

# Clash Of Clusters

Joseph Fahnestock

## Introduction

### Clustering and Project Task

A common workhorse found in the toolbox of data scientists is the notion of clustering. The core of clustering is constructing a rule that partitions observations into some number of groups. The simplest example of this tool is an algorithm like k-Nearest Neighbors which looks at data in a euclidean metric space (x-y coordinates) and tries to ensure that each cluster or group is composed of the most compact set of points possible. This problem is intuitive in the case of two dimensional point clouds (i.e. scatter plots) but the underlying principle works well in general if one can construct an appropriate metric space for a set of observations. This project aims to understand how using concepts from spectral graph theory allows a network of linked online forums (subreddits) to be partitioned into groups. The Excalibur that will cut through the abstract nature of a network diagram and give way to a meaningful clustering algorithm is the graph Laplacian.

### Laplacians and Spectral Theory

That begs the question, “What is a Laplacian and what more specifically is a graph Laplacian?” A Laplacian is a mathematical operator like the derivative. It is precisely defined as the divergence of the gradient. This is a clean, mathematical way of saying that the Laplacian takes a function of interest and tells you how much the speed of change is changing for a given point. How is that useful for clustering? The graph Laplacian can be thought of as a way to measure a cognitive phenomena attributable to the Gestalt Law of Closure where a person sees connected components in a point cloud or other image. Humans are able to make the educated guess that despite being further than other points, points which seem to form a continuous shape are likely connected in some way. The graph Laplacian allows this intuition to be abstracted into a tool which works on network diagrams and other scenarios where one cannot see the space which the data exists in.

The specific definition of the graph Laplacian is the difference of the degree matrix and the weight or similarity matrix of a network. The weight or similarity matrix is like an adjacency matrix where the entries indicate how strongly connected two nodes in our graph are. The degree matrix is just a diagonal matrix where each nonzero entry is the sum of the edge weights for the node it represents. This is a looser definition than you would find in a textbook, but it is fine for the purposes of this report

## Data

The data for this project comes from Stanford University’s SNAP datasets specifically their Social Network: Reddit Hyperlink Network Data. The data is comprised of two data sets which record hyperlinks from one subreddit to another in either the title or body of a post on Reddit. This project only considers the data set of body links for the sake of being manageable.

## Description

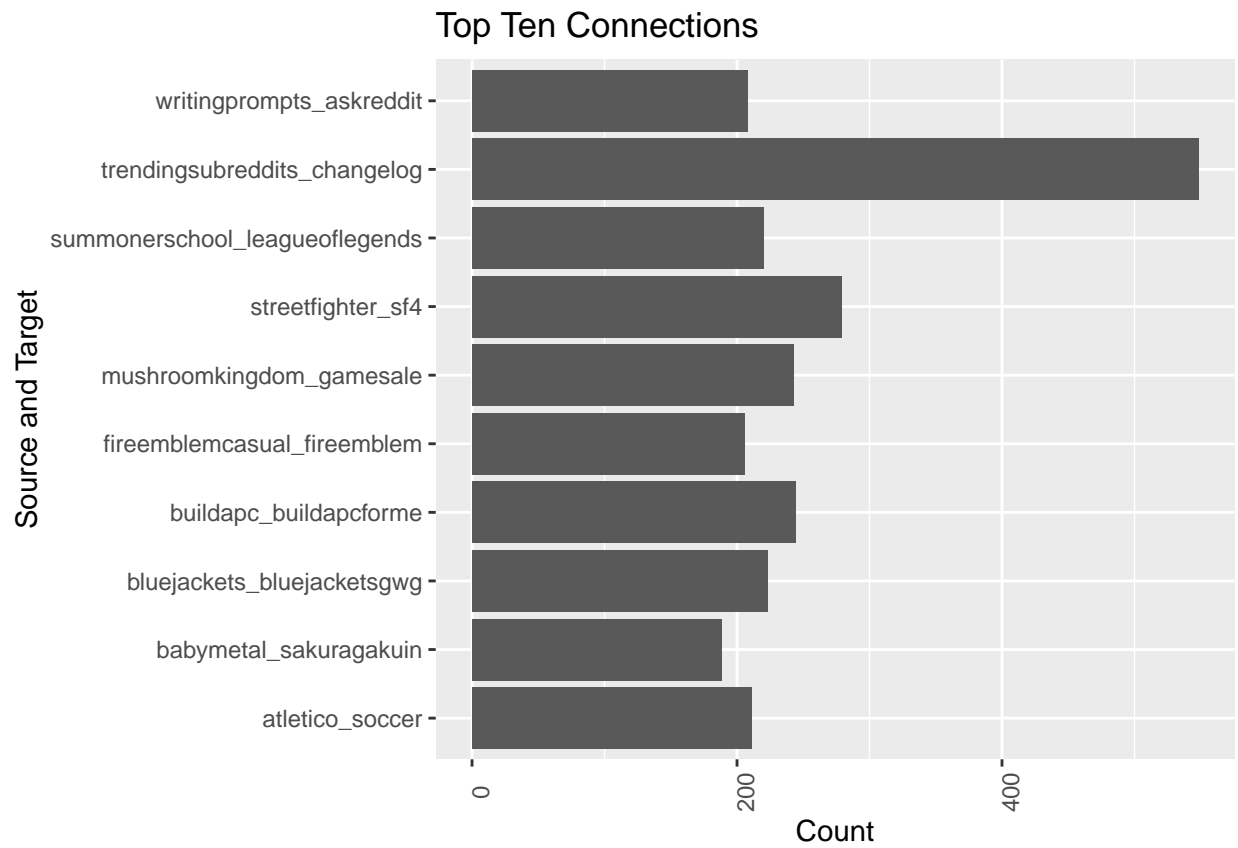
Variables Imported:

