

Money Moves the Pen

Link Prediction in Congress Bill Co-Sponsorship Networks Using Political Donor Network Information
Yi Zhong, Eddie Chen — CS224W Fall 2018 II Class Project

Introduction

Political collaboration is an important part of legislative life, and congress bill cosponsorships provide a rich source of information about the social network between legislators, and serving as a proxy to understand legislators’ ”connect-edness” and collaboration graph. Moreover, according to Mark Twain, ”we have the best government that money can buy” - money and politics have already been intertwined. In this project, we applied social network analysis tools on political donation networks and congress bill cosponsorship networks, and framed our research problem as a link prediction task on congress bill cospon-sorship networks using political campaign donation records for the US (Congress and Presidential Campaigns) with its network characteristics. We modeled and presented graph characteristics of the two political networks, and showed inves-tigation results of link prediction using various supervised learning techniques for this project. We then compared models’ performance to a naive baseline to come up with evaluations.

Method

Our project is made up of two parts: graph modeling, and link prediction. For graph modeling, we aim to construct a tripartite graph of political committees, legislators (we will ignore those failed to get elected to office), and the bills those legislators worked together on. A sample graph can be found in Figure 1. With the graph constructed, we provide a set of statistics and descriptions of the graph structure (including their one-mode projections, for both bills-legislators and committees-legislators subgraphs). After that, we construct a link prediction problem by dividing graph into different years of congress, and select the suitable years for model training and evaluation. Lastly, we report our learnings from the entire exercise.

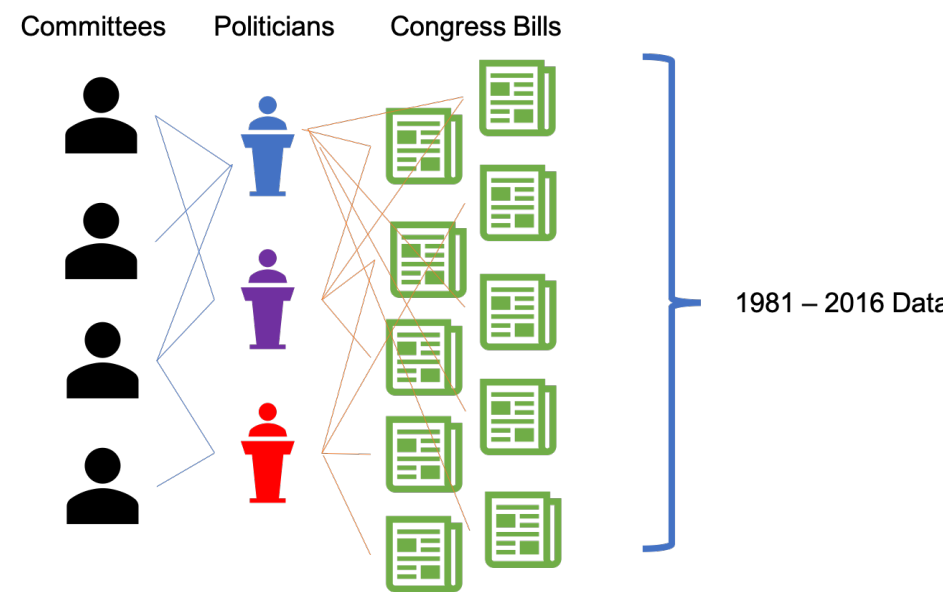


Figure 1: Illustration of the Congress Political Network

Link Prediction

We frame our link prediction problem as follows: predict the link between leg-islators, where a link exists if two legislators cosponsored a bill together, for a specific Congress term.

We have constructed three link prediction models:

- Naive baseline predictor: a baseline model based on graph density
- Legislator only predictor: a baseline model based on candidate attributes (party, state)
- Campaign only predictor: prediction models using features generated from campaign graph network attributes

Features and Algorithms

We constructed features from the campaign subgraph for Campaign only pre-dictor. Features include:

- Common Neighbors, Union of Neighbors, Jaccard Index
- Degree Difference in a pair of legislator nodes
- Contribution Amount (sum and absolute difference)
- Clustering Co-efficient (sum, absolute difference, mean)
- Degree Centrality difference
- Shortest Distance between two legislator nodes
- Spectral Clusters from Clauset-Newman-Moore greedy modularity maximiza-tion
- node2vec embeddings

For features used in the Legislator only predictor, we collected legislators’ affil-iated party and home state information, by congressional term.

