

# PageRank

Caroline, Victoria and Magnus

January 2026

Network Visualization (Node size = PageRank)

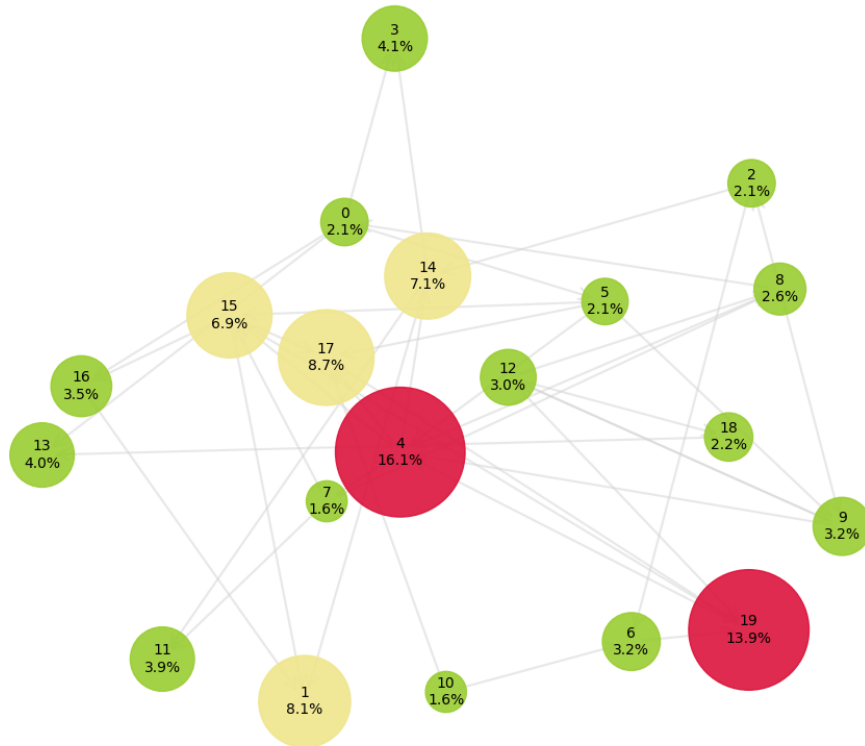


Figure 1: Plot of matrix\_pagerank model

## 1 Introduction

In this project we explore the PageRank algorithm. This algorithm considers a finite amount of webpages connected by directed links. It ranks the pages based on the probability of ending at a certain page  $k$  at a given time  $t$ , if you start at a randomly selected page, and move through the links after the rule, that every time you end up at a page you move to the next page by randomly selecting a link on that page. The algorithm also has a damping factor, which allows the simulation to jump to a random page, even if it lands on a page with no outgoing links.

There are different approaches to the algorithm.

The first we looked at is the Random surfer model. In this model we start by choosing a random page

to start at, then after that choose a link between the ones on the page. Leading to the next page, and keeps repeating that  $t$  times. If you end up at a page, that don't link to any other page, you choose a random page among all the pages. To get the probability you count how many times you visit the page  $N_p$  over the number of pages visited  $n$ . They should all add up to 1.

Another approach is the Recursive model. This model looks at the probability of being at a page on a certain time by first setting all the probabilities to  $1/N$ , representing starting on a page in the Random Surf. Thereafter the formula,  $PR(p) = (1 - d) * \frac{1}{N} + d * \sum_{q \in Inbound(p)} \frac{PR(q)}{Outbound}$ , is applied to all the pages to calculate the probability of reaching that page on the next step. This step is done a sufficient number of times to get the probability distribution, which the function will output as the PageRank.

The last models are matrix formulations. There are two different approaches in calculating the PageRank from the matrix, but both of them depend on formulating a modified adjacency matrix, which describes how the pages are linked, where if a page isn't linked to any other, it is set as linking to all pages.

The first matrix approach finds an eigenvector corresponding to the eigenvalue 1 and multiply it to ensure, that the entries are positive and sum to 1. The other matrix approach calculates the PageRank by raising the modified adjacency matrix to a high enough power, and choosing the first column, which will correspond to the eigenvector described above.

All the different approaches also include a damping factor, which implements a way to get out of parts of the web, even if it is not connected to the rest of the web.

The code implementation is in the Appendix.

## 2 Theoretical problems

### Problem 1

Assuming that  $\mathbf{x}$  is a probability vector, i.e. its elements sum to 1, show that Equation (2.3):  $\mathbf{x} = \frac{(1-d)}{N} \cdot \mathbf{e} + dL\mathbf{x}$  can be written as  $\mathbf{x} = M_d\mathbf{x}$  (Equation 2.4)

#### Proof

From Equation 2.3 the components of  $\mathbf{x}$  will be:

$$\mathbf{x}_i = \frac{(1-d)}{N} \cdot 1 + d \cdot \sum_{k=1}^N (L_{ik}\mathbf{x}_k)$$

According to Equation 2.4 the components of  $\mathbf{x}$  will be:

$$\begin{aligned} \mathbf{x}_i &= \sum_{k=1}^N \left( \frac{(1-d)}{N} + d \cdot L_{ik} \right) \mathbf{x}_k = \sum_{k=1}^N \left( \frac{(1-d)}{N} \cdot \mathbf{x}_k + d \cdot L_{ik} \cdot \mathbf{x}_k \right) \\ &= \sum_{k=1}^N \left( \frac{(1-d)}{N} \cdot \mathbf{x}_k \right) + \sum_{k=1}^N (d \cdot L_{ik} \cdot \mathbf{x}_k) = \frac{(1-d)}{N} \cdot \sum_{k=1}^N (\mathbf{x}_k) + d \cdot \sum_{k=1}^N (L_{ik} \cdot \mathbf{x}_k) \end{aligned}$$

