

Aalto University  
School of Science  
Master's Programme in Computer, Communication and Information Sciences

Ananth Mahadevan

# Inferring Voting Networks in Online Elections

Master's Thesis  
Espoo, March 14, 2020

**DRAFT! — April 13, 2020 — DRAFT!**

Supervisor: Professor Aristides Gionis  
Advisor: Blank M.Sc. (Tech.)

Aalto University  
School of Science

Master's Programme in Computer, Communication and  
Information Sciences

ABSTRACT OF  
MASTER'S THESIS

<b>Author:</b>	Ananth Mahadevan		
<b>Title:</b>	Inferring Voting Networks in Online Elections		
<b>Date:</b>	March 14, 2020	<b>Pages:</b>	??
<b>Major:</b>	Computer Science	<b>Code:</b>	SCI3042
<b>Supervisor:</b>	Professor Aristides Gionis		
<b>Advisor:</b>	Blank M.Sc. (Tech.)		
abstract			
<b>Keywords:</b>	signed networks, balance, status, elections, Wikipedia, voting, graphs		
<b>Language:</b>	English		

Aalto-yliopisto

Perustieteiden korkeakoulu

Tieto-, tietoliikenne- ja informaatiotekniikan maisteriohjelma

DIPLOMITYÖN

TIIVISTELMÄ

<b>Tekijä:</b>	Ananth Mahadevan		
<b>Työn nimi:</b>	Äänestysverkkojen päätelmät online-valeissa		
<b>Päiväys:</b>	20. maaliskuuta 2020	<b>Sivumäärä:</b>	??
<b>Pääaine:</b>	Tietotekniikka	<b>Koodi:</b>	SCI3042
<b>Valvoja:</b>	Professori Aristides Gionis		
<b>Ohjaaja:</b>	Diplomi-insinööri Blank		
Finnish Abstract			
<b>Asiasanat:</b>	Finnish Keywords		
<b>Kieli:</b>	Englanti		

Aalto-universitetet

Högskolan för teknikvetenskaper

Magisterprogrammet i data-, kommunikations- och infor- SAMMANDRAG AV  
mationsteknik DIPLOMARBETET

<b>Utfört av:</b>	Ananth Mahadevan		
<b>Arbetets namn:</b>	Avsluta områdesnätverk i onlineval		
<b>Datum:</b>	Den 20 mars 2020	<b>Sidantal:</b>	??
<b>Huvudämne:</b>	Datateknik	<b>Kod:</b>	SCI3042
<b>Övervakare:</b>	Professor Aristides Gionis		
<b>Handledare:</b>	Diplomingenjör Blank		
Swedish abstract			
<b>Nyckelord:</b>	Swedish Keywords		
<b>Språk:</b>	Engelska		

# Acknowledgements

Espoo, March 14, 2020

Ananth Mahadevan

# Contents

# Chapter 1

## Introduction

In recent years, researchers have become increasingly interested in understanding the behaviour of voters in social networks. Knowledge of the factors that motivate voters is of great importance in selecting successful policies or candidates. This is a classic problem and has been studied in the fields of game theory and political science [? ? ?]. More recently, there is a focus on using information from the network of voters to model their behaviour. This provides an insight into the interactions and effect of influence on voters in a community. For example, voting for bills in the United States Congress [?] or electing administrators in Wikipedia [? ? ?].

Votes can be represented as a *signed* network with positive or negative links. Finding groups using clustering techniques [? ? ?] and predicting signed links [? ? ?] in these networks is well researched. These approaches provide an ability to understand the group dynamics at play and predict votes and in such a network. However, they do not consider the iterative and chronological nature of the voting that takes place in these networks. In cases where research does focus on voter models, they rely on external features to build machine learning models that are task-specific and static [? ? ?].

In this thesis, we propose a model that creates a local signed network consisting of the neighbours of the current voter and the preceding votes. It will then predict the vote, which when added to the network will comply the most with concepts of balance and status in signed networks. After all the votes are cast in a session, the model can be easily updated to improve quality and is, therefore, iterative and dynamic. The model is also flexible and can incorporate external features to build the local signed network of a voter. The results for Wikipedia administrator elections shows that our model outperforms machine learning based models and traditional signed link prediction solutions.

## 1.1 Thesis Outline

The rest of the thesis is organized in the following manner. We discuss the background relating to signed graphs, hierarchy in directed networks in Section ???. In Section ?? we describe the vote prediction problem and approaches to solving it. Section ?? provides a comprehensive view of Wikipedia and the election process for administrators. In Section ?? we explain the datasets used, construction of the model and evaluation criteria. We report our findings in Section ?? and discuss their implications. Finally, we conclude the thesis and present future work in Section ??.



# Graph Theory

## 2.1 Signed Graphs, Balance and Status

- Discuss basic terms related to graph theory
- Define terms such as Nodes, Edges, direction, edge weight, successor, predecessor and neighbors
- Signed graphs and restrictions
- Explain relevance in real world settings

- Explain balance theory origin and significance.
- Illustrate with triads and examples `iiiiiii HEAD`
- Define mathematical background to measure balance through the Eigen-decomposition of the graph Laplacian `=====`

- Define mathematical background to measure balance through the Eigen decomposition of the graph Laplacian [93f6b9e59b44e614844e797662563c7499148e27](#)

### 2.1.3 Status Theory

- Describe the nature of the directed setting
- Illustrate the differences to Balance theory
- Mention existing ways to measure violation to status in a network

## Chapter 3

# Vote Prediction

iiiiiii HEAD ===== llllllll 93f6b9e59b44e614844e797662563c7499148e27  
In this section we first provide the motivation of choosing independent vote prediction as our target and the differences from predicting the result of an election. Next we describe the available techniques and methods to predict individual votes or signed edges in a network and how it relates to the problem at hand. We then provide two novel methods of using user information long with past election results to predict votes.