For algorithms, we used supervised learning models like **logistic regression** and **decision trees**. For logistic regression, we used scikit-learn’s default im-plementation with -1,1 notation for labels and $L2$ regularization. A decision tree is a tree where each node represents a feature, each branch represents a decision/rule and each leaf represents a classification in our case. We used scikit-learn’s default implementation which uses Gini Index as the metric.

Evaluation Method

We used accuracy as our main success measure:

$$Accuracy = \frac{NumberOfCorrectPredictions}{TotalNumberOfPredictionsMade}$$

We define our **naive baseline predictor** as follows: given a pair of nodes v_1, v_2 , we will always predict there will be an edge between these two pairs, i.e. as a complete graph. That is,

$$Accuracy_{NaiveBaseline} = \frac{2||E||}{||V||(||V|| - 1)}$$

Using the 100th Congress (1987 - 1988) as the training set and the 101th Congress (1989-1990) as the test set. The baseline accuracy is calculated as $96,052/138,075 = 0.695$.

Results and Findings

The basic stats of the tripartitie graph are included below:

- Legislator count: 1,919 (1813 of which are found in campaign financial net-work)
- Bill count: 221,726
- Committee count: 14,326
- Edges between legislators and bills: 3,086,039
- Edges between committees and candidates: 911,965
- Overall tripartite graph node count: 237,971, and edge count: 3,998,004

Network Description

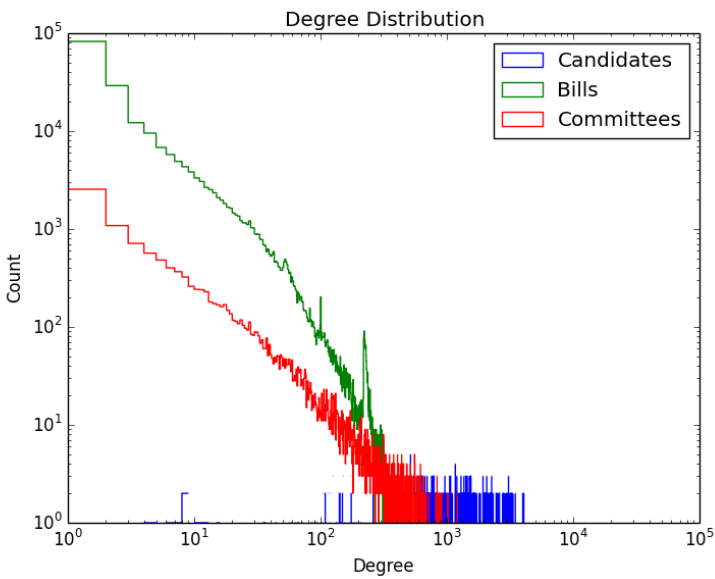


Figure 2: Overall Tripartite Graph Degree Distribution on log-log scale

Bill Co-authorship Graph resembles the academic collaboration graph with a power law pattern (long tail) - the most frequent degrees are the smallest degrees, and it has a very high clustering co-efficient. For the **folded campaign contribution graph, it represents a typical small world pat-tern with high clustering coefficient**.

