

# PageRank - Lord of the Rings

Ksenia Lepikhina

Kylee Bennett

Michael Feller

Section: 002

Recitation: 002

Recitation: 001

05/05/2017

## 1 Abstract

The topic we will be investigating in this paper is the analysis of the relationships between characters in The Lord of the Rings by J.R.R. Tolkien using PageRank. We have adapted our idea from a project done by Andrew Beveridge and Jie Shan called the “Network of Thrones”. PageRank was used by Google in order to determine which pages to show first. It ranks each page by how many ”votes” it has (links directing to it). Using this, we can rank characters in a book by how many characters mention, or are in close proximity ( $x$  number of words) to, a certain character. The character that has the most relationships with other important characters will be ranked as the most important character in the book.

## 2 Attribution

Kylee Bennett and Michael Feller performed the cleaning of the data as well as the programming portion of this project using C++ and Python. Ksenia Lepikhina researched the mathematics behind the project. All three contributed to the preparation of this report.

## 3 Introduction

When the internet started to become available to the public in the early 1990s, search engines were created to help locate relevant documents. These early search engines displayed results by determining whether or not text keywords appeared on web pages, and in what quantities. As a result, many garbage web pages were created, filled with garbage keyword text, in order to spoof the search engines and bring in more traffic in order to obtain advertising revenue. As a remedy to this, Sergey Brin and Larry Page developed PageRank.

PageRank is defined as an algorithm used by Google Search to rank pages in their search engine results. Larry Page, a founder Google, was the first to profit off wide-scale use of the algorithm. The algorithm is used to measure the importance of pages by measuring the quantity/ quality of the links to the page. The idea is that the most important page will have more websites linking or referring to it, as opposed to what text appears on the website. PageRank assigns a numerical weight to each link. In our analysis of character importance in Lord of the Rings, we use characters as web pages. We assume that when two characters are mentioned within the vicinity of one another, each one gives the other a "vote" of importance. Though this algorithm does take into account a scenario where the nodes (websites, characters) are not connected, we will be investigating the most simple case where each node is connected to at least one other node since in the trilogy, each character interacts with at least one other.

## 4 Mathematical Formulation

One method of determining relative importance of a node, whether it be a website or a character, is determine their *degree centrality*. The degree centrality is the number of edges incident to a given vertex. This method assumes that a node with more connections is more important than a node with fewer connections. This number will be the baseline value we will use to compare our PageRank results.

In order to begin investigating PageRank in terms of the Lord of the Rings characters, we need to first understand the process of the algorithm in terms of linear algebra. Let  $A$  be our  $n \times n$  edge matrix, where each element is nonnegative. Let  $A^T$  be the transpose of  $A$ . We, also, must create our network such that the matrix,  $A$  will be column stochastic (meaning each element in a column sums to 1). Since  $A$  is column stochastic, then  $A^T$  is row stochastic. If  $\vec{v}$  is a stochastic vector, then  $A\vec{v}$  is also stochastic. Since  $A^T$  is row stochastic, it will always have an eigenvalue of  $\lambda = 1$  with a corresponding (stochastic) eigenvector:

$$\begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} \quad (1)$$

Since  $A^T$  and  $A$  have the same eigenvalues,  $\lambda = 1$  is also an eigenvalue for  $A$ . To find the PageRank vector, the final step is to find the corresponding stochastic eigenvector. The largest value in the PageRank vector will correspond to the highest ranked node.

PageRank uses the idea that we use eigenvectors to adjust the weighted degree centrality based on the importance of connected vertices. We can calculate the importance  $x_i$  of vertex  $i$  as the weighted sum of the importance of its neighboring vertices:

$$x_i = \sum_{j \in V} a_{ji} x_j$$

for every  $i \in V$ . Solving this linear system for  $x_i$  gives us the eigenvector centrality, where we find an eigenvector for eigenvalue  $\lambda = 1$  of the matrix  $A$  with entries  $a_{ij}$ .

There are a couple problems with this algorithm. The first is when a character does not mention other characters, even if they are mentioned by others. The second results when a subnetwork of characters is not connected. To address these problems, we fix a positive constant  $p$  to a value called the *damping factor*, which in research is typically 0.15, and define the new PageRank matrix as:

$$M = (1 - p) \cdot A + p \cdot B$$

where  $B_{n \times n} = \frac{1}{n} \begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix}$ .