- **SOURCE\_SUBREDDIT:** A string containing the name of the subreddit where the post was made
- **TARGET\_SUBREDDIT:** A string containing the name of the subreddit linked to in the body of the post

Basic Information:

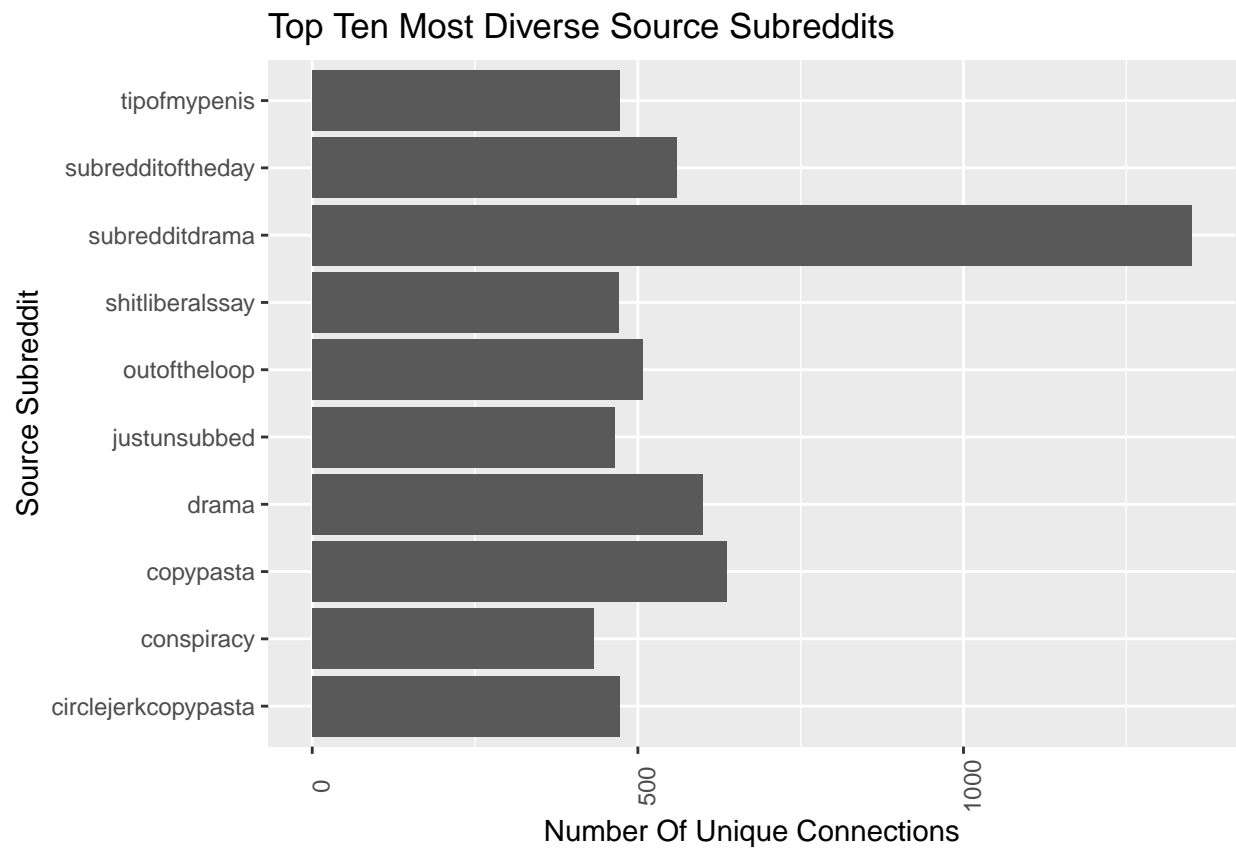
- **Number of Observations:** 286561
- **Number of Unique Source Subreddits:** 27863
- **Number of Unique Target Subreddits:** 20606

