

Question 1

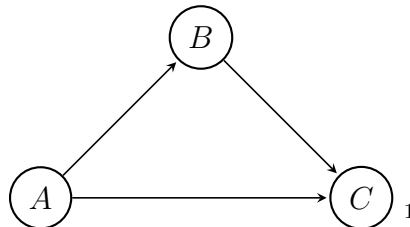
LJ Brown

May 1, 2018

1 Question

1. (25/25 points) Consider the following web pages and the set of web pages that they link to: Page A points to pages B and C. Page B points to page C. All other pages have no outgoing links. Apply the PageRank algorithm to this subgraph of pages. Assume $\alpha = 0.15$. Simulate the algorithm for three iterations.

Disconnected Webgraph



2 PageRank Algorithm

The "pageRank algorithm" attempts to gauge the importance of each website in a webgraph by ranking each node according to the number and "quality" of its inlinks.

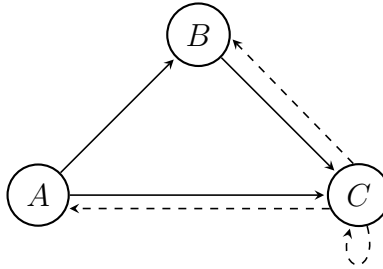
Imagine a web surfer who begins at some webpage or node in the "Disconnected Webgraph" above, and performs a random walk by clicking on different links. At every node the surfer picks a random outlink or edge to

¹Webpage C, in the "Disconnected Webgraph" is called a "dangling node".

follow, each outlink from the same webpage having an equal chance being selected. In the "Disconnected Webgraph" above, the surfer will eventually become trapped at node C .

The "pageRank algorithm" removes the "dangling node" trap by adding connections to nodes with no outlinks to every single other node in the graph. By connecting the "dangling nodes" to every other node, not only would the surfer be able to click forever, but the surfer would continue to revisit every single node in the graph during this infinite walk.

Strongly Connected Webgraph

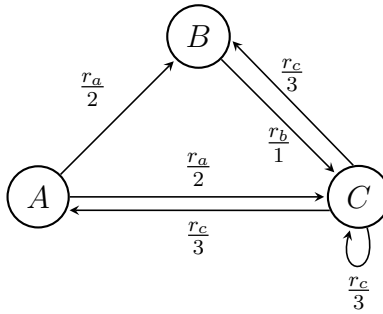


However, although the surfer would eventually revisit every single website, the probability that any particular website is visited next may be different. This probability is the websites "rank".

A nodes rank, r_i , is defined by summing the inlink values of all nodes pointing to it,

$$r_i = \sum_{k \rightarrow i} \frac{r_k}{o_k}$$

Where o_k is the number of outlinks for node k .



Ignoring α or the damping factor in the question above, the definition above for a websites rank, r_i , gives the following equations,

$$r_a = \frac{r_c}{3} \quad (1)$$

$$r_b = \frac{r_a}{2} + \frac{r_c}{3} \quad (2)$$

$$r_c = \frac{r_a}{2} + \frac{r_b}{1} + \frac{r_c}{3} \quad (3)$$

These equations have no unique solution (3 equations with 3 unknowns). But since we are looking for the probability of a websites selection in a sense we can add the additional constraint that all the ranks sum to one in order to solve the system,

$$r_a + r_b + r_c = 1 \quad (4)$$

Writing this system in matrix form, performing Gauss-Jordan elimination and solving for website ranks r_a , r_b , and r_c ,