This method is very effective in determining the importance of a node, and led to the enormous success of Google as a search engine, and allows us to see interesting results in social networks.

## 5 Examples and Numerical Results

To analyze the importance of the characters in The Lord of the Rings, we found a copy of all three books on the internet. Using python, we prepared the text by making all words lowercase and removed all punctuation. We also replaced non-UTF-8 characters with their UTF-8 equivalent. We then wrote a program in C++ that first would read in a name file and create a node for each character in the file. Then the program would read through the text word by word, and if it found a character within fifteen words of another character, it would add an edge between those two characters. If the program encountered those two characters within fifteen

words of one another again, it would strengthen the edge between them. This creates a square, symmetric matrix the size of the number of characters.

Next, we divide each entry of column  $j$  by the sum of all elements in that column, creating a stochastic matrix  $S$ . We then run our iterative algorithm on this matrix, which gives us our final rank vector. The code for this final step is below:

```
//Matrix multiply S_matrix*rank = updated rank [iterations times]
for(int a=0; a < iterations; a++) {
    for(int i = 0; i < size; i++) {
        entry = 0;
        for(int j = 0; j < size; j++) {
            entry = entry + S_matrix[i][j]*ranks[j];
        }
        temp[i] = 0.85*entry+0.15;
    }
    ranks = temp;
}
```

Based purely on number of mentions, Frodo (unsurprisingly) takes the top spot at 1,982 mentions, with Sam coming in second at 1,361 mentions, and Gandalf in third with 1,170.

However, using the PageRank method of assessing importance, we do see some changes. Frodo is again first with a PageRank of 5.207, but now Gandalf is second with a PageRank of 3.803, with Sam in third with a PageRank of 3.605. The figure below shows the difference between each characters importance determined by number of connections to other characters versus their PageRank value.

## 6 Discussion and Conclusion

Ultimately, the results of our Page Rank analysis illustrated trends that coincided well with our knowledge of the characters and their interactions. This graph illustrates how character ranking changes when applying the Page Ranking algorithm, as opposed to a ranking based on sheer number of connections:

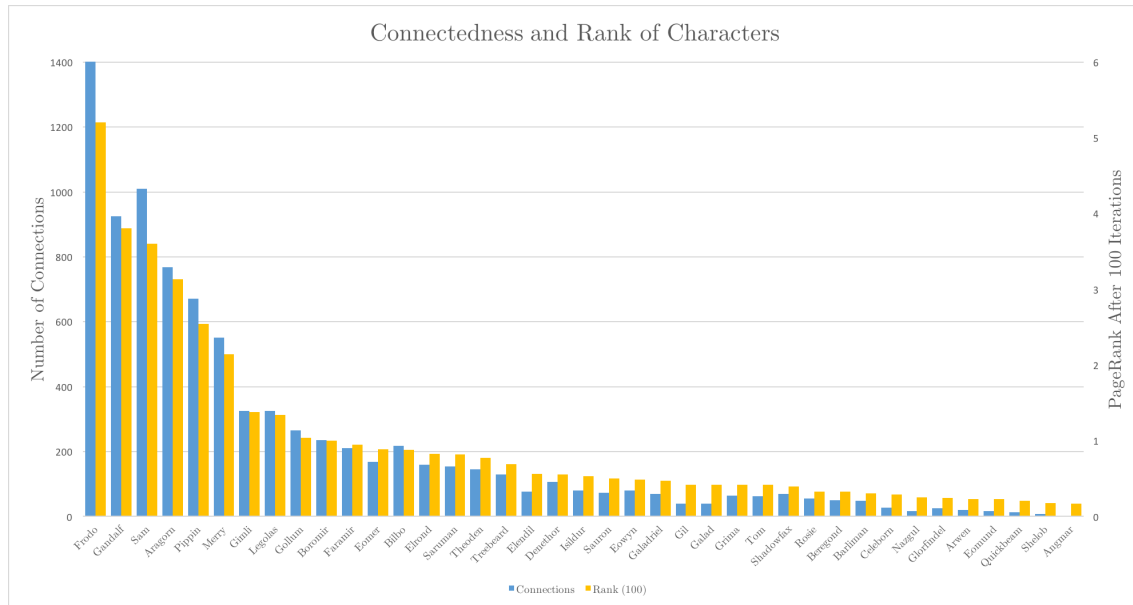
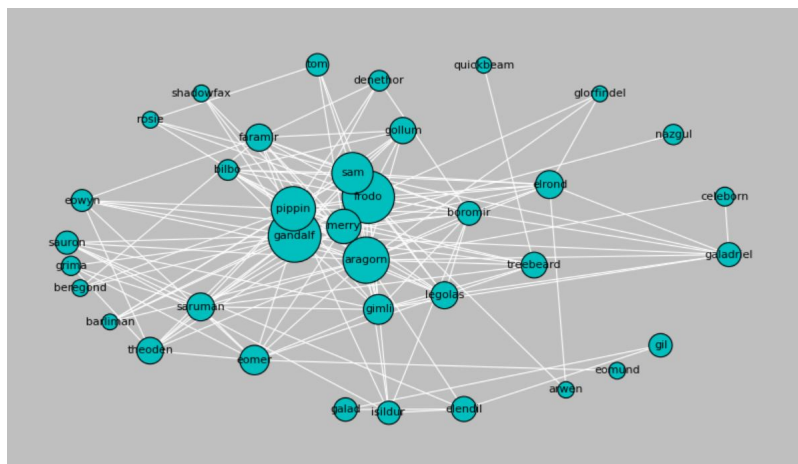