Model Performance

Model	Train Accuracy	Test Accuracy
1 - Naive Baseline	0.695	0.695
2 - Candidate Party/State, Logistic Regression	0.695	0.698
3.1 - Campaign only, Logistic Regression	0.786	0.774
3.2 - Campaign only, Decision Tree	0.854	0.794
3.3 - Campaign only, Logistic Reg w/ node2Vec	0.728	0.728

Table 1: Model Performance for Limited Dataset

Model	Train Accuracy	Test Accuracy
1 - Naive Baseline	0.697	0.691
2 - Candidate Party/State, Logistic Regression	0.695	0.698
3.1 - Campaign only, Logistic Regression	0.748	0.714
3.2 - Campaign only, Decision Tree	0.795	0.740

Table 2: Model Performance for All Datasets Combined

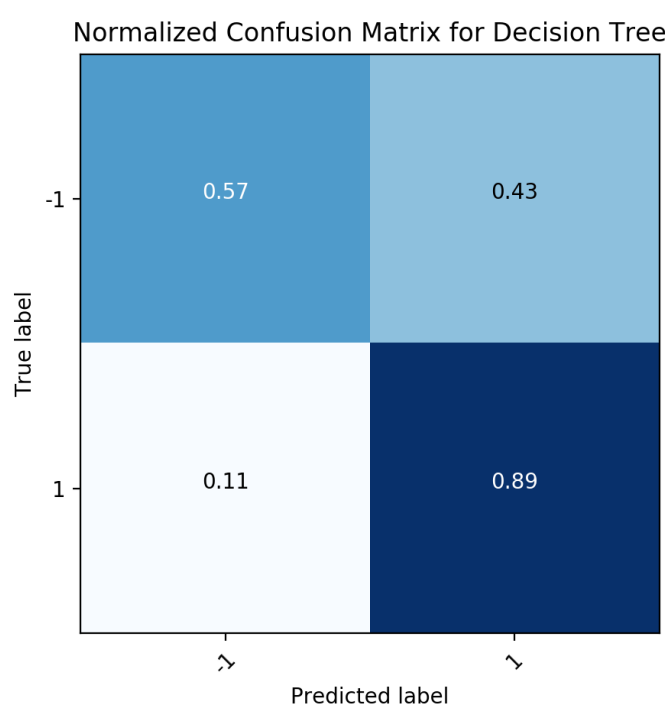


Figure 3: Confusion Matrix for Deci-sion Tree in Limited Dataset

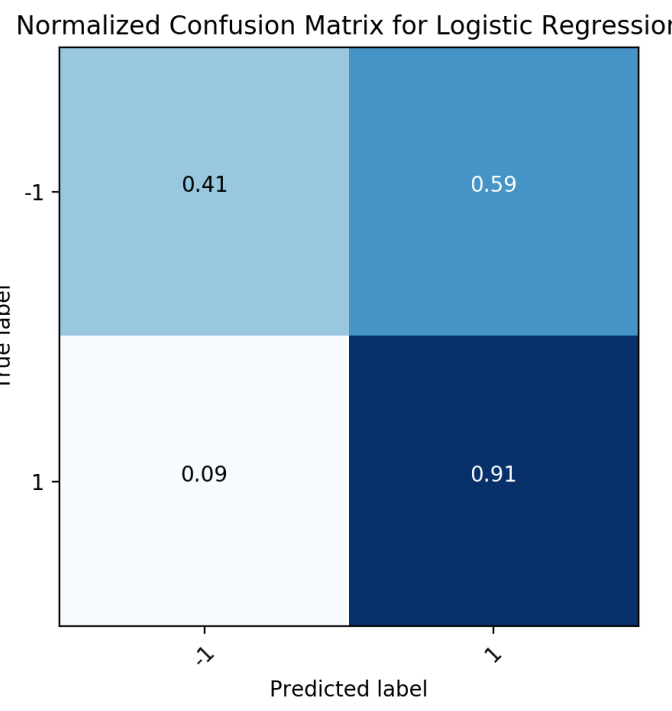


Figure 4: Confusion Matrix for Logis-tic Regression in Limited Dataset

Conclusion

US Congressional Politics is indeed a small world: legislators are connected to other legislators via common donors and co-authorship on bills. We have identified the academic collaboration network-like pattern for bill co-authorship data, and a ”small world” pattern among legislators, with consistently high clustering coefficients. Moreover, it does appear that ”money moves politics”: using features learned from campaign donation networks, we can confidently predict if two legislators will later collaborate on bills together - beating a naive baseline. In paricular, decision tree model performed very well to give us 79% accuracy.