Since  $\mathbf{x}$  is a probability vector:  $\sum_{k=1}^N x_k = 1$ :

$$\frac{(1-d)}{N} \cdot \sum_{k=1}^N (\mathbf{x}_k) + d \cdot \sum_{k=1}^N (L_{ik} \cdot \mathbf{x}_k) = \frac{(1-d)}{N} \cdot 1 + d \cdot \sum_{k=1}^N (L_{ik} \cdot \mathbf{x}_k)$$

And thus  $\mathbf{x}_i$  is the same for both Equation 2.3 and 2.4, which shows that Equation 2.3 can be written as Equation 2.4.  $\square$

## Problem 2

Lemma 3.1: If  $A$  is a Markov matrix then  $A^m$  is also a Markov matrix for  $m = 0, 1, 2, \dots$

### Proof

$A$  and  $B$  are both Markov matrices, therefore we know that the sum of each of their columns is 1:

$$\sum_i A_{ij} = 1, \quad \sum_i B_{ij} = 1$$

We can now write the sum of each column in the product  $AB$  by using the definition of a matrix product:

$$\begin{aligned} \sum_i (AB)_{ij} &= \sum_i \left( \sum_k A_{ik} B_{kj} \right) \\ &= \sum_i A_{i1} B_{1j} + \sum_i A_{i2} B_{2j} + \dots + \sum_i A_{in} B_{nj} \\ &= B_{1j} \sum_i A_{i1} + B_{2j} \sum_i A_{i2} + \dots + B_{nj} \sum_i A_{in} \\ &= B_{1j} + B_{2j} + \dots + B_{nj} \\ &= \sum_i B_{ij} = 1 = 1 \end{aligned}$$

The product of two Markov matrices is therefore also a Markov matrix.  $A^2$  is a product of two Markov matrices  $A$  and  $A$ , and is therefore also a Markov matrix.  $A^m$  for  $m = 0, 1, 2, \dots$ , can be seen as a product of the Markov matrices  $A$  and  $A^{m-1}$ , and is therefore also a Markov matrix.  $\square$

## Problem 3

Lemma 3.2: If  $A \in \mathbb{R}^{n \times n}$  is a Markov matrix then the spectral radius  $\text{rad}_A = 1$

### Proof

We split the proof into three main steps: showing  $A$  has  $\lambda = 1$ , proving norm inequality, and finally proving the spectral radius is  $\text{rad}_A = 1$ .

#### Step 1: $A$ has $\lambda = 1$

We start by proving that  $A$ , has an eigenvalue equal to 1. To do that, we find the eigenvalues for the transposed of  $A$ . We know that the transposed of  $A_k$  has an eigenvalue, that is equal to 1 as shown in Equation (3.2). This means that  $A$  also has an eigenvalue that is equal to 1, since:

$$\det(A - \lambda \cdot I_n) = \det(A^T - \lambda \cdot I_n)$$

As you can expand a determinant by a row or a column.

#### Step 2: Norm inequality

Now we show the inequality  $\|Xv\|_\infty \leq \|X\|_\infty \cdot \|v\|_\infty$ . Let the rows in  $X$  be denoted by  $w_i$  for  $i = 1, 2, \dots, m$ , then we have

$$\begin{aligned} \|Xv\|_\infty &= \max_i (|w_i \cdot v|) \\ &\leq \max_i (|(w_i)_1| \cdot |v_1| + |(w_i)_2| \cdot |v_2| + \dots + |(w_i)_n| \cdot |v_n|) \\ &\leq \max_i \left( \sum_j |(w_i)_j| \cdot \|v\|_\infty \right) \\ &= \|v\|_\infty \cdot \max_i \sum_j |(w_i)_j| = \|X\|_\infty \cdot \|v\|_\infty \end{aligned}$$

Which shows the inequality holds.

**Step 3: Spectral radius is  $rad_A = 1$**

Now we want to prove the spectral radius is  $rad_A = 1$ . We will prove by contradiction. Suppose that  $\lambda$  is an eigenvalue with  $|\lambda| > 1$ . Let  $v$  be a corresponding eigenvector and consider  $A^k v$  as  $k \rightarrow \infty$ . We have that  $A^k v = \lambda^k v$ , then take the norm

$$\|A^k v\|_\infty = \|\lambda^k v\|_\infty = |\lambda^k| \cdot \|v\|_\infty$$

By Step 2, we have  $\|\lambda^k v\|_\infty \leq \|A^k\|_\infty \cdot \|v\|_\infty$ . Furthermore, we know that  $\|A^k\|_\infty \leq n$ , so

$$\begin{aligned} \|\lambda^k v\|_\infty &\leq \|A^k\|_\infty \cdot \|v\|_\infty \leq n \cdot \|v\|_\infty \\ \Rightarrow |\lambda^k| \cdot \|v\|_\infty &\leq n \cdot \|v\|_\infty \\ \Rightarrow |\lambda^k| &\leq n \end{aligned}$$

Since, for any  $k$ , our eigenvalue  $\lambda$  has to be  $-1 \leq \lambda \leq 1$ , because any value greater would grow to infinity as  $k \rightarrow \infty$ , which will eventually be greater than  $n$ , since our  $n$  is finite. We have reached a contradiction and thus the maximum absolute value among the eigenvalues of  $A$  is  $rad_A = 1$ .  $\square$

**Problem 4**

If  $A$  is Markov matrix and  $0 < \tau < 1$  prove that  $A_\tau$  is also a Markov matrix and has strictly positive entries,  $(A_\tau)_{ij} > 0$ .