$$\begin{aligned} r_a + 0 - \frac{r_c}{3} &= 0 \\ \frac{r_a}{2} - r_b + \frac{r_c}{3} &= 0 \\ \frac{r_a}{2} + r_b - \frac{2r_c}{3} &= 0 \\ r_a + r_b + r_c &= 1 \end{aligned} \rightarrow \begin{bmatrix} 1 & 0 & \frac{-1}{3} \\ \frac{1}{2} & -1 & \frac{1}{3} \\ \frac{1}{2} & 1 & \frac{-2}{3} \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} r_a \\ r_b \\ r_c \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & \frac{-1}{3} & \left| 0 \right. \\ \frac{1}{2} & -1 & \frac{1}{3} & \left| 0 \right. \\ \frac{1}{2} & 1 & \frac{-2}{3} & \left| 0 \right. \\ 1 & 1 & 1 & \left| 1 \right. \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & \frac{-1}{3} & \left| 0 \right. \\ \frac{1}{2} & -1 & \frac{1}{3} & \left| 0 \right. \\ \frac{1}{2} & 1 & \frac{-2}{3} & \left| 0 \right. \\ 1 & 1 & 1 & \left| 1 \right. \end{bmatrix} \rightarrow \begin{bmatrix} 1 & 0 & 0 & \left| \frac{2}{11} \right. \\ 0 & 1 & 0 & \left| \frac{3}{11} \right. \\ 0 & 0 & 1 & \left| \frac{6}{11} \right. \\ 0 & 0 & 0 & \left| 0 \right. \end{bmatrix}$$

$$r_a = \frac{2}{11}, r_b = \frac{3}{11}, r_c = \frac{6}{11}$$

If there was no damping factor ($\alpha = 0$) then these solutions for r_a , r_b , and r_c would be the page ranks corresponding to the "Strongly Connected Webgraph".

3 Link Matrix

Another way to represent this problem involves a "Link Matrix" or adjacency matrix representation of the "Disconnected Webgraph". The Link Matrix, L_0 , is an $n \times n$ Matrix (n is the number of nodes) where a 1 is placed at element l_{ij} if there is an edge from node i to node j ,

$$L_0 = \begin{matrix} & \begin{matrix} a & b & c \end{matrix} \\ \begin{pmatrix} 0 & 1 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{pmatrix} & \begin{matrix} a \\ b \\ c \end{matrix} \end{matrix}$$

The elements of L_0 containing 1's represent a path of length 1 from from node i to node j . This matrix has the interesting property that by raising L to the k th power we find all paths of length k within this graph. For example to find paths of length 2,

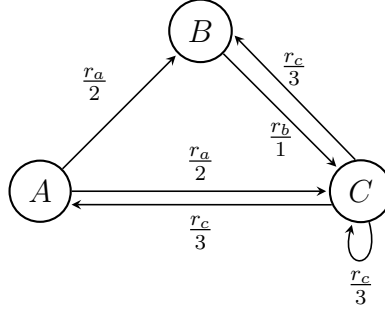
$$L_0^2 = \begin{pmatrix} 0 & 1 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} 0 & 1 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{pmatrix} = \begin{matrix} \begin{matrix} a & b & c \end{matrix} \\ \begin{pmatrix} 0 & 0 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \begin{matrix} a \\ b \\ c \end{matrix} \end{matrix}$$

Meaning the only path of length 2 in the "Disconnected Webgraph" is from node a to node c . If you instead raise L_0 to the 3rd power instead you will find a matrix of all zeros, implying that there are no paths of length 3. If you create another link matrix, L , from the "Strongly Connected Webgraph" and perform the same operation you will notice that the resulting matrix is positive for any value of $k \geq 1$.

$$L^2 = \begin{pmatrix} 0 & 1 & 1 \\ 0 & 0 & 1 \\ 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} 0 & 1 & 1 \\ 0 & 0 & 1 \\ 1 & 1 & 1 \end{pmatrix} = \begin{matrix} \begin{matrix} a & b & c \end{matrix} \\ \begin{pmatrix} 1 & 1 & 2 \\ 1 & 1 & 1 \\ 1 & 2 & 3 \end{pmatrix} \begin{matrix} a \\ b \\ c \end{matrix} \end{matrix}$$

This is the reason that the surfer can continue clicking forever in the "Strongly Connected Webgraph".

