

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! — March 28, 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

Author:	Ananth Mahadevan		
Title:	Inferring Voting Networks in Online Elections		
Date:	March 14, 2020	Pages:	18
Major:	Computer Science	Code:	SCI3042
Supervisor:	Professor Aristides Gionis		
Advisor:	Blank M.Sc. (Tech.)		
abstract			
Keywords:	signed networks, balance, status, elections, Wikipedia, voting, graphs		
Language:	English		

Aalto-yliopisto

Perustieteiden korkeakoulu

Tieto-, tietoliikenne- ja informaatiotekniikan maisteriohjelma

DIPLOMITYÖN
TIIVISTELMÄ

Tekijä:	Ananth Mahadevan		
Työn nimi:	Äänestysverkkojen päätelmät online-vaaleissa		
Päiväys:	20. maaliskuuta 2020	Sivumäärä:	18
Pääaine:	Tietotekniikka	Koodi:	SCI3042
Valvoja:	Professori Aristides Gionis		
Ohjaaja:	Diplomi-insinööri Blank		
Finnish Abstract			
Asiasanat:	Finnish Keywords		
Kieli:	Englanti		

Aalto-universitetet

Högskolan för teknikvetenskaper

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

Utfört av:	Ananth Mahadevan		
Arbetets namn:	Avsluta omröstningsnätverk i onlineval		
Datum:	Den 20 mars 2020	Sidantal:	18
Huvudämne:	Datateknik	Kod:	SCI3042
Övervakare:	Professor Aristides Gionis		
Handledare:	Diplomingenjör Blank		
Swedish abstract			
Nyckelord:	Swedish Keywords		
Språk:	Engelska		

Acknowledgements

Espoo, March 14, 2020

Ananth Mahadevan

Contents

1	Introduction	8
1.1	Thesis Contribution	8
1.2	Thesis Outline	8
2	Graph Theory	9
2.1	Signed Graphs, Balance and Status	9
2.1.1	Graphs and Signed Graphs	9
2.1.2	Balance Theory	9
2.1.3	Status Theory	10
2.2	Elections in Wikipedia	10
3	Vote Prediction	11
3.1	Election versus Vote Prediction	11
3.2	Signed Edge Prediction	12
3.3	Linear Combination of Graphs	12
3.4	Local Signed Network	12
4	Wikipedia	13
4.1	Structure and hierarchy in Wikipedia	13
4.2	Elections in Wikipedia	13
5	Experiments	14
5.1	Datasets	14
5.2	Graphs	14
5.3	Models	15
5.3.1	Linear Combination of Graphs	15
5.3.1.1	Iterative Mode	15
5.4	Evaluation	15
6	Results and Discussion	16
6.1	Linear Combination of Graphs	16

6.2	Local Signed Network	16
6.3	Comparison	17
7	Conclusions and Future Work	18

Chapter 1

Introduction

1.1 Thesis Contribution

1.2 Thesis Outline

Chapter 2

Graph Theory

Will provide background and concepts related to graph theory. Then we will focus on the task of predicting the vote of an individual user given their voting history and current state of the election. Next, we briefly explain the Wikipedia hierarchy and election process in Wikipedia. We then present the local signed network based model that accurately predicts a user's vote in a Wikipedia election.

2.1 Signed Graphs, Balance and Status

2.1.1 Graphs and Signed Graphs

- 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

2.1.2 Balance Theory

- Explain balance theory origin and significance.
- Illustrate with traids and examples
- Define mathematical background to measure balance through the Eigen-decomposition of the graph Laplacian

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

2.2 Elections in Wikipedia

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

Chapter 3

Vote Prediction

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 along 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

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