**Proof**

First we show that  $A_\tau$  is also a Markov matrix. Let  $0 < \tau < 1$ . For  $A \in \mathbb{R}^{n \times n}$ , set:

$$A_\tau = \tau \cdot A + \frac{1-\tau}{n} \cdot E_n$$

where  $E_n$  is a  $n \times n$  matrix consisting entirely of 1's.

If  $A$  is a Markov matrix, then each of the columns in  $A$  sum up to 1:

$$\sum_i A_{ij} = 1$$

To show that  $A_\tau$  is also a Markov matrix, we need to show that each of the columns in  $A_\tau$  sum up to 1:

$$\begin{aligned} \sum_i (A_\tau)_{ij} &= \sum_i \left( \tau \cdot A_{ij} + \frac{1-\tau}{n} \cdot (E_n)_{ij} \right) \\ &= \sum_i (\tau \cdot A_{ij}) + \sum_i \left( \frac{1-\tau}{n} \cdot (E_n)_{ij} \right) \\ &= \tau \cdot \sum_i (A_{ij}) + \sum_i \left( \frac{1-\tau}{n} \cdot 1 \right) \text{ (since all components in } E_n \text{ is 1)} \\ &= \tau \cdot 1 + \frac{1-\tau}{n} \cdot n \end{aligned}$$

(since the columns in  $A$  each add up to 1, and the second term is added  $n$  times as  $E_n$  is an  $n \times n$  matrix)

$$= \tau + 1 - \tau = 1$$

Thus  $A_\tau$  is a Markov matrix.

Now we need to prove that  $A_\tau$  has strictly positive entries,  $(A_\tau)_{ij} > 0$ . We will prove by contradiction. Since we know  $A_\tau$  is a Markov matrix, let there be an element  $(A_\tau)_{ij} = 0$ .

For  $(A_\tau)_{ij} = 0$ , it must be the case that

$$(A_\tau)_{ij} = \tau A_{ij} + \left( \frac{1-\tau}{n} \cdot E_n \right)_{ij} = 0$$

We can split it up into two equations:

$$\tau A_{ij} = 0, \quad \left( \frac{1-\tau}{n} \cdot E_n \right)_{ij} = 0$$

Since every element in  $E_n$  consists of 1s, to get a 0-element, we must have  $\tau = 1$  such that  $(1-1)/n = 0$ . This is a contradiction, since  $0 < \tau < 1$ . Therefore  $A_\tau$  is a Markov matrix, and  $(A_\tau)_{ij} > 0$  for all  $i, j$ .  $\square$

## Problem 5

For a Markov matrix  $A$  with strictly positive entries, we need to prove that the transposed matrix  $A^T$  has exactly one eigenvalue  $\lambda$  with  $|\lambda| = 1$ , namely  $\lambda = 1$ , and that the corresponding eigenspace is  $\mathcal{E}_1 = \text{span}\{\mathbf{e}\}$ , where  $\mathbf{e} = [1, 1, \dots, 1]^T$ . Furthermore that any other eigenvalue has modulus less than 1.

### Proof

From Problem 3 we know that  $A$  and  $A^T$  have the same eigenvalues and that  $\text{rad}_A = 1$ . Therefore  $\text{rad}_{A^T} = 1$ , which means that for  $\lambda$  an eigenvalue for  $A^T$  then  $|\lambda| \leq 1$ . We have also shown that  $A$  and  $A^T$  does in fact have an eigenvalue  $\lambda = 1$  as shown in Equation (3.2)

Suppose  $\lambda$  is an eigenvalue of  $A^T$  with  $|\lambda| = 1$  and  $\mathbf{v}$  is a corresponding eigenvector. Let  $k$  denote an index s.t.  $\|\mathbf{v}\|_\infty = |v_k|$  and noting  $(A^T)_{kj} = (A)_{jk} = a_{jk}$ . Now we can write:

$$\|v\|_\infty = |v_k| = |\lambda| \cdot |v_k| = |\lambda \cdot v_k| = \left| \sum_{j=1}^n a_{jk} \cdot v_j \right| \quad (3.3)$$

(Since  $\mathbf{v}$  is an eigenvector for  $A^T$ :  $A^T \mathbf{v} = \lambda \mathbf{v} \rightarrow |A^T \mathbf{v}| = |\lambda \mathbf{v}|$ . Then by only looking at the  $k$ 'th row in  $A^T$  we get:  $|\sum_{j=1}^n (A^T)_{kj} \cdot v_j| = |\lambda v_k| \rightarrow |\sum_{j=1}^n a_{jk} \cdot v_j| = |\lambda v_k|$ .)

$$\leq \sum_{j=1}^n |a_{jk} \cdot v_j| \quad (\text{By the triangle inequality}) \quad (3.4)$$

$$\begin{aligned} &= \sum_{j=1}^n a_{jk} \cdot |v_j| \quad (\text{All components in } A \text{ has positive real values}) \\ &\leq \sum_{j=1}^n a_{jk} \cdot \|v\|_\infty \quad (\text{Since } |v_j| \leq \|v\|_\infty \text{ for all } j = 1, 2, \dots) \\ &= \|v\|_\infty \end{aligned} \quad (3.5)$$

$\|v\|_\infty$  does not depend on  $j$  now, and can be put outside the summation:  $= \|v\|_\infty \cdot \sum_{j=1}^n a_{jk}$ . This summation is a column in the Markov matrix  $A$  and add to 1.

The beginning and end are the same ( $\|v\|_\infty$ ), therefore the inequalities will have to be equalities.

By using the first inequality:  $|\sum_{j=1}^n a_{jk} \cdot v_j| = \sum_{j=1}^n |a_{jk} \cdot v_j|$  and part 2 of the triangle inequality, we now know that  $a_{jk}v_j$  for  $j = 1, 2, \dots$  are non-negative real multiples of each other. By also using the second inequality:  $\sum_{j=1}^n a_{jk} \cdot |v_j| = \sum_{j=1}^n a_{jk} \cdot \|v\|_\infty$  and  $A$  has strictly positive entries, we can conclude that  $|v_j| = \|v\|_\infty = |v_k|$  for all  $j = 1, 2, \dots$  and we must have  $v = Ce$ , where  $|C| = \|v\|_\infty = |v_k|$ .