4 Transition Matrix



The "Transition Matrix", T_0 , or weighted matrix is like the link matrix except takes into account the outlinks, o_i , for a given edge. o_i is the number of non zero entries in a row, i , of the link matrix. And the non zero values of the link matrix are replaced by $\frac{1}{o_i}$ for every element in T_0 .

$$T_0 = \begin{pmatrix} a & b & c \\ 0 & \frac{1}{2} & \frac{1}{2} \\ 0 & 0 & 1 \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{pmatrix} \begin{matrix} a \\ b \\ c \end{matrix}$$

Every row of the "Transition Matrix" sums to 1, making it a "Row Stochastic Matrix". This matrix can be used to model the surfers movement. Given a starting vector, \vec{v}_a , with the surfer starting on webpage a , we can find the probability that the surfer will be on webpage a on his first visit (webpage a),

$$\vec{v}_a T_0^1 = (1 \ 0 \ 0) \begin{pmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ 0 & 0 & 1 \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{pmatrix} = (1 \ 0 \ 0)$$

Or the probability of websites after 2 clicks,

$$\vec{v}_a T_0^2 = (1 \ 0 \ 0) \begin{pmatrix} \frac{1}{6} & \frac{1}{6} & \frac{2}{3} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{9} & \frac{5}{18} & \frac{11}{18} \end{pmatrix} = (\frac{1}{6} \ \frac{1}{6} \ \frac{2}{3})$$

Or finally the website ranks, or,

$$\lim_{k \rightarrow \infty} \vec{v}_a T_0^k = (\frac{2}{11} \ \frac{3}{11} \ \frac{6}{11})$$

Notice the answer comes out to the same one obtained from the system of equations above,

$$r_a = \frac{2}{11}, r_b = \frac{3}{11}, r_c = \frac{6}{11}$$

$$\vec{r} = \left(\frac{2}{11} \quad \frac{3}{11} \quad \frac{6}{11} \right)$$

\vec{r} will be the same regardless of the starting vector (in this case \vec{v}_a). \vec{r} is also be the eigenvector for the eigenvalue, $\lambda = 1$, for the Transition Matrix, T_0 , created from the "Strongly Connected Webgraph's" Link Matrix, L .

$$L = \begin{matrix} & \begin{matrix} a & b & c \end{matrix} \\ \begin{pmatrix} 0 & 1 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{pmatrix} & \begin{matrix} a \\ b \\ c \end{matrix} \end{matrix}, T_0 = \begin{matrix} & \begin{matrix} a & b & c \end{matrix} \\ \begin{pmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ 0 & 0 & 1 \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{pmatrix} & \begin{matrix} a \\ b \\ c \end{matrix} \end{matrix}, \vec{r} = \left(\frac{2}{11} \quad \frac{3}{11} \quad \frac{6}{11} \right)$$

5 Damping Factor

However if a damping factor is added, $\alpha = 0.15$, is included then the flow equations and graph are changed and the the Transition Matrix, $T_{0.15}$, becomes,

$$T_{0.15} = (0.85) \begin{pmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ 0 & 0 & 1 \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{pmatrix} + (0.15) \begin{pmatrix} \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{pmatrix} = \begin{pmatrix} 0.05 & 0.475 & 0.475 \\ 0.05 & 0.05 & 0.09 \\ 0.33 \dots & 0.33 \dots & 0.33 \dots \end{pmatrix}$$

If α were set to 1 then each website would have the same rank. α can be thought of as the probability the surfer jumps to a random webpage.

6 Question 1 Answer

1. (25/25 points) Consider the following web pages and the set of web pages that they link to: Page A points to pages B and C. Page B points to page C. All other pages have no outgoing links. Apply the PageRank algorithm to this subgraph of pages. Assume $\alpha = 0.15$. Simulate the algorithm for three iterations.

If we choose the starting vector, \vec{v} , to be the eigenvector of $T_{0.15}$ corresponding to the eigenvalue, $\lambda = 1$, then the probabilities will not change for each iteration.

$$\vec{v} \approx (0.19757929 \quad 0.28155074 \quad 0.52086996)$$

Then for $k = 3$,

$$\vec{v}T_{0.15}^3 = \vec{v} \approx (0.19757929 \quad 0.28155074 \quad 0.52086996)$$