $\varepsilon$  in the following identity:

$$(I - A + \varepsilon C) \left( \frac{1}{\varepsilon} X_{-1} + X_0 + \varepsilon X_1 + \cdots \right) = I,$$

which results in

$$(I - A)X_{-1} = 0, \quad (6.4)$$

$$(I - A)X_0 + CX_{-1} = I, \quad (6.5)$$

$$(I - A)X_1 + CX_0 = 0. \quad (6.6)$$

From equation (6.4) we conclude that

$$X_{-1} = \mathbf{1}\mu_{-1}, \quad (6.7)$$

where  $\mu_{-1}$  is some vector. We find this vector from the condition that (6.5) has a solution. In particular, (6.5) has a solution if and only if

$$\mu(I - CX_{-1}) = 0.$$

By substituting the expression (6.7) into the above equation, we obtain

$$\mu - \mu C\mathbf{1}\mu_{-1} = 0,$$

and consequently,

$$\mu_{-1} = \frac{1}{\mu C\mathbf{1}}\mu,$$

which together with (6.7) gives (6.2).  $\square$

**Acknowledgments.** This work was supported by EGIDE ECO-NET grant no. 10191XC and by NWO Meervoud grant no. 632.002.401.

## References

- [Avrachenkov 99] K. Avrachenkov. “Analytic Perturbation Theory and Its Applications.” PhD thesis, University of South Australia, 1999.
- [Avrachenkov and Litvak 06] K. Avrachenkov and N. Litvak. “The Effect of New Links on Google PageRank.” *Stoch. Models* 22:2 (2006), 319–331.
- [Avrachenkov et al. 01] K. Avrachenkov, M. Haviv, and P. Howlett. “Inversion of Analytic Matrix Functions That Are Singular at the Origin.” *SIAM Journal on Matrix Analysis and Applications* 22:4 (2001), 1175–1189.
- [Avrachenkov et al. 07] K. Avrachenkov, N. Litvak, D. Nemirovsky, and N. Osipova. “Monte Carlo Methods in PageRank Computation: When One Iteration Is Sufficient.” *SIAM J. Numer. Anal.* 45:2 (2007), 890–904.
- [Berkhin 05] P. Berkhin. “A Survey on PageRank Computing.” *Internet Math.* 2 (2005), 73–120.
- [Bianchini et al. 05] M. Bianchini, M. Gori, and F. Scarselli. “Inside PageRank.” *ACM Trans. Inter. Tech.* 5:1 (2005), 92–128.
- [Boldi et al. 05] P. Boldi, M. Santini, and S. Vigna. “PageRank as a Function of the Damping Factor.” In *Proc. of the Fourteenth International World Wide Web Conference, Chiba, Japan*. New York: ACM Press, 2005.
- [Broder et al. 00] A. Broder, R. Kumar, F. Maghoul, P. Raghavan, S. Rajagopalan, R. Statac, A. Tomkins, and J. Wiener. “Graph Structure in the Web.” *Computer Networks* 33 (2000), 309–320.

- [Chen et al. 06] P. Chen, H. Xie, S. Maslov, and S. Redner. “Finding Scientific Gems with Google’s PageRank Algorithm.” *Journal of Informatics* 1:1 (2007), 8–15.
- [Dill et al. 02] S. Dill, R. Kumar, K. S. Mccurley, S. Rajagopalan, D. Sivakumar, and A. Tomkins. “Self-Similarity in the Web.” *ACM Trans. Inter. Tech.* 2:3 (2002), 205–223.
- [Eiron et al. 04] N. Eiron, K. McCurley, and J. Tomlin. “Ranking the Web Frontier.” In *Proceedings of the 13th International Conference on the World Wide Web*, pp. 309–318. New York: ACM Press, 2004.
- [Fortunato and Flammini 06] S. Fortunato, and A. Flammini. “Random Walks on Directed Networks: The Case of PageRank.” *International Journal of Bifurcation and Chaos* 17:7 (2007), 2343–2353.
- [Kleinberg 99] J. M. Kleinberg. “Authoritative Sources in a Hyperlinked Environment.” *Journal of the ACM* 46:5 (1999), 604–632.
- [Korolyuk and Turbin 93] V. S. Korolyuk and A. F. Turbin. *Mathematical Foundations of the State Lumping of Large Systems*, Mathematics and Its Applications 264. Dordrecht: Kluwer, 1993.
- [Kumar et al. 00] R. Kumar, P. Raghavan, S. Rajagopalan, D. Sivakumar, A. Tomkins, and E. Upfal. “The Web as a Graph.” In *Proceedings of the 19th ACM SIGACT-SIGMOD-AIGART Symposium on Principles of Database Systems*, pp. 1–10. New York: ACM Press, 2000.
- [Langville and Meyer 03] A. N. Langville and C. D. Meyer. “Deeper inside PageRank.” *Internet Math.* 1 (2003), 335–380.
- [Langville and Meyer 06] A. N. Langville and C. D. Meyer. *Google’s PageRank and Beyond: The Science of Search Engine Rankings*. Princeton: Princeton University Press, 2006.
- [Lempel and Moran 00] R. Lempel and S. Moran. “The Stochastic Approach for Link-Structure Analysis (SALSA) and the TKE Effect.” *Comput. Networks* 33:1–6 (2000), 387–401.
- [Moler and Moler 03] C. Moler, and K. Moler. *Numerical Computing with MATLAB*. Philadelphia: SIAM, 2003.
- [Page et al. 98] L. Page, S. Brin, R. Motwani, and T. Winograd. “The PageRank Citation Ranking: Bringing Order to the Web.” Technical report, Stanford University, 1998.
- [Pervozvanskii and Gaitsgori 88] A. A. Pervozvanskii and V. G. Gaitsgori. *Theory of Suboptimal Decisions*, Mathematics and Its Applications (Soviet Series) 12. Dordrecht: Kluwer, 1988.
- [Seneta 06] E. Seneta. *Non-negative Matrices and Markov Chains*, Springer Series in Statistics; revised reprint of the second (1981) edition. New York: Springer, 2006.
- [Yin and Zhang 05] G. G. Yin and Q. Zhang. *Discrete-Time Markov Chains*. Applications of Mathematics (New York) 55. New York: Springer, 2005.

Konstantin Avrachenkov, INRIA Sophia Antipolis, 2004, Route des Lucioles, 06902, France (k.avrachenkov@sophia.inria.fr)

Nelly Litvak, University of Twente, Dept. of Applied Mathematics, P.O. Box 217, 7500AE Enschede, the Netherlands (n.litvak@ewi.utwente.nl)

Kim Son Pham, St. Petersburg State University, 35, University Prospect, 198504, Peterhof, Russia (sonsecure@yahoo.com.sg)

Received February 14, 2008; accepted May 6, 2008.