Thus the corresponding eigenvector must be a multiple of  $\mathbf{e}$ , and there is only one eigenvalue for which  $|\lambda| = 1$  as the corresponding eigenvectors for distinct eigenvalues of the same matrix will be linearly independent.

The corresponding eigenspace for the eigenvalue  $\lambda = 1$  is thus  $\mathcal{E}_1 = \text{span}\{\mathbf{e}\}$ , where  $\mathbf{e} = [1, 1, \dots, 1]^T$ .  $\square$

## Problem 6

**Theorem 3.4** Let  $n \geq 2$  and  $A \in \mathbb{R}^{n \times n}$  be a Markov matrix with strictly positive entries. Then there is a unique vector  $x \in \mathbb{R}^n$  such that  $Ax = x$ , all entries of  $x$  are strictly positive and sum to 1, and:

$$\lim_{k \rightarrow \infty} (A)^k = xe^t = [x, x, \dots, x]$$

If  $A$  is Markov matrix with strictly positive entries, the unique probability vector  $x$  satisfying  $Ax = x$  is called the stationary distribution for  $A$ .

**Explain the relevance of Theorem 3.4 to the PageRank algorithm.**

In programming task 3, we managed to compute the PageRank by taking the power of the matrix  $M_d$  (which we showed to be equivalent to the recursive model in Problem 1), and then taking the first column of the matrix, to get the PageRank vector  $x$ .

Theorem 3.4 is relevant for computing the PageRank algorithm, because it shows that we can find the ranking of a site by taking the power  $(M_d)^k$  as  $k \rightarrow \infty$ . It will cause it to converge to a matrix with the PageRank vector  $x$  in its columns. It also shows that there is a unique ranking for the web. This method causes our simulation to converge faster than computing all the eigenvectors of the matrix and simulating it with the random surfer model. This is crucial, since, for large networks of pages, we want to be able to compute the PageRank fast and efficient.

We can test the runtime by implementing a time function in our three different scenarios. We will generate a random web by the function `make_web`, and count the seconds that elapses for the functions: `random_surf`, `recursive_pagerank`, `eigenvector_pagerank`, and `matrix_pagerank` (Thm. 3.4.):

```
>>> web = make_web(100, 20)
>>> ranking1 = random_surf(web, 100000)
>>> ranking2, iterations = recursive_pagerank(web, 0.0000001)
>>> ranking3 = eigenvector_pagerank(web)
>>> ranking4 = matrix_pagerank(web, 20)

Total runtime of random_surf is 3.743 seconds
Total runtime of recursive_pagerank is 0.0216 seconds
Total runtime of eigenvector_pagerank is 0.009 seconds
Total runtime of matrix_pagerank is 0.001 seconds
```

This is merely an example of 100 pages where each page links to between 0-20 (If you want to see comparisons between our implementation on larger networks, see Comparing our implementations) This shows the relevance of the theorem in the algorithm.

## Problem 7

Prove Theorem 3.4 for the case that  $A$  has  $n$  distinct eigenvalues.

### Proof

Since  $A$  has  $n$  distinct eigenvalues, we know that it is diagonalizable and we can write  $A = VDV^{-1}$  where  $D$  is a diagonal matrix of our unique eigenvalues, and  $V$  is a matrix of the corresponding eigenvectors. We can write  $A^k = VD^kV^{-1}$ . Since

$$D^k = \begin{bmatrix} \lambda_1^k & & 0 \\ & \lambda_2^k & \\ & & \ddots \\ 0 & & & \lambda_n^k \end{bmatrix}$$

and  $|\lambda| \leq 1$ , then  $D^k$  must converge to a matrix consisting of 1's and 0's. Therefore  $A^k$  must also converge and  $\lim_{k \rightarrow \infty} (A)^k = W$  exists.

To explain why every column of  $W$  must be an eigenvector of  $A$  with an eigenvalue 1, we denote the columns of  $W$  by  $w_1, w_2, \dots, w_n$ , and consider the product  $AW = A \cdot [w_1, w_2, \dots, w_n] = [w_1, w_2, \dots, w_n]$ .

By looking at an arbitrary column in  $W$ ,  $w_i$ , the matrix vector product will be:  $Aw_i = w_i$ .

This shows that the  $i$ 'th column in  $W$  is an eigenvector of  $A$  with eigenvalue 1, and this applies to all the columns of  $W$ .

We know that  $A$  has an eigenvalue of 1 and that there is a corresponding eigenvector  $w_i$ . In Problem 5 we have shown that the eigenspace  $\mathcal{E}_1 = \text{span}\{\mathbf{e}\}$  for  $A^T$ . In this problem we then have  $\mathcal{E}_1 = \text{span}\{w_i\}$ .

This means that the corresponding eigenvectors of  $A$  to the eigenvalue 1 is multiple of each other. One of these eigenvectors will sum to 1, which will be a probability vector denoted by  $\mathbf{x}$ .

By using the probability vector  $\mathbf{x}$  as a basis, all the columns of  $W$  will then be a scalar multiple of  $\mathbf{x}$ .

From Problem 2 we know that  $A^k$  is a Markov matrix, therefore  $W$  must be a Markov matrix. And furthermore  $W$  must be a regular Markov matrix, as  $A$  is a regular Markov matrix, and every entry of the product of two regular Markov matrices can be described by:  $(AB)_{ij} = \sum_k (A_{ik}B_{kj})$ , where all terms in the sum are strictly positive, thus the sum is strictly positive.