Figure 1: Connectivity versus PageRank

Upon inspection, we can see that the Page Rank Algorithm increases the importance of some characters more than others, which is particularly true for Gandalf, Gimli, and Legolas in the upper sections of the plot. Within the context of the novels these adjustments make sense; Gandalf acts as a bridge between multiple important characters and even from a subjective reading of the books we can assume that if someone is significant, he interacts with them. Gimli and Legolas are not as well connected, but rise in importance because of their strong friendship with Aragorn, the fourth most important character in the novels.



To further showcase these connections, we built a network graph of the relationships in the adjacency matrix. This portrays character relationships by basing node distances on relationship strength and node

size on calculated rank. This visualization technique was useful because it modeled social circles as well as character interactions; we recognized clusters of characters that tended to form parties throughout the novel, such as Aragorn with Legolas and Gimli, Frodo with Sam and Gollum, and Gil-galad with Elendil and Isildur. This last group consists of lesser-known, ancient characters of legend, so it was significant to see them grouping together so accurately on our graph since they tend to have fewer mentions overall.

Ultimately, we were happy with the results but were curious how they would have changed if we had made our graph directed, a characteristic that would be more in line with Page Rank's original algorithm since web page links cannot be two-way. We built a symmetric Matrix because determining the direction of a relationship is incredibly difficult - classifying character focus in any particular interaction would end up being more similar to a language processing project - and while we considered the approach of making smaller rankings of characters within chapters and deriving direction from connections to characters higher or lower in the rankings, we were unable to implement it within the scope of this project. We hope to improve our process even further by refining our interaction-detection methods and even applying this Method to books other than the Lord of the Rings.

## 7 References

Page L, Brin S, Motwani R, Winograd, T. (1999) *The PageRank Citation Ranking: Bringing Order to the Web*. Technical Report. Stanford InfoLab.

Tolkien J.R.R. 1954. *The Fellowship of the Ring*. Greylib; [accessed April 2, 2017]. [http://ae-lib.org.ua/texts-c/tolkien\\_the\\_lord\\_of\\_the\\_rings\\_1\\_en.htm](http://ae-lib.org.ua/texts-c/tolkien_the_lord_of_the_rings_1_en.htm).

Tolkien J.R.R. 1955. *The Return of the King*. Greylib; [accessed April 2, 2017]. [http://ae-lib.org.ua/texts-c/tolkien\\_the\\_lord\\_of\\_the\\_rings\\_3\\_en.htm](http://ae-lib.org.ua/texts-c/tolkien_the_lord_of_the_rings_3_en.htm).