### 3.1 Election versus Vote Prediction

- Discuss existing election result predictions schemes
- Discuss the limitations in understanding election dynamics through just predicting election results
- Describe the process as an information cascade, discuss the potential Game Theory settings
- Show the two parts of the problem from an information cascading perspective
  - Who is going to vote next
  - How they are going to vote
- Discuss the assumptions in usual Independent Cascade (IC) models
- Explain the difficulty of both aspects in the domain of an election
- Motivate the selection of the problem as an **Independent Vote Prediction**

## 3.2 Signed Edge Prediction

- Discuss the existing edge predictions work
- Directly using signed triads as features
- Using triads along with network features
- Using user information and interaction data for predicting votes and/or elections
- The main drawbacks in these methods when considering an election setting

## 3.3 Linear Combination of Graphs

- Describe the linear combination of graphs derived from user and election data
- Explain topic similarity, follows network, interaction networks and other features
- How it can also incorporate signed features as additional features in prediction

## 3.4 Local Signed Network

- Explain the concept of the local signed network for a particular user
- Motivate the definition with respect to elections and influence
- Describe how to use balance and status theory to predict the vote
- Clarify the differences to signed edge prediction efforts
- Mention Agony as a way to measure status compliance here?

## Chapter 4

# Wikipedia

In this section we provide an overview of how Wikipedia is structured, the hierarchy that exists withing editors. We then explain the election process of getting administrator rights in Wikipedia.

### 4.1 Structure and hierarchy in Wikipedia

### 4.2 Elections in Wikipedia

- Explain Editors and Administrators in Wikipedia
- Describe the Request for Administrator(RfA) process
- Discuss general trends and patters
- Mention research interest and possible current works?

## Chapter 5

# Experiments

In this section we first describe the datasets that will be used in building our vote prediction models. Then we discuss the various linear and graphical models that we consider and their implementations details. Lastly we define the metrics and other means of evaluating the models and the results.

### 5.1 Datasets

- Maybe a short description of existing SNAP datasets and their limitations
- The details of the *Wiki-RfA* data and the *User-Contribution* datasets

### 5.2 Graphs

- Discuss the process of extraction of the various graphs discussed in the previous sections
- **Agree Graphs and Follows Graph**, where we measure the degree to which one user agrees and follows another user in previous elections
- **Topic similarity** from the top 100 articles edited for each user and the pairwise Jaccard similarity
- **Talk and Interaction graphs**, measures communication between users on their respective user talk pages
- **Signed Graphs**, triad encoding and extracting the triad counts for each voter

## 5.3 Models

### 5.3.1 Linear Combination of Graphs

- Discuss the various linear models considered for Graph Combinations
  - Linear Regression
  - Support Vector Classifier
  - Extreme Gradient Boosting (XGBOOOOST)
- Discuss how each graph contributes features and the problem is a linear classification problem

#### 5.3.1.1 Iterative Mode

- Discuss the motivation behind an iterative model versus a static prediction model
- Describe how balance is derived from the Agree Graph in a local signed network
- Discuss how the Agree graph is updated in terms of Balance
- Describe how status is derived from the Follows graph in a local signed network
- Discuss how the Follows graph is updated after every election
- Describe how to make the predictions
  - Deterministic : just decide based on eigen value or agony as support or oppose
  - Probabilistic : provide a probability for predicting a support vote

## 5.4 Evaluation

- Discuss the issues with the imbalance in the datasets
- Illustrate the issues with pure measures of accuracy
- Define Precision, Recall and Macro F1 score
- Discuss ROC AUC and Precision Recall curves for probability based predictions

## Chapter 6

# Results and Discussion

In this section we will present the results of the models and discuss their implications.

### 6.1 Linear Combination of Graphs

- Present results for each linear classifier
- Discuss the different splits of the dataset to check for robustness and chronological consistency
- Show the feature importances and discuss their relevance
- Compare the raw accuracy versus the macro f1 scores
- Highlight the difficulty of predicting negative votes

### 6.2 Local Signed Network

- Present the Iterative Balance model results
- Discuss quality of predictions using evaluation metrics
- Mention the difference between deterministic and probabilistic prediction accuracies
- Explain the Iterative Status model results
- Discuss the issues with local model of status and the potential reasons for lower score and quality



## 6.3 Comparison

- Compare results from signed edge prediction and Iterative signed models
- Discuss Static Linear combination predictions versus Iterative signed predictions
- Discuss the assumptions used in the models and limitations

## Chapter 7

# Conclusions and Future Work

- Explain the quality of results with the election perspective
- Future work is to extend this to other election settings and investigate generality of this approach
- Possible future work in congressional voting data
- Can also tackle the other problem in information cascade theory of how to predict who is most likely to vote next
- This can lead to a complete model of election dynamics and could incorporate elements of game theory and network inference