Every column in a Markov matrix sum to 1, therefore every column in  $W$ ,  $w_i$ , is a probability vector with strictly positive entries. The only scalar multiple of  $\mathbf{x}$ , which sum to 1, is  $\mathbf{x}$  itself. All the columns in  $W$  must therefore be identical to  $\mathbf{x}$ . Therefore  $W$  takes the form:  $W = [w_1, w_2, \dots, w_n] = [\mathbf{x}, \mathbf{x}, \dots, \mathbf{x}]$ .  $\square$

## Problem 8

Prove Theorem 3.4 for all cases, using the Jordan Canonical Form.

### Proof

Since  $A$  is an  $n \times n$  matrix,  $A$  has  $n$  eigenvalues if counted with multiplicity. Since  $A$  is a square matrix we can write  $A = PJP^{-1}$ , for some matrix  $P$ , where  $J = \begin{bmatrix} J_1 & & 0 \\ & J_2 & \\ & & \ddots \\ 0 & & & J_r \end{bmatrix}$  and each block  $J_m$  is a square matrix that has an eigenvalue of  $A$  down the diagonal, 1's immediately above the diagonal and 0's everywhere

else:  $J_m = \begin{bmatrix} \lambda & 1 & & 0 \\ & \lambda & 1 & \\ & & \ddots & 1 \\ 0 & & & \lambda \end{bmatrix}$ . We can write  $A^k = PJ^kP^{-1}$ , where  $J^k = \begin{bmatrix} J_1^k & & 0 \\ & J_2^k & \\ & & \ddots \\ 0 & & & J_r^k \end{bmatrix}$  and each block with a size

$$l \times l: J_m^k = \begin{bmatrix} \binom{k}{0}\lambda^k & \binom{k}{1}\lambda^{k-1} & \binom{k}{2}\lambda^{k-2} & \dots & \binom{k}{l-1}\lambda^{k-(l-1)} \\ & \binom{k}{0}\lambda^k & \binom{k}{1}\lambda^{k-1} & \ddots & \vdots \\ & & \binom{k}{0}\lambda^k & \ddots & \binom{k}{2}\lambda^{k-2} \\ & & & \ddots & \binom{k}{1}\lambda^{k-1} \\ 0 & & & & \binom{k}{0}\lambda^k \end{bmatrix}$$

We want to show that  $W$  exist for  $A$  using the Jordan normal form. For  $W$  to exist the matrix  $A^k$  has to be converging. We split it up into two cases.

**Case 1:**  $|\lambda| < 1$

First we look at all the eigenvalues that have absolute values smaller than 1. The non-zero values in the Jordan blocks with these eigenvalues will have values:  $\binom{k}{0}\lambda^k$  which converges to 0,  $\binom{k}{1}\lambda^{k-1}$  which also converges to 0, ...,  $\binom{k}{l-1}\lambda^{k-(l-1)}$ , which converges to 0 for  $k > l$ . Thus the Jordan blocks with eigenvalues that have absolute values smaller than 1 will all converge to 0.

**Case 2:**  $|\lambda| = 1$

We will prove by contradiction. Suppose that  $am(1) > 1$ . We know that  $gm(1) = 1$  from Problem 5, then there will be a Jordan block, where at least one value is  $k\lambda^{k-1}$  when  $k \rightarrow \infty$ , that will go to infinity. This cannot be true since it is a Markov matrix, which by definition has columns that sum to 1. A column with a value that go to infinity won't have this property of a Markov matrix.

Thus  $am(1) = 1$  which means that it converges to 1, as the Jordan block only consists of the non-zero value  $\lambda^k$ . Therefore  $J^k$  converges so  $A^k$  must also converge, and the limit exists.

Now we have the same case as in Problem 7, and the rest of the proof follow from there.  $\square$

## 3 Further questions

### Damping factor

The damping factor has an effect on the range of values the algorithm gives. If the damping factor is very low, it is more likely that a page will be chosen at random, than it will be reached via a link. The standard value used in our project is 0.85.

Example: Comparing PageRank calculation with damping factor 0.2, 0.85 and 0.99:

```
>>> web={0: {1, 7}, 1: {3, 6}, 2: {0, 1, 3}, 3: set(), 4: set(), 5: {3, 4, 6}, 6: {0}, 7: {4,
↪ 6}}
>>> ranking1, iterations = recursive_pagerank(web,0.00001, 200, 0.2)
PageRank: {0: 0.1413, 1: 0.1277, 2: 0.1065, 3: 0.1334, 4: 0.1256, 5: 0.1065, 6:
↪ 0.1384, 7: 0.1206}
Iterations: 3
>>> ranking2, iterations = recursive_pagerank(web,0.00001, 200, 0.85)
PageRank: {0: 0.2077, 1: 0.146, 2: 0.045, 3: 0.1325, 4: 0.1144, 5: 0.045, 6:
↪ 0.1764, 7: 0.1332}
Iterations: 21
>>> ranking3, iterations3 = recursive_pagerank(web,0.00001, 200, 0.99)
PageRank: {0: 0.2247, 1: 0.1519, 2: 0.0305, 3: 0.1259, 4: 0.1108, 5: 0.0305, 6:
↪ 0.186, 7: 0.1418}
Iterations: 198
```

The PageRank has a more even distribution when the damping factor is set to 0.2 compared to the ones set to 0.85 and 0.99. The PageRank of a page will depend more on the amount of inbound links and the structure of



links in the web, if the damping factor is closer to 1. The algorithm also requires more iterations to calculate if the damping factor is closer to 1, but the algorithm captures the link structure better for higher values of  $d$ .

For  $d = 1$  the algorithm doesn't necessarily give a unique pagerank.