Tolkien J.R.R. 1955. *The Two Towers*. Greylib; [accessed April 2, 2017]. [http://ae-lib.org.ua/texts-c/tolkien\\_the\\_lord\\_of\\_the\\_rings\\_2\\_en.htm](http://ae-lib.org.ua/texts-c/tolkien_the_lord_of_the_rings_2_en.htm).

Cornell University. 2014. (<https://blogs.cornell.edu/info2040/2014/11/03/more-than-just-a->

web-search-algorithm-googles-pagerank-in-non-internet-contexts/). Retrieved April 2017.

Princeton University. Ian Rogers. "The Google Pagerank Algorithm and How It Works". (<http://www.cs.princeton.edu/~chazelle/courses/BIB/pagerank.htm>). Retrieved April 2017.

Cornell University. "Lecture #3: PageRank Algorithm - The Mathematics of Google Search". (<http://www.math.cornell.edu/~mec/Winter2009/RalucaRemus/Lecture3/lecture3.html>). Retrieved April 2017.

Wikipedia. "PageRank". 2017. (<https://en.wikipedia.org/wiki/PageRank#Algorithm>). Retrieved April 2017.

Wikipedia. "Centrality". 2017. (<https://en.wikipedia.org/wiki/Centrality>). Retrieved April 2017.

Harvard University. "Lecture 33: Markov matrices". ([http://www.math.harvard.edu/~knill/teaching/math19b\\_2011/handouts/lecture33.pdf](http://www.math.harvard.edu/~knill/teaching/math19b_2011/handouts/lecture33.pdf)). April 2017.

North Carolina State University. "The Mathematics Behind Google's PageRank". ([http://www4.ncsu.edu/~ipsen/ps/slides\\_man.pdf](http://www4.ncsu.edu/~ipsen/ps/slides_man.pdf)). April 2017.

# 8 Appendix

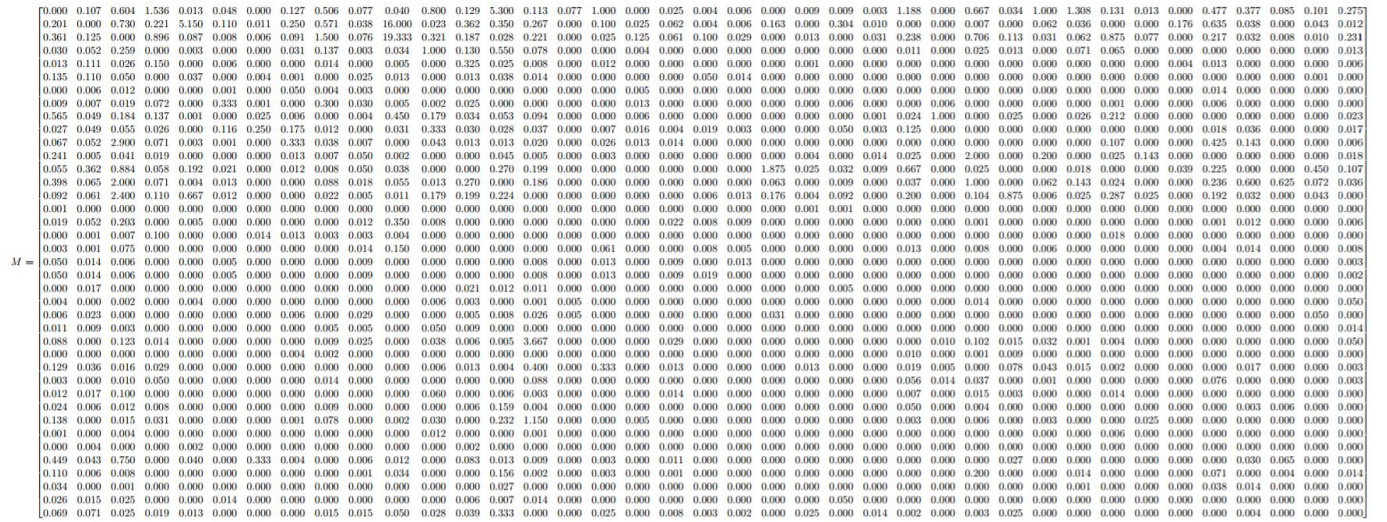


Figure 2: Adjacency Matrix of the dataset