### Frequency of Connections

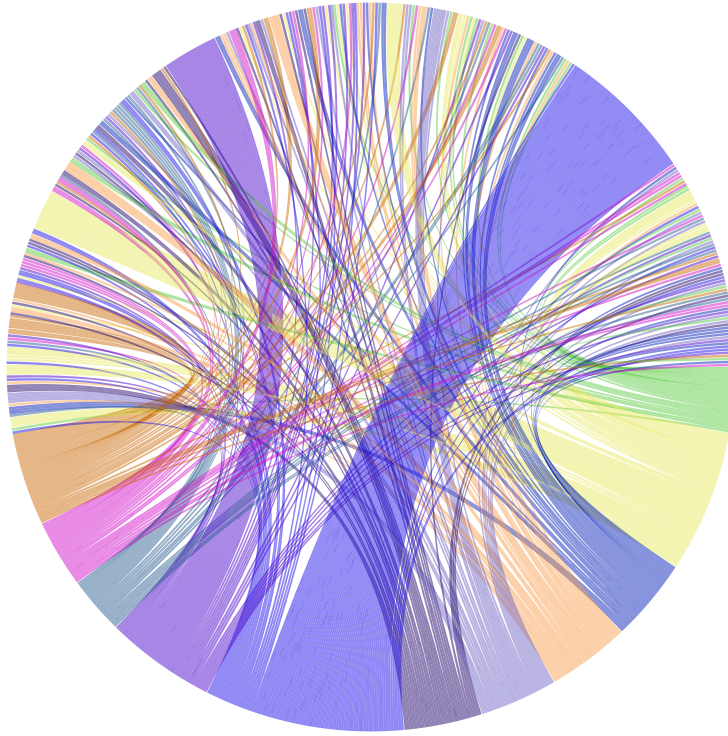




Diversity of Connections



## 1000th – 1010th Most Diverse Sources



The most diverse connections provides some great insights into how difficult this data could be to work with. The most diverse source has over 1250 unique targets. The overall trend amongst highly diverse sources is their relationship to drama or discussion of other subreddits. Again, the metaforums seem to dominate the activity and may prove to be outliers that need to be pruned. The unwieldy nature of the data is highlighted by the chord diagram which needed to use the 1000-1010th most diverse sources to render. The highly connected data would be difficult to analyze outside of proper graph analysis.

## Analysis

### Constructing the Similarity and Degree Matrices

The similarity matrix can be calculated in this case as a simple weighted adjacency matrix using **igraph**. The **Matrix** sparse matrix utilities makes the calculations feasible and taking the difference of the degree and adjacency matrix gives the graph Laplacian. Calculating the smallest nonzero eigenvalue and its eigenvector for the graph Laplacian yields a Fiedler vector which can be used to classify the data.

```
nodes = unique(unlist(body_links))
body_links_graph = graph_from_data_frame(d = body_links, vertices = nodes, directed = FALSE)
adj = get.adjacency(body_links_graph)
deg = Diagonal(35776, degree(body_links_graph))
graphLap = deg - adj
e = eigs(graphLap, k=20, which="SM")
```

```
## Warning in do.call(.Call, args = dot_call_args): only 2 eigenvalue(s) converged,
## less than k = 20
```