Example: For pages 0,1,2,3,4, where 1 and 4 links to each other, 0 and 2 links to each other, 0 links to 3, and 3 links to 2. The pagerank can for this example be computed to two different outcomes:

```
>>> web = {0: {2,3}, 1: {4}, 2: {0}, 3: {2}, 4: {1}}
>>> ranking1 = random_surf(web, 100000, 1)
Pagerank: {0: 0.40, 1: 0.0, 2: 0.40, 3: 0.20, 4: 0.0}
>>> ranking2 = random_surf(web, 100000, 1)
Pagerank: {0: 0.0, 1: 0.50, 2: 0.0, 3: 0.0, 4: 0.50}
```

## Extending our algorithm

We want to create a function that updates the web and page ranking of a site simultaneously. That means we want to be able to add pages and links, as well as delete them (see appendix 4.4 for the function). The parameters of the function is the web, a dictionary of new pages and links you wish to add, a dictionary of links you want to remove from the web, and a list of pages you want to remove.

We then start by fetching the initial pagerank through the `matrix_pagerank` implementation, and then go through each case if there are elements in the respective lists. When we add the new dictionary, if a page is not already in the web, we also initialize a pagerank of  $1/(N + 1)$ . Once we have added/removed from the web, we use the recursive implementation to loop through and update the rankings of each page based on the previous ranking. We also make sure to save the pagerank before going through each page, and push it into the pagerank once we have gone through each page.

## Comparing our implementations

During the comparisons, our implementations `recursive_pagerank`, `eigenvector_pagerank` and `matrix_pagerank` have the same precision when rounding off at 4 decimal places. Because of the nature of our `random_surf` simulation, we can expect it to varie slightly from the rest. We can calculate the margin of error by comparing each implementation against `eigenvector_pagerank` as the true evaluation (as it has no parameters other than web and damping), and summing up the absolute value of the difference. Looking at the error on a web of 100 pages we see that

```
Error: random_surf: 0.0242, recursive_pagerank: 0.0000, matrix_pagerank: 0.0000
```

We can see that the difference is minimal and within an acceptable margin of error of  $\leq 0.025$ . With this in mind, let us compare our four implementations on larger networks. We will use the time module, and log `time.time()` from the beginning and end of the functions, and return the difference. Let us first define the high precision parameters:

```
>>> random_surf(web, 100000)
>>> recursive_pagerank(web, 0.0000001)
>>> eigenvector_pagerank(web)
>>> matrix_pagerank(web, 20)
```

Similarly, here are our low precision parameters:

```

>>> random_surf(web, 10000)
>>> recursive_pagerank(web, 0.00001)
>>> eigenvector_pagerank(web)
>>> matrix_pagerank(web, 10)

```

Now, let us first take a look at our web of 10.000 pages with links between 0-2000 for each:

Method	Runtime (High)	Error (High)	Runtime (Low)	Error (Low)	Relation
random surf	186.98s	0.2507	18.49s	0.7416	10.11x
recursive	39.12s	0.0001	19.34s	0.0001	2.02x
eigenvectors	135.87s	-	135.42s	-	1.00x
matrix	35.46s	0.0000	30.29s	0.0000	1.17x

Table 1: Runtimes for 10.000 pages varying in precision (high vs low).

The relative time between **random\_surf** in high vs low precision is significant, but we also see a large increase in error. This is because of a scaling problem we will touch more on below. Notice the matrix method is faster than the recursive method at high precision, but on low precision the recursive method becomes faster. We also see that the relation of recursion is  $2.02x$  faster, but still maintaining the same error in both high and low precision. This means we could have probably lowered the precision slightly to reduce the overall runtime of the method whilst maintaining accurate results.

Below is a comparison for a web of 20.000 pages and between 0-4000 links each:

Method	Runtime (High)	Error (High)	Runtime (Low)	Error (Low)	Relation
random surf	384.68s	0.35042	37.68s	1.2180	10.21x
recursive	155.00s	0.0000	95.20s	0.0000	1.62x
eigenvectors	988.65s	-	964.14s	-	1.02x
matrix	275.41s	0.0000	214.73s	0.0000	1.28x

Table 2: Runtimes for 20.000 pages varying in precision (high vs low).

Notice a scaling problem we have for the random surf implementation. The margin of error is significantly higher when run on a larger network, considering the same amount of iterations. We also observe that, upon lowering the precision (less iterations), the error margin becomes so high that the ranking is completely wrong. This is because the number of iterations needed for the method to converge depends on the amount of pages we have. If there are 100.000 iterations on a web of 100 pages it may be sufficient with a 0.0242 error margin. This does not scale for a web of 10.000 or 20.000 pages, as we have the same amount of iterations, but the amount of pages are much larger, so we aren't able to go through each page enough times to calculate the pagerank. Therefore it is evident that the random surf model would not work well on large networks, as the runtime needs to be exceedingly higher for larger webs, otherwise we lose precision. This is very inefficient.

We also see that the recursive implementation becomes faster than the matrix implementation on a larger network. This likely happens because the recursive method only updates the pages and links necessary, but the matrix method has to multiply incredibly large matrices together which might take a lot of memory on larger networks. We can see in Table 2, when we lower the precision, that the recursive method runs  $1.62x$  faster whereas the matrix method only runs  $1.28x$  faster than at high precision, but with the exact same error margins.

The matrix method is indeed faster than the eigenvector method as expected from our solutions, but the overall best performance is from the recursive model on larger networks as seen in Table 2.

## 4 Appendix

### 4.1 pagerank\_script1.py

```
1
2 import numpy as np
3 ##### Helper Functions #####
4 def make_web(n,k,kmin=0):
5     assert(k<n)
6     keys = np.array(range(n))
7     web = dict()
8     for j in keys:
9         numlinks = np.random.choice(range(kmin,k+1))
10        web[j] = set(np.random.choice(keys[keys!=j],numlinks,replace=False))
11    return web
12
13 ##### Programming Task 1 #####
14 def surf_step(web, page, d=0.85):
15     """
16     Return a probability distribution over which page to visit next,
17     given a current page.
18     Input: -web is a dictionary of webpages and links
19            - page is a key of the dictionary, the page from which to
20              compute the next step
21            - d is the damping factor
22
23     - If page does not link to any page, then choose any page in web at random
24     - Otherwise:
25         With probability `d`, choose a page at random linked to by `page`,
26         With probability `1 - d`, choose a page at random from all pages in the web.
27     Output: is a dictionary with the same keys as web, and the
28            value for each key is the probability of choosing that page next.
29            (The sum of all these should be 1)
30     """
31    distribution = dict() # the distribution dictionary
32    N = len(web.keys())
33
34    # initialize that every key in web needs a place in distribution
35    for key, val in web.items():
36        distribution[key] = 0
37
38    # If page does not link to any page, then choose any page in web at random
39    if len(web[page]) < 1:
40        for key, val in web.items():
41            distribution[key] += 1 / N
42
43    # Otherwise
44    if len(web[page]) > 0:
45        # With probability `d`, choose a page at random linked to by `page`,
46        for link in web[page]:
47            distribution[link] += d / (len(web[page]))
48        # With probability `1 - d`, choose a page at random from all pages in the web.
49        for key, val in web.items():
50            distribution[key] += (1-d) / N
```

```

51
52     return distribution
53
54
55 def random_surf(web,n,d=0.85):
56     """
57     Return PageRank values for each page by sampling `n` pages
58     according to surf_step.
59     Input: web is a dictionary of webpages and links,
60           n is an integer, the number of steps in the simulation
61           d is the damping factor,
62
63     Returns a dictionary with the same keys as web (the pages), and
64     the value for key k is the page rank of page k. The sum of all PageRank values
65     should be 1.
66     """
67     ranking = dict() # the ranking for each page
68     all_pages = list(web.keys())
69
70     # initialize that every key in web needs a place in distribution
71     for key, val in web.items():
72         ranking[key] = 0
73
74     # Choose the first page randomly
75     p = np.random.choice(all_pages)
76     # Loop through the simulation n times
77     for _ in range(n):
78         # Add 1 each time we visit a site
79         ranking[p] += 1
80         # Get our probability distribution of the page p
81         probdist = surf_step(web, p, d)
82         # Make a list of the probabilities of visiting the next site, given we're at page
83         ↪ p
84         probs = []
85         for key, val in probdist.items():
86             probs.append(val)
87         # Choose a new page based on that probability
88         p = np.random.choice(all_pages, p=probs)
89
90     # Divide by n to get the pagerank
91     for key, val in ranking.items():
92         ranking[key] /= n
93
94     return ranking

```

## 4.2 pagerank\_script2.py

```

1 ##### Programming Task 2
2
3 def rank_update(web,pageranks,page,d):
4     """
5     Updates the value of the pagerank for page based on the formula
6      $PR(p) = (1-d)/N + d \cdot \sum_j (PR(q)/OB(q))$ 
7     where the sum is over all pages q that link to page, PR(q) is the current

```

```

8     pagerank of page q (from "pageranks") and OB(q) is the number of pages
9     outbound from page q. Sinks are treated as linking to all pages in web.
10
11     Input: web and pageranks are dictionaries as in the output of "make_web" and
12     "random_surf", page is the key to the page whose rank we wish to update,
13     and d is the damping factor.
14     Output: The (mutable) dictionary "pageranks" is updated according to the above
15     ↪ formula,
16     and this function returns a float "increment", the (absolute) difference
17     between the previous value and the updated value of PR(p).
18     '''
19     increment = 0
20     # Pagerank of initial page
21     PR_old = pageranks[page]
22     inbound_p = set()
23     N = len(web.keys())
24
25     # add elements to inbound_p
26     # loop through and check if each links to p
27     for key, val in web.items():
28         # if the page is in the links of key, add it to inbound
29         if page in val:
30             inbound_p.add(key)
31
32     # sink, if page does not link to any pages, then add to inbound_p
33     elif not web[key]:
34         inbound_p.add(key)
35
36     # split the equation up between first and second part
37     del1 = (1 - d) * (1 / N)
38     del2 = 0
39
40     for q in inbound_p:
41         # pagerank and outbound of q
42         PR_q = pageranks[q]
43         OB_q = len(web[q])
44
45         # sink for sinks q, outbound is N
46         if OB_q < 1:
47             OB_q = N
48
49         del2 += (PR_q) / (OB_q)
50
51     # modify the existing pageranks dictionary
52     pageranks[page] = del1 + (d * del2)
53     increment = abs(PR_old - (del1 + (d * del2)))
54
55     return increment
56
57 def recursive_pagerank(web, stopvalue, max_iterations=200, d=0.85):
58     """
59     Implements the recursive version of the PageRank algorithm by first creating a
60     pagerank of 1/N to all pages (where N is the total number of pages)
61     then applying "rank_update" repeatedly until either of two stopping conditions is

```

```

61     reached:
62     stopping condition 1: the maximum change from step n to step (n+1) over all pageranks
63     is less than stopvalue,
64     Stopping condition 2: the number of iterations has reached "max_iterations"
65
66     Input: web is a dictionary as in the output of "make_web", d is the damping constant,
67     stop value is a positive float, max_iterations is a positive integer
68     """
69     pageranks = dict()
70     iteration = 0
71     N = len(web.keys())
72
73     for key, val in web.items():
74         pageranks[key] = 1 / N
75
76     # Loop through until max iterations
77     for n in range(max_iterations):
78         # Log the change diffs
79         max_change = 0
80
81         # Go through each page and update the rank
82         for page in web.keys():
83             change = rank_update(web, pageranks, page, d)
84             if change > max_change:
85                 max_change = change
86
87         # Break out if it reaches below stopvalue
88         if(max_change < stopvalue):
89             break
90
91         iteration += 1
92
93     return pageranks, iteration

```

### 4.3 pagerank\_script3.py

```

1  ##### Programming Task 3 #####
2
3  def modified_link_matrix(web, pagelist, d=0.85):
4      """
5      Create a modified link matrix from web
6      Input: web is a dictionary whose keys are contained in the list pagelist,
7      pagelist is just a list of the keys (to give the keys and ordering)
8      d is the damping factor
9      Output:  $d \cdot A^T + (1-d) \cdot E/N$ 
10     where A is an  $N \times N$  numpy array, for which
11     row j has non-zero entries only in the columns corresponding
12     to pagelist to which page j links. In a given row, all non-zero entries have
13      $\hookrightarrow$  the
14     same value, and these values sum to 1. Special case: if page j does not link
15     to any page, then all entries in row j are given the same value (1/N).
16     E is  $\text{np.ones}([N, N])$ 
17     """
18     N = len(web.keys())

```

```

18     E = np.ones([N,N])
19     A = np.zeros([N,N])
20
21     # Modified adjacency matrix
22     for j in range(N):          # cols
23         for i in range(N):      # rows
24             # if p_j links to p_i
25             if pagelist[i] in web[pagelist[j]]:
26                 N_j = len(web[pagelist[j]])
27                 A[j,i] = 1 / N_j
28
29             # if p_j is a sink
30             if len(web[pagelist[j]]) < 1:
31                 A[j, i] = 1 / N
32
33     return d * A.T + (1 - d) * (E / N)
34
35 def eigenvector_pagerank(web,d=0.85):
36     """
37     Returns the pagerank of web as the eigenvector of the modified link matrix
38     Input: web is a dictionary of web pages and lines.
39           d is a positive float, the damping constant
40     Output: A dictionary with the same keys as web, and the values of the pageranks of
41             ↪ the keys
42     """
43     ranking = dict()
44
45     # Init dict with same keys as web
46     for key, val in web.items():
47         ranking[key] = 0
48
49     pages = list(web.keys())
50     # Create our adjacency matrix M from function before
51     M = modified_link_matrix(web,pages, d)
52
53     # Use np eigenvector function to find our eigenvals and vects
54     lamda, V = np.linalg.eig(M)
55     # Choose the first vector
56     V1 = V[:,0:1]      # the eigenvalue method
57     # normalize the eigenvector so that it's components sum to 1
58     V1 = np.real(V1)
59     ranking3 = V1/(np.sum(V1))  # ranking from eigenvector 1
60
61     # Loop through the vector and input them into the ranking dict.
62     for k in range(len(pages)):
63         ranking[k] = float(ranking3[k])
64
65     return ranking
66
67 def matrix_pagerank(web,power,d=0.85):
68     """
69     Returns the pagerank as the first column of the power'th power of the modified link
70     ↪ matrix

```

```

70     Input: web is a dictionary of web pages and lines.
71         d is a positive float, the damping constant
72     Output: A dictionary with the same keys as web, and the values the pageranks of the
73         → keys
74     """
75     ranking = dict()
76
77     for key, val in web.items():
78         ranking[key] = 0
79
80     # Create pages and M adjacency matrix
81     pages = list(web.keys())
82     M = modified_link_matrix(web, pages, d)
83
84     # Use method (b) by computing powerth power of M
85     K = np.linalg.matrix_power(M, power)
86     # Choose the first col
87     new_rank = K[:, 0]
88
89     # Loop through the vector and input it into the ranking dict.
90     for k in range(len(pages)):
91         ranking[k] = float(new_rank[k])
92
93     return ranking

```

#### 4.4 update\_web function (inside pagerank\_script3.py)

```

1  def update_web(web, webadd={}, removelinks={}, removepages=[], d=0.85):
2      pageranks = matrix_pagerank(web, 20)
3      N = len(web.keys())
4
5      # Add webadd to web
6      # Initialize a pagerank if not already in web
7      if len(webadd) > 0:
8          for key, val in webadd.items():
9              if key not in web:
10                 web[key] = set()
11                 pageranks[key] = 1 / (N + 1)
12
13             for num in val:
14                 web[key].add(num)
15
16     # Go through removelinks dict
17     # and remove the links associated with the key
18     if len(removelinks) > 0:
19         for key, val in removelinks.items():
20             for num in val:
21                 web[key].remove(num)
22
23     # Go through removepages list and remove each page
24     # and all its inbound and outbound links
25     if len(removepages) > 0:
26         for page in removepages:
27             # Remove page from web and pageranks

```



```

28         if page in web:
29             del web[page]
30         if page in pageranks:
31             del pageranks[page]
32
33         # Go through every single link in the web and remove page
34         for key, val in web.items():
35             if page in val:
36                 web[key].remove(page)
37
38         # Now we loop through (from recursive model)
39         # and update rank based on previous rankings
40         for _ in range(200):
41             # Log the change diffs to break out if it reaches below stopvalue
42             max_change = 0
43             # Remember to save the ranks whilst going through each page
44             new_ranks = pageranks.copy()
45             # Go through each page and update the rank
46             for page in web.keys():
47                 change = rank_update(web, new_ranks, page, d)
48                 if (change > max_change):
49                     max_change = change
50
51             pageranks = new_ranks
52             # Break out if it reaches below stopvalue
53             if(max_change < 0.0000001):
54                 break

```